

Correspondence

Wavelet Thresholding for Multiple Noisy Image Copies

S. Grace Chang, Bin Yu, and Martin Vetterli

Abstract—This correspondence addresses the recovery of an image from its multiple noisy copies. The standard method is to compute the weighted average of these copies. Since the wavelet thresholding technique has been shown to effectively denoise a single noisy copy, we consider in this paper combining the two operations of averaging and thresholding. Because thresholding is a nonlinear technique, averaging then thresholding or thresholding then averaging produce different estimators. By modeling the signal wavelet coefficients as Laplacian distributed and the noise as Gaussian, our investigation finds the optimal ordering to depend on the number of available copies and on the signal-to-noise ratio. We then propose thresholds that are nearly optimal under the assumed model for each ordering. With the optimal and near-optimal thresholds, the two methods yield similar performance, and both show considerable improvement over merely averaging.

Index Terms—Filter noise, image denoising, image restoration, wavelet thresholding.

I. INTRODUCTION

Denoising via wavelet thresholding proposed by Donoho and Johnstone [3] is a simple nonlinear yet effective technique which outperforms linear techniques in theory and practice (cf. [5, ch. 10]). Since this seminal work, there have been many extensions. Most of these works are for situations where there is only one set of observations (e.g., one time series sequence or one still image). However, in many applications there are multiple copies of the same or similar images, thus it is necessary to investigate denoising techniques which remove noise from multiple corrupted copies of the same signal. For a corrupted video sequence, suppose we choose a few consecutive frames in which the motion is not significant and that we have already taken care of the registration problem, one can view the frames as multiple noisy copies of the same image. Another example is when one scans a picture, but with unsatisfactory result, thus one does multiple scans, and then combines these copies to obtain the most noise-free copy possible. Since wavelet thresholding has worked well for one copy of corrupted image (cf. [1], [3], [5], and [6]), we consider its extension to multiple copies in this paper.

The standard method for combining the multiple copies is to compute their weighted average. One can only do better by incorporating

a thresholding step. The question is, which ordering is better, thresholding first or averaging first, and what is the threshold value for each ordering? The answer is not clear at all *a priori* because thresholding is a nonlinear technique that reduces variance at the expense of increasing bias. We address these issues in this paper. With the coefficients of each subband modeled as samples of a Laplacian random variable and the noise as samples of a Gaussian variable, we will show that the optimal ordering (in the mean squared error sense) depends on the number of available copies and the proportion between the noise power and the signal power. Moreover, we propose near-optimal subband adaptive thresholds for both orderings. Results show that with the optimal or the proposed near-optimal thresholds, the two methods yield very similar performance, and both outperforms weighted averaging substantially.

II. DENOISING ALGORITHM FOR MULTIPLE NOISY COPIES

Let $\mathbf{f} = \{f_{ij}\}$ denote the $M \times M$ matrix of the original image to be recovered. (Note that we use the boldfaced letters for the *matrix* of coefficients, and regular letters f_{ij} for *individual* pixels.) The signal \mathbf{f} has been transmitted over a Gaussian additive noise channel N times, and at the receiver we have N copies of noisy observations, $\mathbf{g}^{(n)} = \mathbf{f} + \varepsilon^{(n)}$, $n = 1, \dots, N$. For the n th copy, $\{\varepsilon_{ij}^{(n)}\}$ are *iid* Gaussian $N(0, \sigma_n^2)$, where σ_n^2 is the noise variance of the n th copy. The noise samples between different copies are assumed independent. The goal is to find an estimator $\hat{\mathbf{f}}$ which minimizes the mean squared error (MSE), $\text{MSE} = (1/M^2) \sum_{i,j=1}^M (\hat{f}_{ij} - f_{ij})^2$.

The recovery of the image is done in the orthogonal wavelet transform domain (the readers are referred to standard wavelet literature such as [4], [7] for details of the two-dimensional (2-D) dyadic wavelet transform). Let the wavelet transform of the noisy observation $\mathbf{g}^{(n)} = \mathbf{f} + \varepsilon^{(n)}$ be denoted by $\mathbf{Y}^{(n)} = \mathbf{X} + \mathbf{V}^{(n)}$. The wavelet coefficients are often grouped into *subbands* of different scale and orientation, with one lowest frequency subband, and the rest called *detail subbands*. It has been found that for a large class of images, the coefficients in each detail subband form a histogram well-described by a generalized Gaussian distribution (cf. [4] and [8]), often simplified to Laplacian distribution for tractability. In this work, we use the Laplacian distribution. Then it follows that the MSE is well approximated by the *expected* squared error, or *risk* under the squared loss. Thus, we wish to find the estimator $\hat{f}(g)$ of the coefficients which minimizes the risk, $\text{Risk} = E[(f - \hat{f})^2]$.

A. Wavelet Thresholding and Threshold Selection

To denoise one copy, the wavelet thresholding operation proposed by Donoho and Johnstone [3] has three steps. First, take the wavelet transform of the noisy observation \mathbf{g} to yield \mathbf{Y} . Then each coefficient Y_{ij} (except in the lowest resolution subband) is thresholded with a chosen threshold. Finally, the thresholded coefficients are transformed back to yield the recovered signal.

There are two popular thresholding functions: the *soft-threshold* function, $\eta_\lambda(t) = \text{sgn}(t) \cdot \max(0, |t| - \lambda)$, which shrinks the input toward zero by amount λ , and the *hard-threshold* function, $\psi_\lambda(t) = t \cdot \mathbf{1}\{|t| > \lambda\}$, which keeps the input only if it is above the threshold λ . Although the soft-thresholding operation tends to smooth the image slightly more than the hard-threshold function, it yields images with better visual quality especially when the noise power is significant. Furthermore, with the chosen probability distribution,

Manuscript received November 2, 1997; revised March 1, 2000. S. G. Chang was supported in part by the NSF Graduate Fellowship and the University of California Dissertation Fellowship. B. Yu was supported in part by ARO Grant DAAH04-94-G-0232 and NSF Grant DMS-9322817. M. Vetterli was supported in part by NSF Grant MIP-93-213002 and Swiss NSF 20-52347.97. Part of this work was presented at the IEEE International Conference on Image Processing (ICIP), Chicago, IL, Oct. 1998. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Patrick L. Combettes.

S. G. Chang was with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA. She is now with Hewlett-Packard Company, Grenoble, France (e-mail: grchang@yahoo.com).

B. Yu is with the Department of Statistics, University of California, Berkeley, CA 94720 USA (e-mail: binyu@stat.berkeley.edu)

M. Vetterli is with the Laboratory of Audiovisual Communications, Swiss Federal Institute of Technology (EPFL), Lausanne, Switzerland and also with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA.

Publisher Item Identifier S 1057-7149(00)06511-8.

optimal soft-thresholding yields a lower risk than optimal hard-thresholding, as was shown in [1]. Thus, soft-thresholding is the preferred operation in this work.

The next issue is the selection of the threshold value. The threshold selection is performed once for each subband. Thus, let us consider coefficients from one detail subband. As in [1], we model the wavelet coefficients in a detail subband of the original image \mathbf{f} as samples from a zero-mean Laplacian random variable with an unknown parameter. That is, each X_{ij} is modeled as $X \sim p(x) = \text{LAP}(\beta) \triangleq (1/\sqrt{2}\beta) \exp\{-(\sqrt{2}/\beta)|x|\}$. Since the noise is Gaussian and the wavelet transform is orthogonal, each wavelet coefficient, Y_{ij} , of the corrupted image has distribution $Y|X \sim p(y|x) = N(x, \sigma^2) = 1/\sqrt{2\pi\sigma^2} \exp\{-(y-x)^2/2\sigma^2\}$. We use the soft-threshold estimate, $\hat{X} = \eta_\lambda(Y)$, and the optimal threshold is defined to be

$$\begin{aligned} \lambda^* &= \arg \min_{\lambda} E_X E_{Y|X} (\hat{X} - X)^2 \\ &= \arg \min_{\lambda} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (\eta_\lambda(y) - x)^2 p(y|x) p(x) dy dx. \end{aligned} \quad (1)$$

As far as we know, λ^* does not have a closed form solution and is thus found numerically (with minimization and numerical integration techniques). A good approximation of λ^* was found in [1] to be $\tilde{\lambda}(\beta) = \sigma^2/\beta$, where β is the standard deviation of X . It results in less than 0.8% deviation from the minimum MSE (with the optimal threshold λ^*). The threshold $\tilde{\lambda}(\beta)$ is simple and effective and has an intuitive explanation. When the noise power is much smaller than the signal power, $\sigma/\beta \ll 1$, the normalized threshold $\tilde{\lambda}/\sigma$ is small to preserve most of the signal features; on the other hand, when $\sigma/\beta \gg 1$, $\tilde{\lambda}/\sigma$ is chosen to be large to remove the noise which has overwhelmed the signal. By allowing the unknown parameter β to be estimated, this method also allows a *data-driven* selection of the threshold which adapts to each subband.

B. Combining Thresholding and Averaging

When there are multiple copies available, the standard method is to use the (pixel-wise) weighted average as the estimate. Let $V^{(n)} \sim N(0, \sigma_n^2)$, $n = 1, \dots, N$, be the noise variable for the n th copy, and Z be the weighted average of $Y^{(n)} = X + V^{(n)}$, $Z = \sum_{n=1}^N \alpha_n Y^{(n)} = X + \sum_{n=1}^N \alpha_n V^{(n)}$, where $\sum \alpha_n = 1$. It is well-known that the optimal α_n are $\alpha_n^* = (1/\sigma_n^2) / \sum_{i=1}^N (1/\sigma_i^2)$, and the resulting MSE is $\sigma_{\text{total}}^2 = \text{Var}(Z - X) = \text{Var}\left(\sum_{n=1}^N \alpha_n^* V^{(n)}\right) = \left(\sum_{n=1}^N (1/\sigma_n^2)\right)^{-1}$, where $\text{Var}(\cdot)$ denotes the variance.

Now let us incorporate thresholding into averaging. Z is a new random variable with $Z|X \sim N(x, \sigma_{\text{total}}^2)$. Since this is exactly the setting for one copy thresholding, we can simply use $\sigma_{\text{total}}/\beta$ as the threshold. However, can we do better than that? More specifically, since we have two operations here, averaging and thresholding, of which one is linear and the other not, the ordering could make a difference. Thus we investigate which ordering is best in the mean squared sense.

Consider the special case when $\sigma_1 = \sigma_2 = \dots = \sigma_N \triangleq \sigma$. Thus, $\alpha_1 = \dots = \alpha_N = 1/N$. To make references more convenient, let $\mathcal{A}(\cdot)$ denote the weighted average operation and $\mathcal{T}(\cdot)$ the thresholding operation, and we give the following notation to the two orderings:

$$\mathcal{A}(\mathcal{T}(Y^{(1)}, \dots, Y^{(N)})): \hat{X}_{\mathcal{AT}}(\lambda) = \frac{1}{N} \sum_{n=1}^N \eta_\lambda(Y^{(n)}) \quad (2)$$

$$\mathcal{T}(\mathcal{A}(Y^{(1)}, \dots, Y^{(N)})): \hat{X}_{\mathcal{TA}}(\lambda) = \eta_\lambda\left(\frac{1}{N} \sum_{n=1}^N Y^{(n)}\right). \quad (3)$$

The risk of the estimator $\hat{X}_{\mathcal{AT}}(\lambda)$ is

$$R_{\mathcal{AT}}(\lambda) = E_X E_{Y^{(1)}, \dots, Y^{(N)}|X} \left(\frac{1}{N} \sum_{n=1}^N \eta_\lambda(Y^{(n)}) - X \right)^2 \quad (4)$$

$$= E_X E_{Y^{(1)}, \dots, Y^{(N)}|X} \left[\frac{1}{N} \sum_{n=1}^N (\eta_\lambda(Y^{(n)}) - X) \right]^2 \quad (5)$$

$$= E_X E_{Y^{(1)}, \dots, Y^{(N)}|X} \left[\frac{1}{N^2} \sum_{n=1}^N (\eta_\lambda(Y^{(n)}) - X)^2 + \frac{1}{N^2} \sum_{n=1}^N \sum_{\substack{m=1 \\ m \neq n}}^N (\eta_\lambda(Y^{(n)}) - X)(\eta_\lambda(Y^{(m)}) - X) \right] \quad (6)$$

$$= \frac{1}{N} E_X E_{Y|X} (\eta_\lambda(Y) - X)^2 + \frac{N-1}{N} E_X \cdot [E_{Y|X} (\eta_\lambda(Y) - X)]^2 \quad (7)$$

where (7) follows from the fact that $\{Y^{(1)}, \dots, Y^{(N)}\}$ conditioned on X are independent. The risk of $\hat{X}_{\mathcal{TA}}(\lambda)$ is

$$\begin{aligned} R_{\mathcal{TA}}(\lambda) &= E_X E_{Y^{(1)}, \dots, Y^{(N)}|X} (\hat{X}_{\mathcal{TA}}(\lambda) - X)^2 \\ &= E_X E_{Z|X} (\eta_\lambda(Z) - X)^2, \end{aligned} \quad (8)$$

where $Z = (1/N) \sum_{n=1}^N Y^{(n)}$, and $Z|X \sim N(x, \sigma^2/N)$. The optimal thresholds are $\lambda_{\mathcal{AT}}^* = \arg \min_{\lambda} R_{\mathcal{AT}}(\lambda)$ and $\lambda_{\mathcal{TA}}^* = \arg \min_{\lambda} R_{\mathcal{TA}}(\lambda)$. Note that $\lambda_{\mathcal{AT}}^*$ and $\lambda_{\mathcal{TA}}^*$ depend on β and σ . We do not have closed form solutions for $\lambda_{\mathcal{AT}}^*$ and $\lambda_{\mathcal{TA}}^*$, thus we resort to numerical methods. Without loss of generality, we can set $\sigma = 1$ (alternatively, one can solve for $\lambda_{\mathcal{AT}}^*/\sigma$ and $\lambda_{\mathcal{TA}}^*/\sigma$ as a function of β/σ). For a fixed value β_0 , $\lambda_{\mathcal{TA}}(\beta_0)$ and $\lambda_{\mathcal{AT}}(\beta_0)$ are found by locating the zero crossing of the derivatives of the risks, $R'_{\mathcal{TA}}(\lambda)$ and $R'_{\mathcal{AT}}(\lambda)$, respectively. Standard numerical integration routines (such as the Romberg integration found in numerical recipes) are used to calculate the expectations.

To compare the risks of these two methods, we look at the scaled risk difference, $(R_{\mathcal{AT}}(\lambda_{\mathcal{AT}}^*) - R_{\mathcal{TA}}(\lambda_{\mathcal{TA}}^*))/\sigma^2$, as a function of N and of the ratio β/σ , illustrated in Fig. 1. For each $N \leq 5$, there is a cutoff point C_N^* below which $R_{\mathcal{AT}}(\lambda_{\mathcal{AT}}^*) > R_{\mathcal{TA}}(\lambda_{\mathcal{TA}}^*)$, and above which $R_{\mathcal{AT}}(\lambda_{\mathcal{AT}}^*) < R_{\mathcal{TA}}(\lambda_{\mathcal{TA}}^*)$. For $N > 5$, however, the $\mathcal{T}(\mathcal{A}(\cdot))$ method is better for any value of β/σ . The cutoff points C_N^* for each N are listed in Table I. This indicates that *the best method depends on the relative power between the noise and signal, and also on the value of N* . With the optimal thresholds, the improvement of one method over the other is small, however, on the order of $10^{-3}\sigma^2$. $\mathcal{T}(\mathcal{A}(\cdot))$ requires much less computation than $\mathcal{A}(\mathcal{T}(\cdot))$ (since the former can be implemented by computing the wavelet transform once, whereas the latter computes it N times). Thus if computation is an issue, $\mathcal{T}(\mathcal{A}(\cdot))$ is preferred. However, thresholding each individual copy of the image may be advantageous. It is possible that a different noisy copy is collected and processed at each receiving station, and only this processed copy is kept. At a later time, these separately processed copies can be collected by a central receiver to yield one better copy.

By plotting the numerically found $\lambda_{\mathcal{TA}}^*/\sigma$ and $\lambda_{\mathcal{AT}}^*/\sigma$ against β/σ and N , we discover that there is a simple analytical expression which well approximates $\lambda_{\mathcal{TA}}^*$ and $\lambda_{\mathcal{AT}}^*$. For the $\mathcal{T}(\mathcal{A}(\cdot))$ estimator, the threshold is simply a modification of $\tilde{\lambda}$ for one copy denoising, but with a change in the noise variance

$$\tilde{\lambda}_{\mathcal{TA}} = \frac{\sigma^2/N}{\beta}. \quad (9)$$

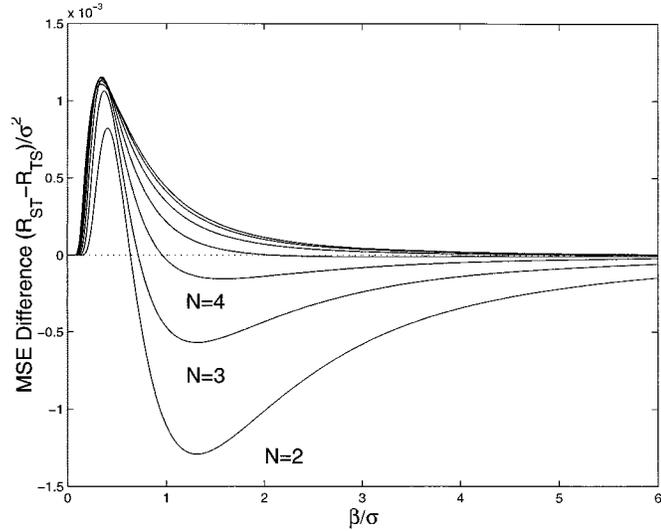


Fig. 1. Scaled MSE difference $(R_{\mathcal{A}T}(\lambda_{\mathcal{A}T}^*) - R_{T\mathcal{A}}(\lambda_{T\mathcal{A}}^*)) / \sigma^2$ as a function of N and β/σ .

TABLE I
CUTOFF VALUES (IN UNIT σ_x/σ) FOR EACH N , WHERE C_N^* IS THE CUTOFF VALUE FOR WHEN USING THE OPTIMAL THRESHOLDS, AND \tilde{C}_N (LISTED ONLY FOR $N \leq 5$) IS THE CUTOFF VALUE WHEN USING THE PROPOSED THRESHOLDS, $\tilde{T}_{\mathcal{A}T}$ AND $\tilde{T}_{T\mathcal{A}}$

	C_N^*	\tilde{C}_N
$N = 2$.6367	.1379
$N = 3$.7154	.7654
$N = 4$.9601	.9466
$N = 5$	1.9768	1.0884
$N > 5$	∞	> 1.23

For the $\mathcal{A}(T(\cdot))$ method, we found the approximation

$$\tilde{\lambda}_{\mathcal{A}T} = \frac{\sigma^2 / N^{(3/4)}}{\beta}. \quad (10)$$

Fig. 2 compares the optimal and approximate thresholds for both methods as a function of N , for $\sigma = 1$ and $\beta = 1$. The thresholds $\tilde{\lambda}_{T\mathcal{A}}$ and $\tilde{\lambda}_{\mathcal{A}T}$ result in less than 0.2% deviation from the minimum MSE (with $\lambda_{T\mathcal{A}}^*$ and $\lambda_{\mathcal{A}T}^*$, respectively). Fig. 3 compares the optimal threshold $\lambda_{\mathcal{A}T}^*$ and the approximation $\tilde{\lambda}_{\mathcal{A}T}$ (scaled by $1/\sigma$) as a function of β/σ for $N = 2, \dots, 6$. The MSE due to $\tilde{\lambda}_{\mathcal{A}T}$ deviates from the optimal MSE by less than 3.5% for $\beta/\sigma < 1$ and less than 0.1% for $\beta/\sigma > 1$. Since typically the signal power is much larger than the noise power (otherwise the image features has been greatly corrupted), inaccurate approximations for small β/σ are acceptable. The thresholds $\tilde{\lambda}_{T\mathcal{A}}$ and $\tilde{\lambda}_{\mathcal{A}T}$ also yield a different set of cutoff values \tilde{C}_N (tabulated in Table I for $N \leq 5$). Notably for $N > 5$, there are some \tilde{C}_N less than ∞ , but even above these \tilde{C}_N , the MSE of the two methods are so close that either one can be chosen. The scaled MSE difference $(R_{\mathcal{A}T}(\tilde{\lambda}_{\mathcal{A}T}) - R_{T\mathcal{A}}(\tilde{\lambda}_{T\mathcal{A}})) / \sigma^2$ is similar to the curves shown in Fig. 2 for optimal thresholds and is of the same order of magnitude. Thus, the use of $\tilde{\lambda}_{T\mathcal{A}}$ and $\tilde{\lambda}_{\mathcal{A}T}$ do not change the previous conclusions.

Notice that the threshold for $\mathcal{A}(T(\cdot))$ decreases as N increases, even though at the thresholding stage, each copy is thresholded independently of the other copies. To explain this dependency on N , we rewrite the inner expectation of $R_{\mathcal{A}T}(\lambda)$

$$E_{Y^{(1)}, \dots, Y^{(N)} | X}(\hat{X}_{\mathcal{A}T}(\lambda) - X)^2 \quad (11)$$

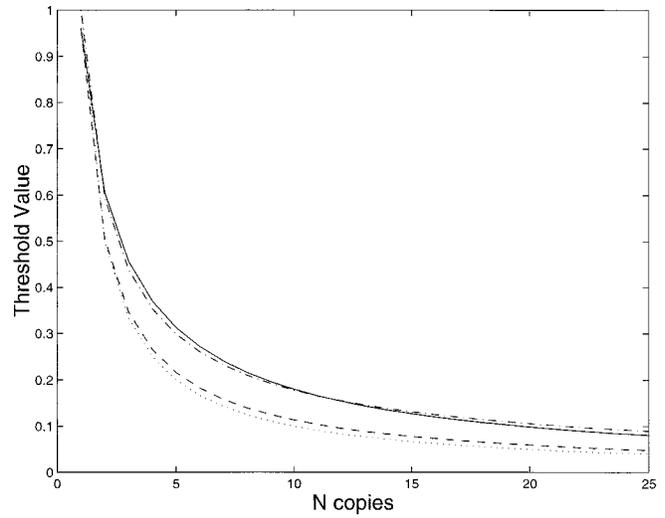


Fig. 2. Comparing $\lambda_{T\mathcal{A}}^*$ (---) versus $\tilde{\lambda}_{T\mathcal{A}}$ (···), and $\lambda_{\mathcal{A}T}^*$ (—) versus $\tilde{\lambda}_{\mathcal{A}T}$ (-·-·), when $\sigma = 1$ and $\beta = 1$.

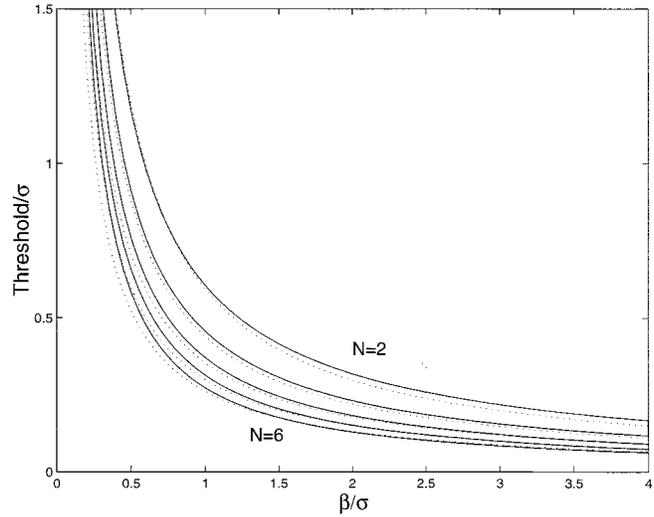


Fig. 3. Comparing $\lambda_{\mathcal{A}T}^*$ (—) and $\tilde{\lambda}_{\mathcal{A}T}$ (···) for $\sigma_1 = \sigma_2 = \dots = \sigma_N \triangleq \sigma$ as a function of β/σ and $N = 2, \dots, 6$.

$$= E_{\{Y^{(n)}\} | X} [\hat{X}_{\mathcal{A}T}(\lambda) - E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda) + E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda) - X]^2 \quad (12)$$

$$= E_{\{Y^{(n)}\} | X} [\hat{X}_{\mathcal{A}T}(\lambda) - E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda)]^2 + (E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda) - X)^2 + 2E_{\{Y^{(n)}\} | X} \cdot [(\hat{X}_{\mathcal{A}T}(\lambda) - E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda)) \cdot (E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda) - X)] \quad (13)$$

$$= E_{\{Y^{(n)}\} | X} [\hat{X}_{\mathcal{A}T}(\lambda) - E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda)]^2 + (E_{\{Y^{(n)}\} | X} \hat{X}_{\mathcal{A}T}(\lambda) - X)^2 \quad (14)$$

$$= E_{\{Y^{(n)}\} | X} \left[\frac{1}{N} \sum_n \eta_\lambda(Y^{(n)}) - E_{Y | X} \eta_\lambda(Y) \right]^2 + (E_{Y | X} \eta_\lambda(Y) - X)^2 \quad (15)$$

$$= \frac{1}{N} E_{Y | X} (\eta_\lambda(Y) - E_{Y | X} \eta_\lambda(Y))^2 + (E_{Y | X} \eta_\lambda(Y) - X)^2 \quad (16)$$

where $\{Y^{(n)}\}$ has been used as a shorthand for $Y^{(1)}, \dots, Y^{(N)}$. The first term is the variance from thresholding, while the second term is the square of the bias. The optimal threshold is obtained from the tradeoff between the variance term (which decreases with increasing λ) and the bias term (which increases with increasing λ). As N becomes large, the variance term decreases due to the $1/N$ factor while the bias term stays the same. Thus, λ needs to be decreased as well to obtain the minimum total.

Up to now we have assumed the knowledge of the noise variance σ^2 and the standard deviation, β , of X . In practice, these two values may not be known and need to be estimated from the noisy observations. For both methods, these two parameters are estimated the same way for a fair comparison. First the noise variance σ_n^2 is estimated by the robust median estimator in the highest subband (also used in [3]), $\hat{\sigma}_n = \text{Median}(|Y_{ij}^{(n)}|)/0.6745$, with all $Y_{ij}^{(n)}$ in the HH₁ subband of the n th copy, then $\hat{\sigma}^2$ is taken to be the average of these N estimates. Since the noise is independent from the signal, $\text{Var}(Z) = \text{Var}(X) + \sigma^2/N = \beta^2 + \sigma^2/N$. Thus, for each subband of $Z = (1/N) \sum_{n=1}^N Y^{(n)}$, the sample variance estimate of $\text{Var}(Z)$, $\hat{\sigma}_Z^2 = \text{Average}(Z_{ij} - \text{Mean}(Z_{ij}))^2$, is calculated, and the estimate of the standard deviation of the Laplacian distribution is

$$\hat{\beta} = \sqrt{\max(0, \hat{\sigma}_Z^2 - \hat{\sigma}^2/N)}. \quad (17)$$

Heterogeneous Noise Variances: Now consider the case when the noise variances σ_n^2 are different. This extension is straightforward in the $T(\mathcal{A}(\cdot))$ case. The multiple copies are averaged with coefficients α_n^* , and the threshold is $\tilde{\lambda}_{T\mathcal{A}}$ in (9) but with σ^2/N replaced by σ_{total}^2 .

For the $\mathcal{A}(T(\cdot))$ method, one needs to find the optimal threshold λ_n for each copy and the optimal weights α_n . By minimizing the risk $E_X E_{Y^{(1)}, \dots, Y^{(N)}|X} \left(\sum_{n=1}^N \alpha_n \eta_{\lambda_n} (Y^{(n)} - X) \right)^2$ with respect to $\alpha_1, \dots, \alpha_N$ subject to $\sum \alpha_n = 1$, and also with respect to $\lambda_1, \dots, \lambda_N$, one can find the optimal values. The optimal α_n are found numerically to be almost identical to the values in α_n^* . We do not have the closed form solution of the optimal thresholds. Thus, we approximate them by $\tilde{\lambda}_{\mathcal{A}T}^{(n)} = \sigma_n^{1/2}/\beta \left(1 / \left(\sum_{i=1}^N (1/\sigma_i^2) \right) \right)^{3/4}$, $n = 1, \dots, N$, which yields $\tilde{\lambda}_{\mathcal{A}T}$ in (10) when $\sigma_1 = \sigma_2 = \dots = \sigma_N$. Fig. 4 compares the optimal thresholds (—) and the thresholds $\tilde{\lambda}_{\mathcal{A}T}^{(n)}$ (· · ·) for $N = 3$ and $\{\sigma_n\} = \{1, 3, 5\}$, against β on the horizontal axis. The plot shows the best fit for the threshold corresponding to $\min\{\sigma_n\}$, and the approximation worsens for thresholds whose corresponding σ_n is further from $\min\{\sigma_n\}$, especially in the region of small β (relative to $\min\{\sigma_n\}$). This inaccuracy is mitigated by the fact that the α_n^* 's for large σ_n 's are small, thus the overall MSE is still close to the optimal MSE (less than 0.3% for $\beta > 1$ for the graph shown).

III. EXPERIMENTAL RESULTS

To validate our proposed methods, we take as the test image a 256×256 block from the grayscale images *Barbara* and *Lena*, with $\sigma_1 = \sigma_2 = \dots = \sigma_N = \sigma = 30$, using Daubechies' least asymmetric wavelet with eight vanishing moments (i.e., Symmlet8) and four scales of wavelet transform. The algorithm is not too sensitive to the choices of wavelet and number of levels, as long as the wavelet is smooth enough (has several vanishing moments) and at least three or four levels are used for this magnitude of σ (for example, using Symmlet4 increases the MSE by 1%–2% and using two levels of decomposition increases the MSE by 2%–4%). The parameters β and σ are estimated as discussed previously. We compare the MSE's of five methods for N ranging from 1 to 25:

- 1) averaging;
- 2) $\mathcal{A}(T(\cdot))$;
- 3) $T(\mathcal{A}(\cdot))$;

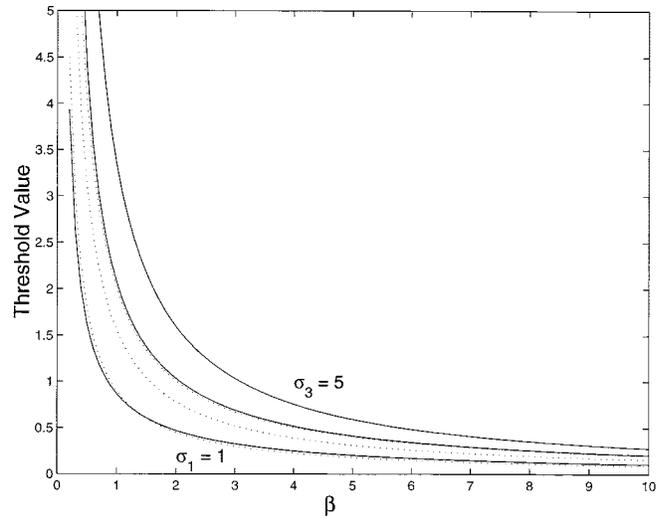


Fig. 4. For $N = 3$ and different noise σ_n ($\{1, 3, 5\}$), compare the optimal threshold for $\mathcal{A}(T(\cdot))$ (—) and $\tilde{\lambda}_{\mathcal{A}T}^{(n)}$ (· · ·).

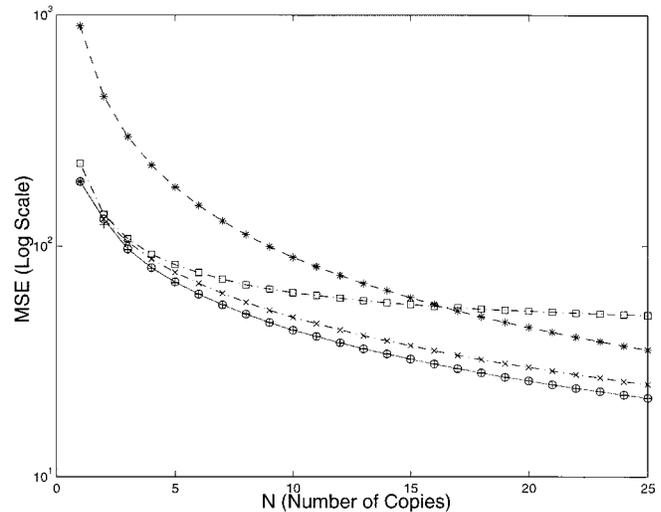


Fig. 5. For the image *Barbara*, comparing as a function of N the MSE of averaging (*), $\mathcal{A}(T(\cdot))$ (x), $T(\mathcal{A}(\cdot))$ (+), switching (o), and Wiener filtering (□), for $\sigma = 30$. Note that the (+) and (o) curves are overlapping.

- 4) switching between the two thresholding methods (only for $N \leq 5$) with cutoff values \tilde{C}_N (thus the switching method becomes $T(\mathcal{A}(\cdot))$ for $N > 5$);
- 5) Wiener filtering of the averaged copy [from (a)].

The Wiener filter is the “wiener2” routine from the MATLAB image processing toolkit, with the input noise power unknown and the adaptation window size set to the default (3×3). The resulting MSE's are shown in Fig. 5 for *Barbara* and Fig. 6 for *Lena*. The three thresholding methods show significant improvement over merely averaging, ranging from 70% to 30% reduction in MSE for N varying from 2 to 25, and are either comparable or up to 5%–15% improvement over Wiener filtering (depending on N). The removal of noise due to thresholding is also significant visually (see Fig. 7 for *barbara* at $N = 4$), especially for small N . For $N = 1$ and $\sigma = 30$, Wiener filtering does poorly because of the large noise power. For larger N , the averaging reduces the noise power substantially, and thus Wiener filtering yields images comparable with those from thresholding, though slightly more blurry. Among the thresholding methods, the $T(\mathcal{A}(\cdot))$ method is the best in terms of MSE, even better than switching, suggesting that perhaps the $\mathcal{A}(T(\cdot))$ method

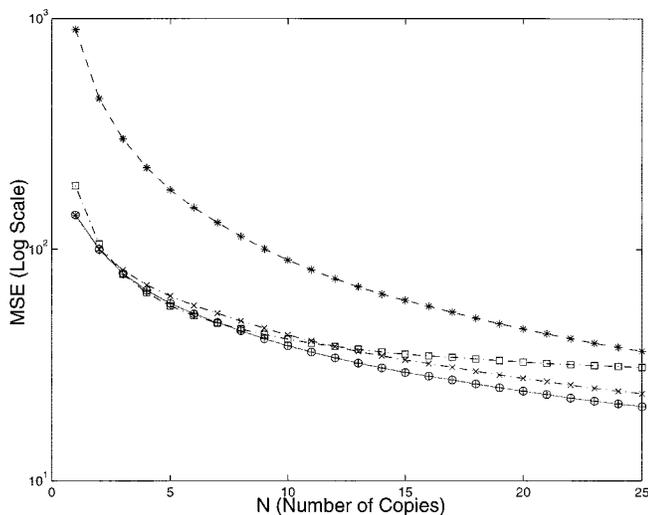


Fig. 6. For the image *Lena*, comparing as a function of N the MSE of averaging (*), $\mathcal{A}(T(\cdot))$ (x), $T(\mathcal{A}(\cdot))$ (+), switching (o), and Wiener filtering (□), for $\sigma = 30$. Note that the (+) and (o) curves are overlapping.

is more sensitive to model errors and threshold estimation errors. For large values of σ (e.g., 30), $\mathcal{A}(T(\cdot))$ produces noticeably more noisy images than $T(\mathcal{A}(\cdot))$, though for smaller values, it yields images of comparable quality as those from $T(\mathcal{A}(\cdot))$. For $1 < N \leq 5$, the switching method yields MSE's that are between those of $\mathcal{A}(T(\cdot))$ and $T(\mathcal{A}(\cdot))$. The $T(\mathcal{A}(\cdot))$ method requires the least amount of computation since it can be implemented with only one wavelet transform and seems to work well for large values of σ as well. Thus, in practice, the $T(\mathcal{A}(\cdot))$ method suffices to combine multiple noisy copies.

It is curious to investigate if an additional stage of thresholding (i.e., performing $T(\mathcal{A}(T(\cdot)))$) can have a significant improvement. It cannot do worse, since we can always choose the latter stage threshold to be zero. To test this idea, we take the output of $\mathcal{A}(T(\cdot))$ and optimally threshold it assuming that we have the original. The resulting MSE is only slightly better than the $T(\mathcal{A}(\cdot))$, suggesting that thresholding of the weighted average yields a sufficiently denoised image already. Furthermore, finding the optimal thresholds of a two-stage thresholding operation is difficult.

IV. CONCLUSION

In this paper, we addressed the issue of image recovery from its multiple noisy copies. We explored the idea of combining the wavelet thresholding technique with the more traditional averaging operation. Our investigation showed that the optimal ordering of these two operations is not so straightforward and is in fact a function of the number of available copies and of the relative energy between noise and signal. We proposed a near-optimal threshold for each ordering. With these thresholds, the performances are similar, and for computational reasons, averaging followed by thresholding is recommended. Furthermore, all of these thresholding methods show substantial improvement over mere averaging and moderate improvement over Wiener filtering, both visually and in the MSE sense.

REFERENCES

[1] S. G. Chang, B. Yu, and M. Vetterli, "Image denoising via lossy compression and wavelet thresholding," in *Proc. IEEE Int. Conf. Image Processing*, vol. 1, Nov. 1997, pp. 604–607.
 [2] I. Daubechies, *Ten Lectures on Wavelets*. Philadelphia, PA: SIAM, 1992.



Fig. 7. Comparison of denoised images, for $N = 4$: (a) original, (b) noisy image with $\sigma = 30$, (c) averaging, (d) switching, (e) $\mathcal{A}(T(\cdot))$, (f) $T(\mathcal{A}(\cdot))$, and (g) Wiener filtering. This image can be found at <http://www-wavelet.eecs.berkeley.edu/~grchang/multiThresh/>.

[3] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation via wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455, 1994.
 [4] S. Mallat, "A theory for multiresolution signal decomposition: The wavelet representation," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, pp. 674–693, July 1989.
 [5] S. Mallat, *A Wavelet Tour of Signal Processing*. New York: Academic, 1998.
 [6] P. Moulin and J. Liu, "Analysis of multiresolution image denoising schemes using generalized-Gaussian and complexity priors," *IEEE Trans. Inform. Theory*, vol. 45, pp. 909–919, Apr. 1999.
 [7] M. Vetterli and J. Kovačević, *Wavelets and Subband Coding*. Englewood Cliffs, NJ: Prentice-Hall, 1995.
 [8] P. H. Westerink, J. Biemond, and D. E. Boekee, "An optimal bit allocation algorithm for sub-band coding," in *Proc. Int. Conf. Acoustic Speech Signal Processing*, Dallas, TX, Apr. 1987, pp. 1378–1381.