

Denoising using wavelets

WTBV

December 19, 2017

- 1 Threshold functions
- 2 Wavelet shrinking
- 3 The VisuShrink method
- 4 The SURE method

- Threshold functions $s_\lambda(t)$
 - are used to suppress parts of a signal with very low amplitudes and (usually) high frequencies (“noise”)
 - examples are

$$\{\text{Hard}, \lambda\} \quad \begin{cases} 0 & |x| \leq \lambda \\ x & |x| > \lambda \end{cases}$$

$$\{\text{Soft}, \lambda\} \quad \begin{cases} 0 & |x| \leq \lambda \\ \text{sgn}(x)(|x| - \lambda) & |x| > \lambda \end{cases}$$

$$\{\text{PiecewiseGarrote}, \lambda\} \quad \begin{cases} 0 & |x| \leq \lambda \\ x - \frac{\lambda^2}{x} & |x| > \lambda \end{cases}$$

$$\{\text{SmoothGarrote}, \lambda, n\} \quad \frac{x^{2n+1}}{x^{2n} + \lambda^{2n}}$$

$$\{\text{Hyperbola}, \lambda\} \quad \begin{cases} 0 & |x| \leq \lambda \\ \text{sgn}(x)\sqrt{x^2 - \lambda^2} & |x| > \lambda \end{cases}$$

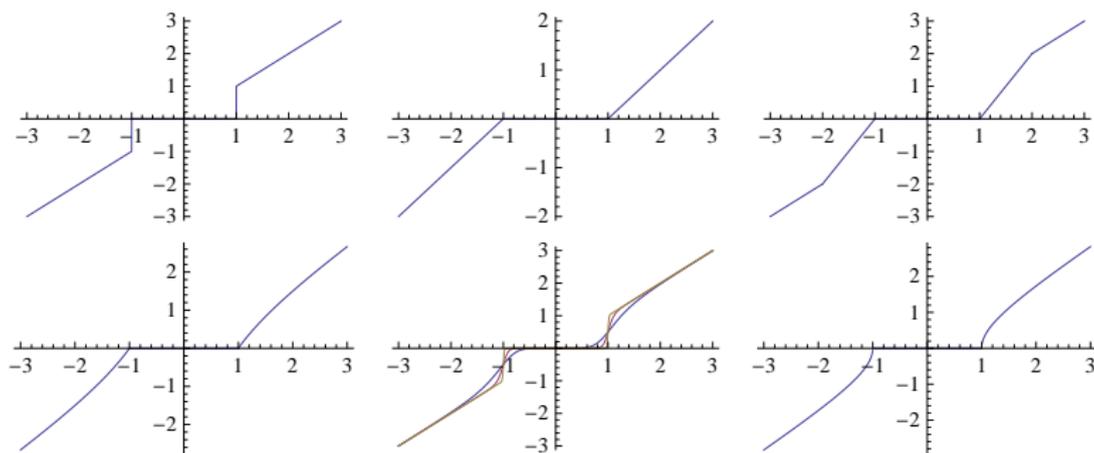


Figure: Examples of threshold functions $s_\lambda(t)$

- Setting and strategy

- “true” signal : $\mathbf{v} = (v_1, v_2, \dots, v_N)$
- noise vector: $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N)$
- noised signal $\mathbf{y} = (y_1, y_2, \dots, y_N) = \mathbf{v} + \boldsymbol{\varepsilon}$
- wavelet filtering (orthogonal transform!)

$$\mathbf{y} \xrightarrow{WT} \mathbf{z} = (\mathbf{a}, \mathbf{d}) = (H\mathbf{y}, G\mathbf{y})$$

- applying thresholding with $s_\lambda(t)$ to the high-pass component

$$\mathbf{z} \mapsto \widehat{\mathbf{z}} = (\mathbf{a}, \widehat{\mathbf{d}}) \quad \text{with} \quad \widehat{\mathbf{d}} = s_\lambda(\mathbf{d})$$

- inverse wavelet transform

$$\widehat{\mathbf{z}} \xrightarrow{WT^{-1}} \widehat{\mathbf{v}} = H^\dagger \mathbf{a} + G^\dagger \widehat{\mathbf{d}}$$

- Heuristic considerations

- Noise modelled as *Gaussian white noise* with noise level σ :
 $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N)$ generated by independent and identically $\mathcal{N}(0, \sigma^2)$ -distributed random variables
- For a vector $\varepsilon = (\varepsilon_1, \dots, \varepsilon_N)$ of independent, $\mathcal{N}(0, \sigma^2)$ -distributed random variables ε_i and an orthogonal ($N \times N$) matrix U the components γ_i of the vector

$$\gamma = (\gamma_1, \dots, \gamma_N) = U\varepsilon,$$

are again independent $\mathcal{N}(0, \sigma^2)$ -distributed random variables

- Due to orthogonality of the wavelet transform the transformed noise term $WT(\varepsilon)$ in

$$\mathbf{y} = \mathbf{v} + \varepsilon \mapsto WT(\mathbf{y}) = WT(\mathbf{v}) + WT(\varepsilon)$$

is still characterized by being white noise with noise level σ

- Heuristic considerations (contd.)
 - In wavelet transformations most energy goes into the approximation (low-pass) component \mathbf{a}
 - Noise of high frequency goes into the detail (high-pass) component \mathbf{d}
 - \implies the detail component mainly (but not exclusively) consists of noise (detail coefficients $\lesssim \sigma$) — that is where to attack!
 - The problem: the noise level σ is not known and has to be estimated from the data to be denoised themselves
 - How to choose λ (depending on the estimate for σ) ?
 - Measure of quality for denoising: *mean squared error* (MSE)

$$E [\|\mathbf{v} - \hat{\mathbf{v}}\|^2] = E \left[\sum_{1 \leq j \leq N} (v_j - \hat{v}_j)^2 \right]$$

- $\mathbf{v} = (v_1, \dots, v_N) \in \mathbb{R}^N$
- $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)$ white noise with variance σ^2 , $\mathbf{y} = \mathbf{v} + \boldsymbol{\varepsilon}$
- $\hat{\mathbf{v}} = (\hat{v}_1, \dots, \hat{v}_N)$ estimate for \mathbf{v} with $A \subseteq \{1, 2, \dots, N\}$ and

$$\hat{v}_j = \begin{cases} y_j & \text{if } j \in A \\ 0 & \text{if } j \notin A \end{cases}$$

- In this case the MSE equals

$$\mathbb{E} [\|\mathbf{v} - \hat{\mathbf{v}}\|^2] = \mathbb{E} \left[\sum_{1 \leq j \leq N} (v_j - \hat{v}_j)^2 \right] = \sum_{j \in A} \mathbb{E} [\varepsilon_j^2] + \sum_{j \notin A} \mathbb{E} [v_j^2]$$

and this is minimized by setting $j \in A \iff v_j^2 > \sigma^2$

- so that the ideal MSE is

$$\mathbb{E} [\|\mathbf{v} - \hat{\mathbf{v}}\|^2] = \sum_{1 \leq j \leq N} \min(v_j^2, \sigma^2)$$

- VISUSHRINK :

- is Wavelet shrinkage with

$$\lambda = \lambda^{\text{univ}} = \sigma \cdot \sqrt{2 \ln N} \quad (\text{"universal tolerance"})$$

- Theorem [DONOHO-JOHNSTONE, 1995]

For $\mathbf{v} \in \mathbb{R}^N$ and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_N)$ white noise with noise level σ , for $\hat{\mathbf{v}} = s_\lambda(\mathbf{v} + \boldsymbol{\varepsilon})$ (soft threshold) one gets

$$E [\|\mathbf{v} - \hat{\mathbf{v}}\|^2] \leq (2 \ln N + 1) \cdot \left(\sigma^2 + \sum_{1 \leq j \leq N} \min(v_j^2, \sigma^2) \right)$$

- σ will be estimated (on the first high-pass component!) using the *mean absolute deviation* (MAD):

- $\mathbf{w} = (w_1, \dots, w_N)$
- $\tilde{\mathbf{w}} = \text{Median of } \mathbf{w}$
- $\mathbf{v} = (|w_1 - \tilde{\mathbf{w}}|, \dots, |w_N - \tilde{\mathbf{w}}|)$
- $\text{MAD}(\mathbf{w}) = \text{Median of } \mathbf{v} = \tilde{\mathbf{v}}$

- Theorem [HAMPEL, 1974]

$$\text{MAD}(\mathbf{w}) \approx 0.6745 \cdot \sigma$$

- SURE method (STEINS unbiased risk estimator, 1981)
- Goal: choice of the λ parameter for soft-shrinking methods

$$s_{\lambda}(x) = \begin{cases} x - \lambda & \text{if } x > \lambda \\ 0 & \text{if } |x| \leq \lambda \\ x + \lambda & \text{if } x < -\lambda \end{cases}$$

- C.M. STEIN, Estimation of the mean of a multivariate normal distribution, *Ann. Stat.* 1981.
- D. DONOHO, I. JOHNSTONE, Adapting to unknown smoothness via wavelet shrinkage, *J. Amer. Stat. Assoc.* 1995.
- P. VAN FLEET, *Discrete Wavelet Transformations*, Wiley, 2008 (ch. 9).

- Lemma:

For a $\mathcal{N}(0, \sigma^2)$ -distributed random variable ε , any $z \in \mathbb{R}$ and any piecewise differentiable function $g : \mathbb{R} \rightarrow \mathbb{R}$ one has

$$\mathbb{E} [\varepsilon \cdot g(z + \varepsilon)] = \sigma^2 \cdot \mathbb{E} [g'(z + \varepsilon)]$$

- This follows from partial integration:

$$\begin{aligned} \mathbb{E} [\varepsilon \cdot g(z + \varepsilon)] &= \frac{1}{\sqrt{2\pi\sigma^2}} \int x \cdot g(z + x) \cdot e^{-\frac{x^2}{2\sigma^2}} dx \\ &= \frac{-1}{\sqrt{2\pi\sigma^2}} \int \sigma^2 \cdot g(y) \cdot \frac{d}{dy} e^{-\frac{(y-z)^2}{2\sigma^2}} dy \\ &= \sigma^2 \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \int \frac{d}{dy} g(y) \cdot e^{-\frac{(y-z)^2}{2\sigma^2}} dy \\ &= \sigma^2 \cdot \mathbb{E} [g'(z + \varepsilon)] \end{aligned}$$

- Consequence:

With z, ε, g as in the Lemma, one gets for the MSE of the estimate $\hat{z} = w + g(w)$ of the random variable $w = z + \varepsilon$:

$$\begin{aligned} \mathbb{E} [(\hat{z} - z)^2] &= \mathbb{E} [(\varepsilon + g(z + \varepsilon))^2] \\ &= \mathbb{E} [(\varepsilon^2 + 2\varepsilon \cdot g(z + \varepsilon) + g(z + \varepsilon)^2)] \\ &= \mathbb{E} [\sigma^2 + 2\sigma^2 \cdot g'(w) + g(w)^2] \end{aligned}$$

- Special case: soft-shrinking with threshold value λ
 - The function is

$$g(z) = \begin{cases} -z & \text{if } |z| < \lambda \\ -\lambda \operatorname{sgn}(z) & \text{if } |z| \geq \lambda \end{cases}$$

and thus

$$\frac{d}{dz}g(z) = \begin{cases} -1 & \text{if } |z| < \lambda \\ 0 & \text{if } |z| \geq \lambda \end{cases}$$

- Therefore

$$\sigma^2 + 2\sigma^2 g'(w) + g(w)^2 = \begin{cases} w^2 - \sigma^2 & \text{if } |w| < \lambda \\ \sigma^2 + \lambda^2 & \text{if } |w| \geq \lambda \end{cases}$$

- General situation:

- $\mathbf{z} = (z_1, z_2, \dots, z_N) \in \mathbb{R}^N$
- $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N)$ vector of $\mathcal{N}(0, \sigma_k^2)$ -distributed random variables (noise, not necessarily independent and identically distributed)
- $\boldsymbol{\sigma} = (\sigma_1, \sigma_2, \dots, \sigma_N)$
- $\mathbf{w} = \mathbf{z} + \boldsymbol{\varepsilon}$ noised vector
- $\mathbf{g} = (g_1, g_2, \dots, g_N)$ with $g_k : \mathbb{R}^N \rightarrow \mathbb{R}$ correcting functions
- $\hat{\mathbf{z}} = \mathbf{w} + \mathbf{g}(\mathbf{w})$ estimate for \mathbf{z}

- Theorem (STEIN)

With the notions just introduced, the MSE can be written as

$$\begin{aligned} \mathbb{E} [\|\hat{\mathbf{z}} - \mathbf{z}\|^2] &= \mathbb{E} \left[\sum_{1 \leq j \leq N} \left(\sigma_j^2 + 2\sigma_j^2 \frac{\partial}{\partial w_j} g_j(\mathbf{w}) + g_j(\mathbf{w})^2 \right) \right] \\ &= \|\boldsymbol{\sigma}\|^2 + 2 \sum_{1 \leq j \leq N} \sigma_j^2 \mathbb{E} \left[\frac{\partial}{\partial w_j} g_j(\mathbf{w}) \right] + \mathbb{E} [\|\mathbf{g}(\mathbf{w})\|^2] \end{aligned}$$

- Special case: soft-shrinking with threshold value λ

$$\begin{aligned} \mathbb{E} [\|\hat{\mathbf{z}} - \mathbf{z}\|^2] &= \mathbb{E} \left[\sum_{1 \leq j \leq N} \left(w_j^2 - \sigma_j^2 + (2\sigma_j^2 - w_j^2 + \lambda^2) \chi_{|w_j| \geq \lambda} \right) \right] \\ &= \mathbb{E} [\|\mathbf{w} - \boldsymbol{\sigma}\|^2] + \mathbb{E} \left[\sum_{1 \leq j \leq N} (2\sigma_j^2 - w_j^2 + \lambda^2) \chi_{|w_j| \geq \lambda} \right] \end{aligned}$$

- The left summand is independent of λ .
Minimizing the MSE can be achieved by choosing λ depending on the sample vector \mathbf{w} so that the integrand in the second summand is minimal!

$$f(\lambda) = \sum_{1 \leq j \leq N} (2\sigma_j^2 - w_j^2 + \lambda^2) \chi_{|w_j| \geq \lambda}$$

- Instead of dealing with $f(\lambda)$ it is more convenient to consider

$$\tilde{f}(\lambda) = \sum_{1 \leq j \leq N} (2\sigma_j^2 - w_j^2 + \lambda^2) \chi_{|w_j| > \lambda}$$

which changes nothing as far as the expectation is concerned

- Assume that the components of the vector $\mathbf{w} = (w_1, \dots, w_N)$ are ordered by increasing absolute value

$$|w_0| = 0 \leq |w_1| \leq |w_2| \leq \dots \leq |w_N|$$

- The function $\tilde{f}(\lambda)$ is continuous from the right for $\lambda \in \mathbb{R}_+$, and if $|w_j| < |w_{j+1}|$ holds, then $\tilde{f}(\lambda)$ is strictly increasing on the half-open interval $[|w_j|, |w_{j+1}|)$, so that it takes its minimum at $|w_j|$:

$$\tilde{f}(|w_j|) = (N - j)w_j^2 + \sum_{j+1 \leq k \leq N} (2\sigma_k^2 - w_k^2)$$

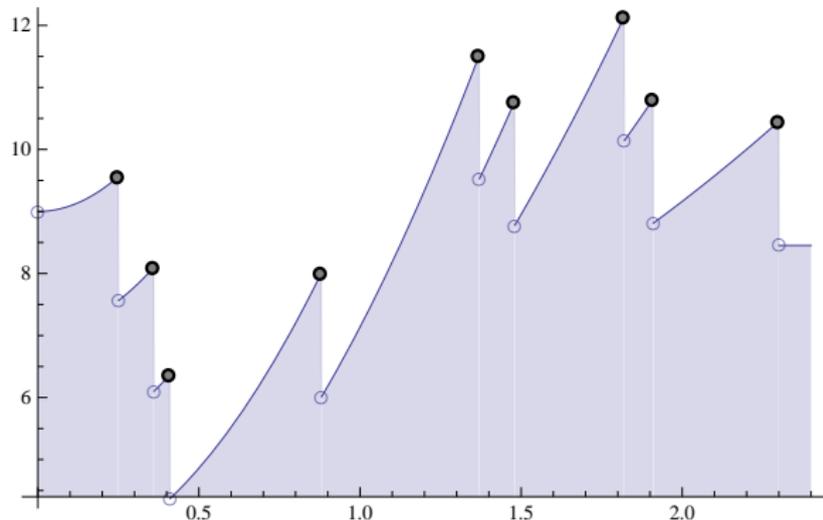


Figure: Example for the computation of the minimum value of $f(\lambda)$ for the sequence $\{0.25, 0.36, 0.41, 0.88, 1.37, 1.48, 1.82, 1.91, 2.3\}$

- From this one gets

$$\min_{\lambda \in \mathbb{R}_+} \tilde{f}(\lambda) = \min_{0 \leq j \leq N} \tilde{f}(|w_j|) \quad \lambda^{sure} := \operatorname{argmin}_{|w_j|} \tilde{f}(|w_j|)$$

- The definition of \tilde{f} yields a (downward) recursion

$$\tilde{f}(|w_j|) = \tilde{f}(|w_{j+1}|) + 2\sigma_{j+1}^2 + (N - j)(w_j^2 - w_{j+1}^2)$$

which starting from

$$\tilde{f}(w_N) = 0$$

gives a fast computation of the minimum!

- The usual assumption $\sigma_1 = \dots = \sigma_N$ somewhat simplifies the formulas and the computations

SURE in the context of the wavelet transform

- Consider a one-level WT

$$\mathbf{z} \mapsto \mathbf{y} = \mathbf{z} + \varepsilon \xrightarrow{WT} (\mathbf{a}, \mathbf{d}) \mapsto (\mathbf{a}, \hat{\mathbf{d}}) \xrightarrow{WT^{-1}} \hat{\mathbf{z}}$$

where

$$\mathbf{d} = G(\mathbf{z} + \varepsilon) = G\mathbf{z} + G\varepsilon \xrightarrow{s_\lambda} \hat{\mathbf{d}}$$

Note: the high-pass component $G\varepsilon$ has the same noise as ε

- σ must be estimated beforehand (as in the VISUSHRINK method)
- One has $\lambda^{sure} \leq \lambda^{univ}$. Recommendation: If \mathbf{y} is sparse, it is better to use λ^{univ} instead of λ^{sure} , based on the criterion

$$\frac{1}{N} \sum_{1 \leq j \leq n} (y_j^2 - \sigma^2) \leq \frac{3}{2\sqrt{N}} \log_2(N)$$

(DONOHOE, JOHNSTONE)