

BIAS-VARIANCE TRADE-OFFS ANALYSIS USING UNIFORM CR BOUND FOR IMAGES

M. Usman, A.O. Hero, and J.A. Fessler

University of Michigan, Ann Arbor MI 48109

ABSTRACT

We apply a uniform Cramer-Rao (CR) bound [1] to study the bias-variance trade-offs in parameter estimation. The uniform CR bound is used to specify achievable and unachievable regions in the bias-variance trade-off plane. The applications considered in this paper are: 1) two-dimensional single photon emission computed tomography (SPECT) system, and 2) one dimensional edge localization.

1. INTRODUCTION

The mean-square error (MSE) is an important measure of precision of a scalar component $\hat{\theta}_1$ of an estimator $\hat{\theta}$. It is well known that the MSE is a function of both the bias, denoted $\text{bias}_{\hat{\theta}}(\hat{\theta}_1)$ and the variance, denoted $\text{var}_{\hat{\theta}}(\hat{\theta}_1)$ of the scalar estimator:

$$\text{MSE}_{\hat{\theta}}(\hat{\theta}_1) = \text{var}_{\hat{\theta}}(\hat{\theta}_1) + \text{bias}_{\hat{\theta}}^2(\hat{\theta}_1).$$

Obviously increases in MSE can be due to increases in either the bias or variance of $\hat{\theta}_1$. Bias and variance are complementary in nature. While bias is due to 'mismatch' between the average value of the estimator and the true parameter, variance is due to statistical fluctuations in the estimator. There usually exists a tradeoff between bias and variance of the estimated parameter. For example in image reconstruction, implementation of the maximum likelihood algorithm with a smoothness penalty reduces the variance only at the expense of introducing bias. Different estimators can be effectively compared by plotting their performance on a bias-variance trade-off plane. The classical or the unbiased CR bound has been previously applied to compare different estimators [2, 3]. However, in most image processing applications the estimators are biased and their variance is not bounded by the unbiased CR bound. For biased estimators a biased CR bound is available [4] which is only applicable to estimators with fixed bias gradient $\nabla_{\hat{\theta}} \text{bias}_{\hat{\theta}}(\hat{\theta}_1)$, hence it is unable to give a meaningful comparison of different biased estimators that have acceptable bias but different bias gradients. We use uniform CR bound [1] on the variance of biased estimators which divides the bias-variance trade-off plane δ - σ into achievable and unachievable regions. Different estimators can be placed in the achievable region of the δ - σ plane and their performance can be effectively compared.

This work was supported in part by National Science Foundation under grant BCS-9024370, a Government of Pakistan Postgraduate Fellowship, NIH grant CA-60711, and DOE grant DE-FG02-87ER60561

2. UNBIASED CR BOUND

Consider the problem of estimation of an n -dimensional parameter $\underline{\theta} = [\theta_1, \dots, \theta_n]^T$ given an observation of a vector of random variables \underline{Y} with probability density function (pdf) $f_{\underline{Y}}(\underline{y}; \underline{\theta})$. The Cramer-Rao lower bound on the variance of unbiased parameter estimator $\hat{\theta}_1$ is given by the upper-left (1,1) element of the inverse of an $n \times n$, symmetric, positive definite Fisher information matrix (FIM) $F_Y = F_Y(\underline{\theta})$:

$$\text{var}_{\hat{\theta}}(\hat{\theta}_1) \geq \underline{e}_1^T F_Y^{-1} \underline{e}_1, \quad (1)$$

where,

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

$\nabla_{\underline{\theta}}$ denotes the (row) gradient vector $[\frac{\partial}{\partial \theta_1}, \dots, \frac{\partial}{\partial \theta_n}]$, and $\underline{e}_1 = [1, 0, \dots, 0]^T$ is an n -element unit vector.

While the unbiased CR bound (1) is known to be asymptotically achievable for large number of independent identically distributed measurements, in practice, most estimation algorithms are biased and the unbiased CR bound is inapplicable.

3. UNIFORM CR BOUND

For a biased estimator $\hat{\theta}_1$ the following form of the biased CR bound is well known [4]:

$$\text{var}_{\hat{\theta}}(\hat{\theta}_1) \geq [\nabla_{\underline{\theta}} m_1] F_Y^{-1} [\nabla_{\underline{\theta}} m_1]^T, \quad (2)$$

where $\nabla_{\underline{\theta}} m_1 = \nabla_{\underline{\theta}} m_1(\underline{\theta}) = \nabla_{\underline{\theta}} b_1 + \underline{e}_1$ is an n element row vector of the gradient of the mean $E_{\hat{\theta}}(\hat{\theta}_1) = m_1(\underline{\theta})$. The application of the biased CR bound (2) is very restricted due to the fact that it is only applicable to estimators with a given bias gradient $\nabla_{\underline{\theta}} b_1$. In [1] Hero gives a 'uniform' CR bound on the variance of a single parameter θ_1 for non-singular F_Y . This bound is applicable to all biased estimators whose bias gradient length $\|\nabla_{\underline{\theta}} b_1\|$ satisfies:

$$\|\nabla_{\underline{\theta}} b_1\|^2 \leq \delta^2 < 1. \quad (3)$$

The following theorem is proven in [1].

Theorem 1 *Let $\hat{\theta}_1$ be an estimator with bias $b_1(\underline{\theta})$ whose n -element bias gradient vector $\nabla_{\underline{\theta}} b_1$ satisfies (3). Assume that the FIM F_Y is non-singular. Then the variance of $\hat{\theta}_1$ is given by:*

$$\text{var}_{\hat{\theta}}(\hat{\theta}_1) \geq B(\underline{\theta}, \delta), \quad (4)$$

where $B(\underline{\theta}, \delta)$ is equal to:

$$\begin{aligned} B(\underline{\theta}, \delta) &= [\underline{\epsilon}_1 + \underline{d}_{\min}]^T F_Y^{-1} [\underline{\epsilon}_1 + \underline{d}_{\min}], \\ &= \lambda^2 \underline{\epsilon}_1^T [I + \lambda F_Y]^{-1} F_Y [I + \lambda F_Y]^{-1} \underline{\epsilon}_1 \end{aligned} \quad (5)$$

where $\underline{\epsilon}_1 = [1, 0, \dots, 0]^T$ is an n -element unit vector and:

$$\underline{d}_{\min} = -[I + \lambda F_Y]^{-1} \underline{\epsilon}_1, \quad (7)$$

and λ is given by the unique non-negative solution of the following equation involving the monotone decreasing, strictly convex function $g(\lambda) \in [0, 1]$:

$$g(\lambda) = \underline{d}_{\min}^T \underline{d}_{\min} = \delta^2 \quad \lambda \geq 0. \quad (8)$$

A more general version of Theorem 1, which will not be required here, is given in [5] and applies to singular F_Y . Note that since $\lambda \geq 0$ and $F_Y \geq 0$, the use of the expression (6) does not suffer from any ill-conditioning of the FIM F_Y . In Theorem 1, \underline{d}_{\min} defined in (7) is an optimal bias gradient in the sense that it minimizes the biased CR bound (2) over all vectors $\nabla_{\underline{\theta}} b_1$.

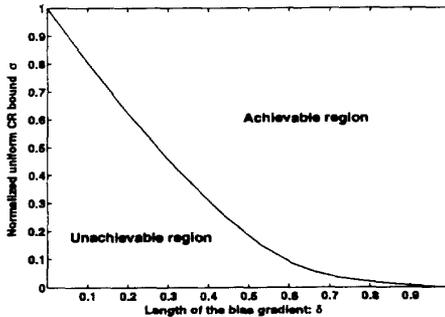


Figure 1: The Normalized Uniform CR bound.

Figure 1 shows a typical bias-variance trade-off curve in the δ - σ plane. The region above and including the curve is the so called 'achievable' region where all the realizable estimators exist. Note that if an estimator lies on the curve then lower variance can only be bought at the price of increased bias and vice versa. At $\delta = 1$ the variance goes to zero. This corresponds to the trivial case $\hat{\theta}_1 = \text{Constant}$ for which $\nabla_{\underline{\theta}} b_1 = \underline{\epsilon}_1$.

3.1. Estimation of the Bias Gradient

To compare a particular estimator to the uniform bound of Theorem 1 we require the length of the estimator bias gradient so that the estimator can be placed somewhere within the achievable region of Figure 1. In most cases the bias and the bias-gradient are analytically intractable. The method of moments is the standard method for experimentally determining bias and covariance of $\hat{\underline{\theta}}$ which is based on forming the sample mean and sample covariance statistics for a sequence of L repeated experiments $\{\underline{Y}_i\}_{i=1}^L$ each generated from the density $f_{\underline{Y}}(\underline{y}_i; \underline{\theta})$. The method of moments for estimating the bias-gradient would require n additional

sequences of L repeated experiments, each generated for a particular perturbation of a different component of the parameter vector $\underline{\theta}$. Such a direct approach is impractical. In [5] a method for experimentally determining the bias-gradient of an estimator $\hat{\underline{\theta}}$ is presented that requires a single simulation of the same type as that commonly used to determine bias and covariance of $\hat{\underline{\theta}}$. The unbiased estimate of the bias gradient for the estimate of $\hat{\theta}_1$ is given by [5]: $\nabla_{\underline{\theta}} \hat{b}_1 =$

$$\frac{1}{L-1} \sum_{i=1}^L \left(\hat{\theta}_1(\underline{Y}_i) - \frac{1}{L} \sum_{j=1}^L \hat{\theta}_1(\underline{Y}_j) \right) \nabla_{\underline{\theta}} \ln f_{\underline{Y}}(\underline{Y}_i; \underline{\theta}) - \underline{\epsilon}_1^T. \quad (9)$$

A few comments about the bias gradient are in order. The bias gradient $\nabla_{\underline{\theta}} b_1$ is a measure of the influence of each component parameter $\theta_1, \dots, \theta_n$ on the mean $m_1(\underline{\theta})$ of the estimator $\hat{\theta}_1$. Ideally, to be close to unbiased one would like $m_1(\underline{\theta})$ to be insensitive to the variations in the other parameters $\theta_2, \dots, \theta_n$. Alternatively, since $b_1(\underline{\theta}) = m_1(\underline{\theta}) - \theta_1$, it is desirable that the components $\frac{\partial}{\partial \theta_k} b_1(\underline{\theta})$ be of small magnitude, $k = 2, \dots, n$. The bias gradient therefore provides important information about the parameter coupling to the estimator mean. The bias gradient is in general only indirectly related to the estimator bias, with the exception that $\nabla_{\underline{\theta}} b_1 = 0$ implies $b_1(\underline{\theta}) = \text{constant}$. An estimator that has a constant bias independent of $\underline{\theta}$ is removable, and therefore $\nabla_{\underline{\theta}} b_1 = 0$ implies that the estimation can be performed without bias. Conversely, a non-zero bias gradient implies non-removable estimator bias that is dependent on the estimator parameters. On the other hand, one can have a large bias gradient even though the bias is very small. Therefore the bias and the bias gradient together give a more complete picture of estimator behavior.

3.2. Bias-Variance Trade-Off Plane

When accurate estimates \hat{b}_1 , $\nabla_{\underline{\theta}} \hat{b}_1$ and $\hat{\sigma}^2$ of the estimator bias, bias gradient, and variance are available for a given estimator $\hat{\theta}_1$ of θ_1 , the uniform CR bound lying in the δ - σ plane can be easily mapped into the b - σ plane of variance and biases. This is accomplished by using the ordered triplet $(\hat{b}_1, \nabla_{\underline{\theta}} \hat{b}_1, \hat{\sigma}^2)$ as a mapping between the δ - σ and the b - σ planes. The uniform CR bound on the variance as a function of bias is simply the ordered pair: $\left(\hat{b}_1, \left[\underline{\epsilon}_1 + \nabla_{\underline{\theta}}^T \hat{b}_1 \right]^T F_Y^{-1} \left[\underline{\epsilon}_1 + \nabla_{\underline{\theta}}^T \hat{b}_1 \right] \right)$, denoted $B(\underline{\theta}; b)$ in the sequel.

4. APPLICATIONS

We will apply the uniform CR bound to study the bias-variance trade-offs for: 1) a particular class of roughness penalized maximum-likelihood (PML) in SPECT image reconstruction, and 2) one-dimensional edge localization.

4.1. SPECT Image Reconstruction

4.1.1. System Description

The system used in this paper is shown in Figure 2 and is called the SPRINT II system [6]. The system was designed specifically for brain imaging and consists of a ring of detectors and a ring of collimators. The function of the collimator is to reduce the uncertainty associated with the emission location of a γ -ray to a line or a strip in the field of view (Figure 2). During imaging time, the collimator ring is rotated through small steps about the source. A γ -ray photon passing through one of the collimator slits at one of the rotation angles is counted as an event acquired in one 'detector bin'. For reconstruction the source domain is divided into n small regions, called pixels. The detection process is governed by Poisson statistics: $\underline{Y} = [Y_1, \dots, Y_d]^T$.

$$f_{\underline{Y}}(\underline{y}; \underline{\theta}) = \prod_{j=1}^d \frac{\mu_j^{Y_j}}{Y_j!} e^{-\mu_j}. \quad (10)$$

In (10) θ_i is the average γ -ray intensity of the i -th pixel; $i = 1, \dots, p$, Y_j is number of γ -rays detected at the j -th detector, and μ_j is the average γ -ray intensity of the j -th detector; $j = 1, \dots, d$; $\underline{\mu} = A \underline{\theta}$, where A is the $d \times p$ system matrix that depends on the tomographic geometry.

The objective is to reconstruct the object intensity of each pixel $\underline{\theta} = [\theta_1, \dots, \theta_n]^T$ given the set of observations \underline{Y} . It can be easily shown that the FIM is of the form:

$$F_Y(\underline{\theta}) = A^T [\text{diag}(\underline{\mu})]^{-1} A \quad (11)$$

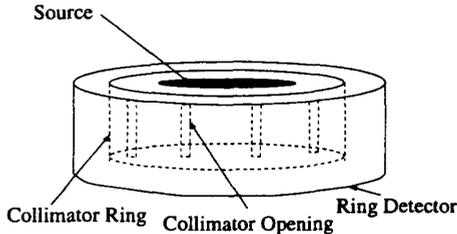


Figure 2: The SPRINT II system. Not drawn to scale.

The system parameters are given in Appendix A and unless otherwise specified are those used in the simulations.

In the following simulations the effect of attenuation was neglected. The total number of detected γ -ray counts were 10^9 . Noise due to scatter were 5% of the total counts. Since the algorithm considered in this section is non-linear, an analytic expression for the bias gradient is intractable, and therefore the bias gradient was estimated using (9). We used $L = 400$ realizations of the projection data \underline{Y} . The object is a disk of uniform intensity 1 with a high intensity region of 4 pixels in the center of uniform intensity 2, called the hot spot. The pixel of interest was the pixel at the upper edge of the hot spot, marked '1'. The diameter of the disk is 32 pixels. In the following simulation, the algorithm was initialized by a uniform disk of intensity 1 and diameter 32 pixels.

4.1.2. Penalized Maximum Likelihood

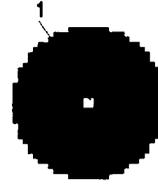


Figure 3: The object used in the simulations. The object dimensions are 32×32 . The black pixels are of intensity 1 while the white pixels are of intensity 2.

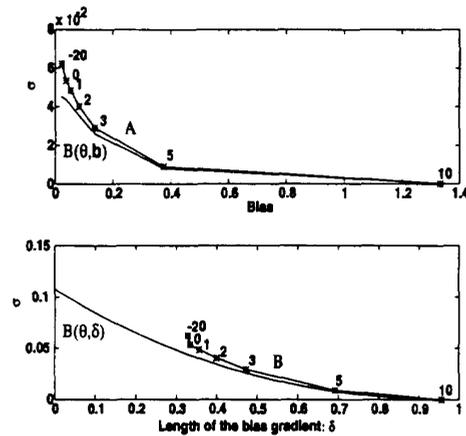


Figure 4: Performance of PML: MAP-SAGE as a function of α

The penalized maximum-likelihood (PML) is penalized for roughness and has the same functional form as a MAP estimator of the image intensities $\underline{\theta}$. The general form of the PML is given by:

$$\hat{\underline{\theta}}(\underline{Y}) = \underset{\underline{\theta} \in \Theta}{\text{argmax}} \left\{ \ln f_{\underline{Y}}(\underline{y}; \underline{\theta}) - \alpha P(\underline{\theta}) \right\},$$

where $P(\underline{\theta})$ is a roughness penalty and α is the smoothing parameter. We use a penalty function described in [7] which is imposed on the 8 neighboring pixels for each pixel of interest. Setting $\alpha = 0$ corresponds to no image smoothing while a large value of α corresponds to a significant amount of smoothing. We have implemented the recursive SAGE algorithm to maximize the PML objective function. SAGE, which stands for space alternating generalized EM, involves an intelligent choice of a 'complete data space' such that the E and M steps are analytically tractable. A detailed description of the PML-SAGE algorithm is given in [7].

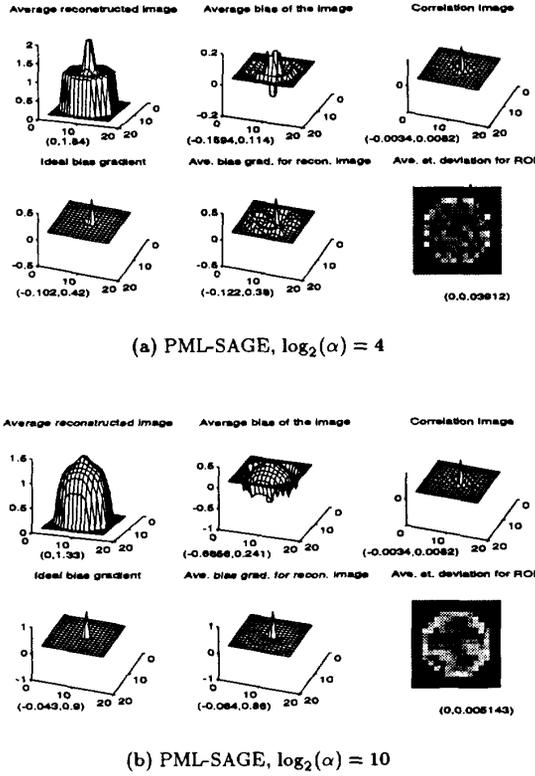


Figure 5: PML-SAGE: different image quantities of interest. An ordered pair with each curve indicates the (minimum, maximum) value associated with that image. The images in (a) and (b), from top left to bottom right, are: Average reconstructed image, average bias of the reconstructed image, correlation image, ideal bias gradient \underline{d}_{min} , average bias gradient for the reconstructed image, and average standard deviation.

It is easy to show that for the Poisson model

$$\nabla_{\underline{\theta}} \ln f_{\underline{Y}}(\underline{y}; \underline{\theta}) = A^T [-\underline{1} + \underline{Y} \oslash \underline{\mu}],$$

where \oslash is a vector operation denoting element-by-element division, and $\underline{1} = [1, 1, \dots, 1]^T$.

For the first set of simulations the smoothing parameter α was varied (Figure 4). Points on the curves in Figure 4 are labeled by the exponent of α . The bias, bias gradient and variance were estimated and the uniform bound was plotted over δ - σ and b - σ domains. The MAP-SAGE algorithms were terminated after 100 iterations for each of the $L = 400$ trials. The 95% ellipsoidal confidence regions are not shown in the figure since they are smaller than the size of the plotting symbol '*'. Note that the bound, denoted by $B(\underline{\theta}; \delta)$ in Figure 4, is achieved for large biases, i.e. large α . For α small, the curve 'B' tends to deviate more from the

lower bound and saturate, i.e. lower α does not decrease the bias gradient. On the other hand the bias decreases to an asymptote near zero. At points close to the unbiased point, i.e. the leftmost corner of the horizontal axis, in curve 'A', maximal reduction in bias is achieved at the price of significant increase in the variance.

Figure 5 shows several image quantities of interest for $\alpha = 2^4$, and $\alpha = 2^{10}$, respectively. For clarity in the figures, we down-sampled all the images by a factor of 2. For each image in Figures 5 the ordered pair at bottom indicates the minimum and maximum values for that image. In Figure 5 (a), the reconstructed image is very close to the true image except around the edges. The correlation image, i.e. the column of F_Y^{-1} corresponding to the pixel of interest, θ_{ROI} , shows a strong correlation with the neighboring pixels. This implies that to estimate θ_{ROI} we must also estimate the strongly correlated neighboring pixels accurately, while the influence of the far pixels can be ignored. Ideally, one would like the correlation between the pixels to be zero so that the estimate of a certain pixel, θ_{ROI} , is independent of the estimates of all other pixels. The plot for the theoretically optimal bias gradient \underline{d}_{min} shows a similar strong influence from the neighboring pixels.

The average bias gradient $\nabla_{\underline{\theta}} b_1$ for the reconstructed image is different from the theoretically optimal bias gradient \underline{d}_{min} . Thus the penalized SAGE image reconstruction algorithm does not take best advantage of its bias allocation since it is only by using the optimal bias gradient \underline{d}_{min} that the minimum bias length is achieved.

Figure 5 (b) shows the same set of images as in Figure 5 (a) but for $\alpha = 2^{10}$. Due to very high regularization, the hot spot is almost entirely smoothed out. Also, neither \underline{d}_{min} nor the average bias gradient $\nabla_{\underline{\theta}} b_1$ for the reconstructed image show significant coupling between the pixel of interest and the neighboring pixels. This is to be expected since in the overly smoothed case the bias is principally determined by the smoothness penalty as opposed to the projection data.

4.2. One-Dimensional Edge Localization

In many imaging applications it is important to determine the location l of an edge along an oriented line segment. In [3] the unbiased CR bound is derived on the localization accuracy of an edge estimator. As in [3] we define an edge by the following 3 parameters (Figure 6): 1) Intensity I , 2) location l , and 3) width σ_s . The edge is modelled as the following function of position x along the oriented line segment:

$$R(x) = I\Phi\left(\frac{x-l}{\sigma_s}\right) + q(x),$$

where $q(x)$ is additive white Gaussian noise of variance n_o^2 , and Φ is the cumulative distribution function of an $\mathcal{N}(0, 1)$ Gaussian random variable. We assume that the width σ_s of the edge is known.

The FIM $F_R(\underline{\theta})$ for the parameter vector $\underline{\theta} = [I, l]^T$ based on the noisy edge observation R is given by [3]:

$$F_{11} = \frac{I^2}{n_o^2 \sqrt{\pi} \sigma_s} \left[\Phi\left(\sqrt{2} \frac{T_x - l}{\sigma_s}\right) - \Phi\left(-\sqrt{2} \frac{T_x + l}{\sigma_s}\right) \right]$$

$$F_{12} = \frac{-I^2}{n_o^2} \left[\Phi^2\left(\frac{T_x - l}{\sigma_s}\right) - \Phi^2\left(-\frac{T_x + l}{\sigma_s}\right) \right] = F_{21}$$

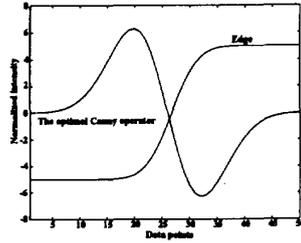


Figure 6: A typical edge profile along with the optimal Canny operator in our example. Edge parameters: Intensity $l=10$, width $\sigma_s = 6$, and location $l=26$.

$$F_{22} = \frac{2}{n_o^2} \int_{-T_x}^{T_x} \Phi^2 \left(\frac{x-l}{\sigma_s} \right) dx,$$

where T_x is the extent of the observation window.

An estimator of l of the edge location was constructed using the Canny operator.

$$h_c(x) = \phi'_{\sigma_c}(x) = \left(\frac{-x}{\sqrt{2\pi}\sigma_c^2} e^{-\frac{x^2}{2\sigma_c^2}} \right) W(x),$$

where ϕ is a Gaussian function with scale σ_c , and $W(x)$ is a window function:

$$W(x) = \begin{cases} 1 & x \in [-\frac{T_x}{2}, \frac{T_x}{2}] \\ 0 & \text{otherwise.} \end{cases}$$

The precess $R(x)$ is filtered by the Canny operator to produce an output $f_c(x)$:

$$f_c(x) = R(x) * h_c(x),$$

where $*$ denotes discrete convolution. The minimum value of $f_c(x)$ determines the location of the edge. It is shown in [3] that the optimal choice of the Canny width σ_c , determined by minimizing the unbiased CR bound, is $\sqrt{5}\sigma_s$ for an unbiased edge localization algorithm.

The length of the data $R(x)$ containing the edge was 1000 points. The edge parameters used were: $l = 15$, $\sigma_s = 6$, and $l = 501$. We used a window T_x of 50 data points, $n_o^2 = 8$. We varied σ_c from 3 corresponding to a difference operator, to 31 corresponding to a ramp filter. For each value of σ_c investigated we generated 100 independent realizations of noisy edge profile $R(x)$. The bias gradient was estimated using (9). The results are shown in Figure 7. The 95% confidence intervals are smaller than the size of the plotting symbol $*$.

The curve 'B' in Figure 7 shows a point of minimum variance at $\sigma_c = 16$, which also corresponds to minimum bias (curve 'A') on the b - σ plane, and hence a point of minimum MSE. Note that the minimum variance is achieved close to the optimal $\sigma_c = \sqrt{5}\sigma_s = 13.5$ determined by minimizing the unbiased CR bound. An interesting point to note is that although the bias and the variance vary non-monotonically with increasing σ_c , the bias gradient length increases monotonically. For σ_c between 4.5 and 16 the estimator standard deviation tracks the uniform CR bound $B(\underline{\theta}, \delta)$, however with an offset of approximately 0.2.

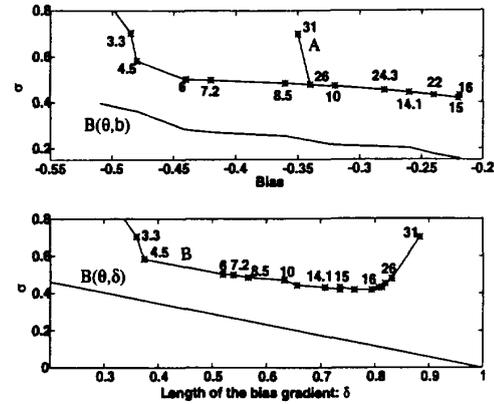


Figure 7: The uniform CR bound and the sample variance for varying σ_c . The numbers associated with the curves 'A' and 'B' indicate σ_c .

A. SYSTEM SPECIFICATIONS

Radius of the detector ring	25 cms
Number of detectors	512
Radius of the collimator ring	17 cms
Number of collimator slits	10 (uniformly spaced)
Slit Width	2.4 mm

5. REFERENCES

- [1] A. O. Hero
A Cramer-Rao Type Lower Bound for Essentially Unbiased parameter Estimation, Technical Report 890, Lincoln Laboratory, MIT 1992.
- [2] M. Usman, A. O. Hero and W. L. Rogers
Performance Gain Analysis for Adding Vertex View to a Standard SPECT, MWSCS, August 1993, Detroit, MI.
- [3] R. Kakarala and A.O. Hero
On Achievable Accuracy in Edge Localization, IEEE Transactions on Pattern Analysis and Machine Intelligence, 14:7, pp. 777-781, July, 1992.
- [4] H.L. van Trees
Detection, Estimation and Modulation Theory, (Part I), John Wiley and Sons. 1968.
- [5] M. Usman, A.O. Hero, and J.A. Fessler
Bias-Variance Trade-offs For Parametric Estimation Problems Using Uniform CR Bound, To be submitted to IEEE Transactions on Signal Processing.
- [6] W.L. Rogers, N.H. Clinthorne, L. Shao, P. Chiao, J. Stamos, and K.F. Koral
SPRINT II, A Second Generation Single Photon Ring Tomograph, IEEE Transactions on Medical Imaging, 7:4, pp. 291-297, 1988.
- [7] J.A. Fessler and A.O. Hero
Space-Alternating Generalized Expectation-Maximization Algorithm, To appear in IEEE Transactions on Signal Processing, Oct., 1994.