

Agenda

- Lecture 1: Sparse representations and what they are good for
- Lecture 2: Overview of Compressive Imaging
- Lecture 3: Proof of uncertainty principle \Rightarrow stable recovery via ℓ_1 (elementary/nontrivial as RV said . . .)
- We will discuss many things in the context of *imaging* to keep things concrete

Sparsity in Compression, Denoising, and Inverse Problems

(Thanks to Emmanuel Candes for some of the slides)

Applied and Computational Harmonic Analysis

- Signal/image $f(t)$ in the time/spatial domain
- Decompose f as a superposition of atoms

$$f(t) = \sum_i \alpha_i \psi_i(t)$$

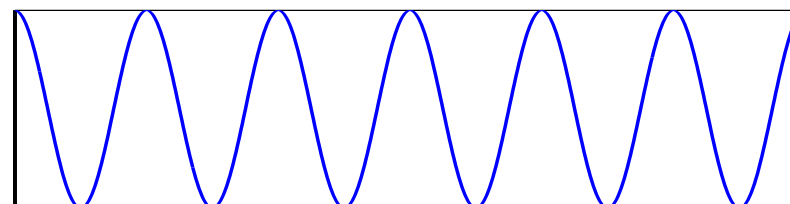
ψ_i = basis functions

α_i = expansion coefficients in ψ -domain

- Classical example: **Fourier series**

ψ_i = complex sinusoids

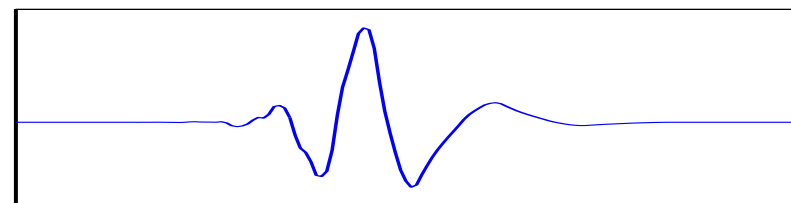
α_i = Fourier coefficients



- Modern example: **wavelets**

ψ_i = “little waves”

α_i = wavelet coefficients



- Cutting-edge example: **curvelets** (more later)

ACHA

- ACHA Mission: construct “good representations” for “signals/images” of interest
- Examples of “signals/images” of interest
 - Classical: signal/image is “bandlimited” or “low-pass”
 - Modern: smooth between isolated singularities (e.g. 1D piecewise poly)
 - Cutting-edge: 2D image is smooth between smooth edge contours
- Properties of “good representations”
 - **sparsifies** signals/images of interest
 - can be computed using **fast algorithms**
($O(N)$ or $O(N \log N)$ — think of the FFT)

Sparse Representations

$$f(t) = \sum_i \alpha_i \psi_i(t)$$

- Perfect S -sparsity: only S of the α_i are nonzero
 $\Rightarrow f$ is only S dimensional (in some sense)
- Approximate sparsity (*compressibility*): α obey a power-law:

$$|\alpha|_{(n)} \lesssim n^{-r} \quad r > 1$$

$|\alpha|_{(n)}$ = coefficients sorted by magnitude

- Sparsity/compressibility \Rightarrow accurate low-order approximations

$$\|f - f_S\|_2^2 \lesssim S^{-2r+1}$$

f_S = best S -term approximation of f

- Sparsity $\Rightarrow f$ is **much simpler** in the ψ -domain than in the time/spatial domain

Sparse Representations

$$f(t) = \sum_i \alpha_i \psi_i(t)$$

- Perfect S -sparsity: only S of the α_i are nonzero
 $\Rightarrow f$ is only S dimensional (in some sense)
- Approximate sparsity (*compressibility*): α obey a power-law:

$$|\alpha|_{(n)} \lesssim n^{-r} \quad r > 1$$

$|\alpha|_{(n)}$ = coefficients sorted by magnitude

- Sparsity/compressibility \Rightarrow accurate low-order approximations

$$\|f - f_S\|_2^2 \lesssim S^{-2r+1}$$

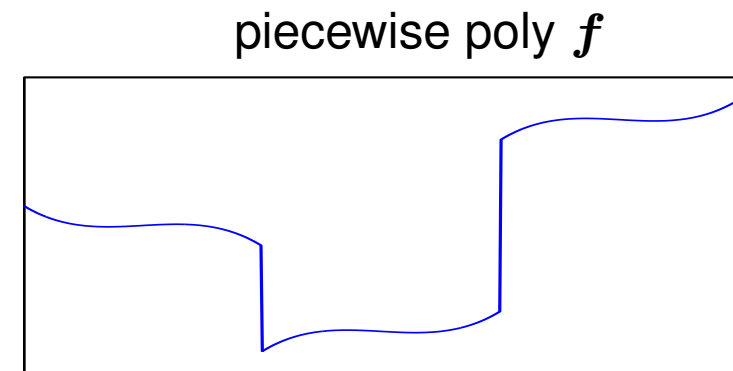
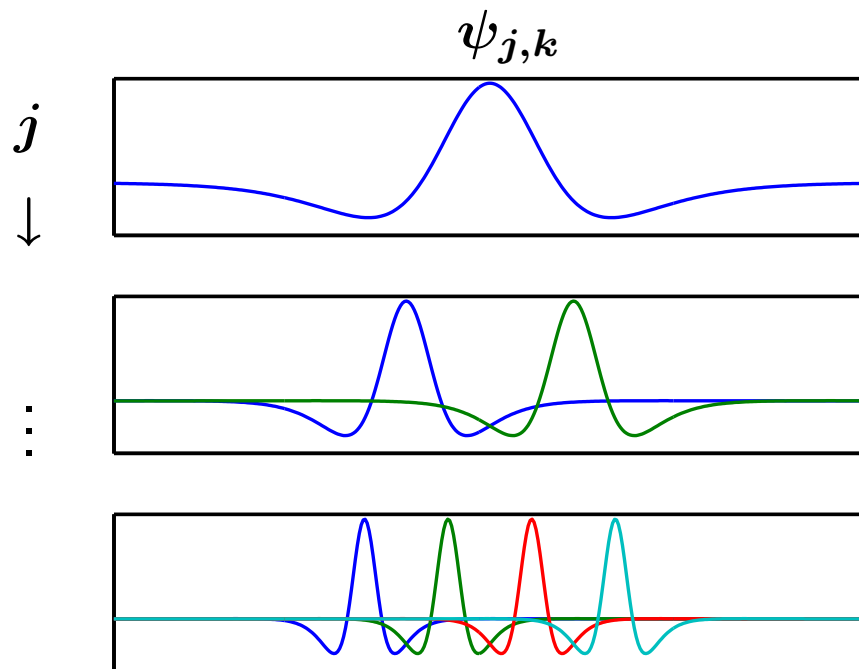
f_S = best S -term approximation of f

- Sparsity $\Rightarrow f$ is **much simpler** in the ψ -domain than in the time/spatial domain

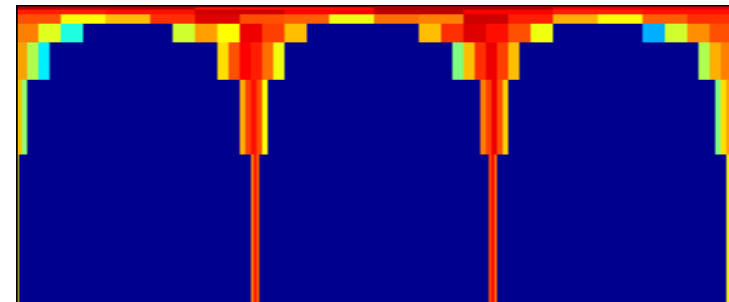
Wavelets

$$f(t) = \sum_{j,k} \alpha_{j,k} \psi_{j,k}(t)$$

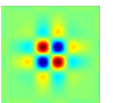
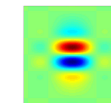
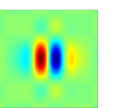
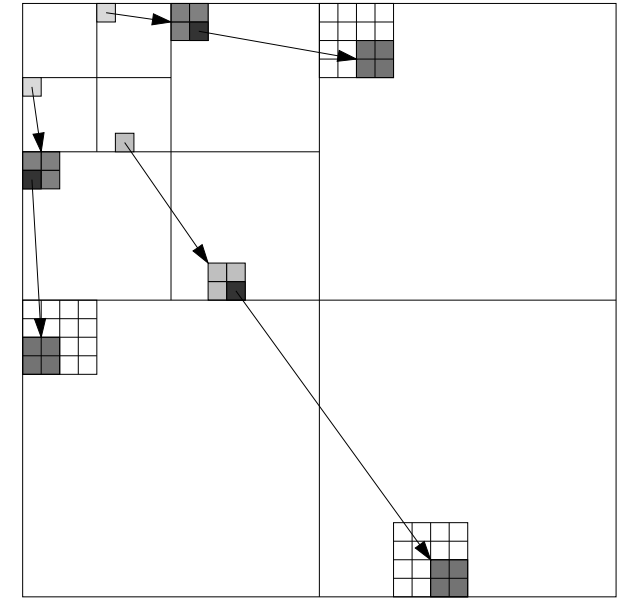
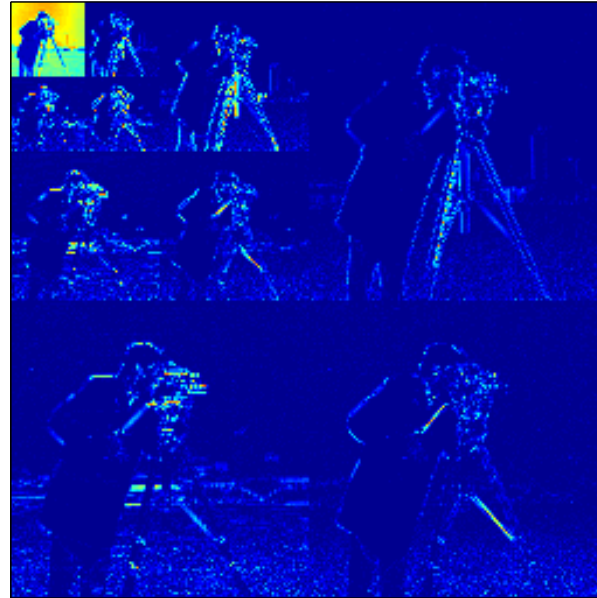
- **Multiscale**: indexed by scale j and location k
- **Local**: $\psi_{j,k}$ analyzes/represents an interval of size $\sim 2^{-j}$
- **Vanishing moments**: in regions where f is polynomial, $\alpha_{j,k} = 0$



wavelet coeffs $\alpha_{j,k}$



2D wavelet transform



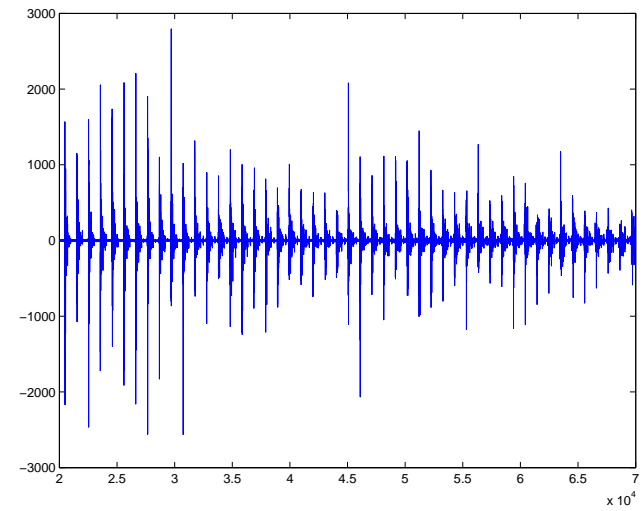
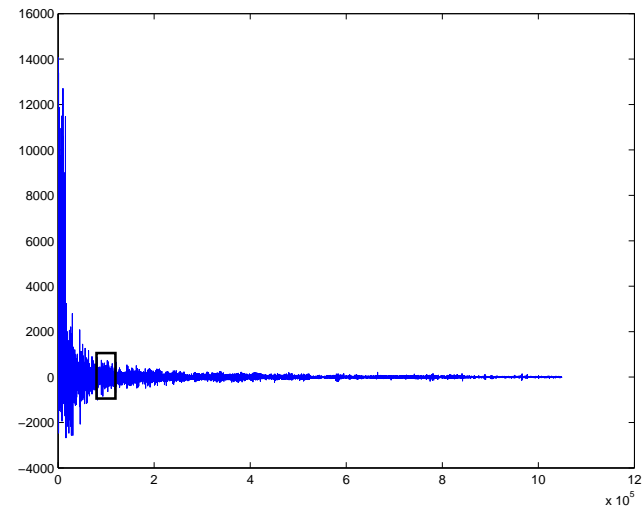
- Important wavelets cluster along edges

Wavelets and Images



1 megapixel image

wavelet coeffs



zoom in

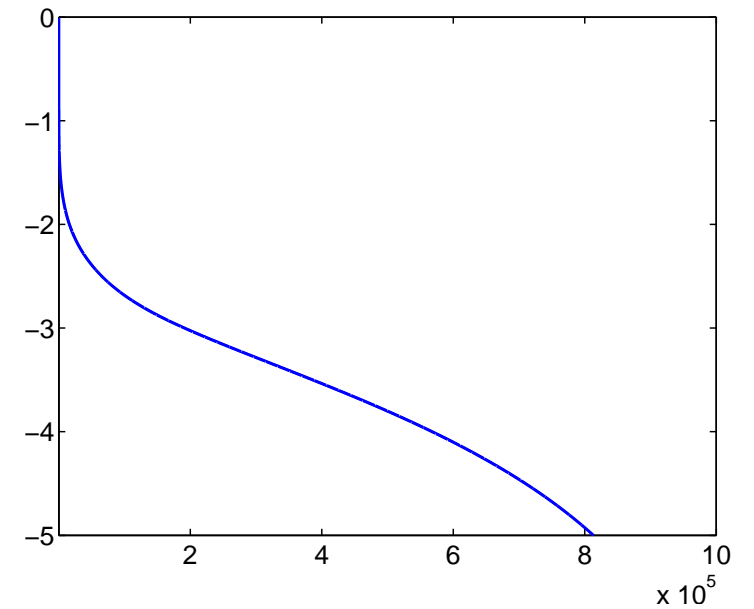
Wavelet Approximation



1 megapixel image



25k term approx



$\log_{10}(S\text{-term approx error})$

The ACHA Paradigm

Sparse representations yield algorithms for (among other things)

1. compression,
2. estimation in the presence of noise (“denoising”),
3. inverse problems (e.g. tomography),
4. acquisition (compressed sensing)

that are

- fast,
- relatively simple,
- and produce (nearly) optimal results

Compression

Transform-Domain Image Coding

- Sparse representation = good compression
Why? Because there are fewer things to code
- Canonical image coder
 1. Transform image into sparse basis
 2. Quantize
Most of the xform coefficients are ≈ 0
 \Rightarrow they require very few bits to encode
 3. Decoder: simply apply inverse transform to quantized coeffs

Image Compression

- Classical example: JPEG (1980s)
 - standard implemented on every digital camera
 - representation = Local Fourier discrete cosine transform on each 8×8 block
- Modern example: JPEG2000 (1990s)
 - representation = wavelets
 - Wavelets are much sparser for images with edges
 - about a factor of 2 better than JPEG in practice
 - half the space for the same quality image

Image Compression

- Key distinction:
JPEG forms DCT approximation in a *fixed* way for all images
JPEG2k *adapts* the wavelet approximation to the image

JPEG vs. JPEG2000

Visual comparison at 0.25 bits per pixel (\approx 100:1 compression)

JPEG



JPEG2000



(Images from David Taubman, University of New South Wales)

Sparse Transform Coding is Asymptotically Optimal

Donoho, Cohen, Daubechies, DeVore, Vetterli, and others . . .

- The statement “transform coding in a sparse basis is a smart thing to do” can be made mathematically precise
- Class of images \mathcal{C}
- Representation $\{\psi_i\}$ (orthobasis) such that

$$|\alpha|_{(n)} \lesssim n^{-r}$$

for all $f \in \mathcal{C}$ ($|\alpha|_{(n)}$ is the n th largest transform coefficient)

- Simple transform coding: transform, quantize (throwing most coeffs away)
- $\ell(\epsilon)$ = length of code (# bits) that **guarantees** the error $< \epsilon$ for all $f \in \mathcal{C}$ (worst case)
- To within log factors

$$\ell(\epsilon) \asymp \epsilon^{-1/\gamma}, \quad \gamma = r - 1/2$$

- For piecewise smooth signals and $\{\psi_i\} =$ wavelets, no coder can do fundamentally better

Statistical Estimation

Statistical Estimation Setup

$$y(t) = f(t) + \sigma z(t)$$

- y : data
- f : object we wish to recover
- z : stochastic error; assume z_t i.i.d. $N(0, 1)$
- σ : noise level
- The quality of an estimate \tilde{f} is given by its **risk** (expected mean-square-error)

$$\text{MSE}(\tilde{f}, f) = E\|\tilde{f} - f\|_2^2$$

Transform Domain Model

$$y = f + \sigma z$$

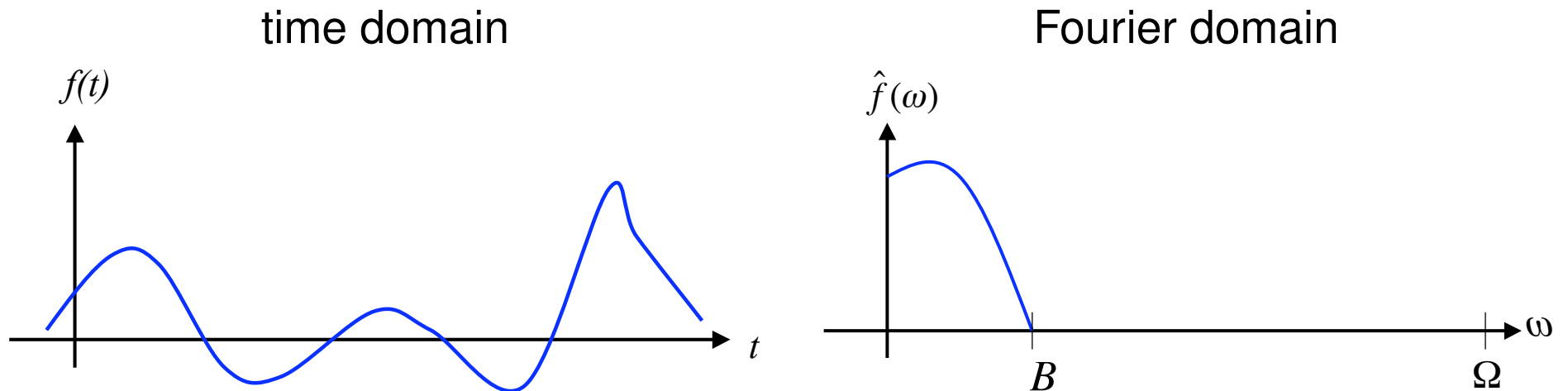
Orthobasis $\{\psi_i\}$:

$$\begin{aligned}\langle y, \psi_i \rangle &= \langle f, \psi_i \rangle + \langle z, \psi_i \rangle \\ \tilde{y}_i &= \alpha_i + z_i\end{aligned}$$

- z_i Gaussian white noise sequence
- σ noise level
- $\alpha_i = \langle f, \psi_i \rangle$ coordinates of f

Classical Estimation Example

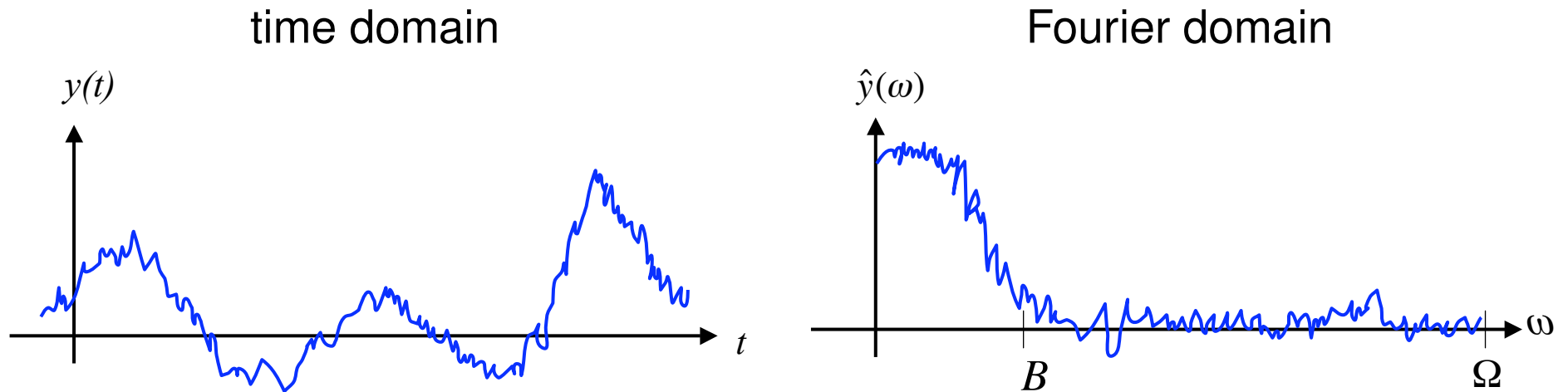
- Classical model: signal of interest f is **lowpass**



- Observable frequencies: $0 \leq \omega \leq \Omega$
- $\hat{f}(\omega)$ is nonzero only for $\omega \leq B$

Classical Estimation Example

- Add noise: $y = f + z$

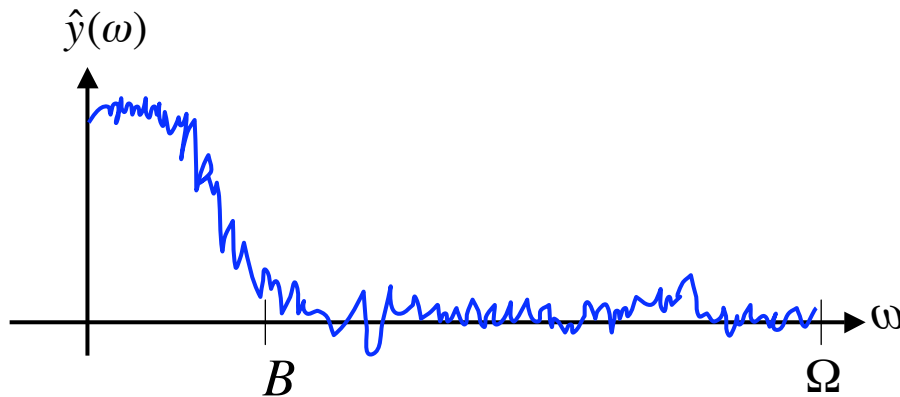


Observation error: $E\|y - f\|_2^2 = E\|\hat{y} - \hat{f}\|_2^2 = \Omega \cdot \sigma^2$

- Noise is **spread out** over entire spectrum

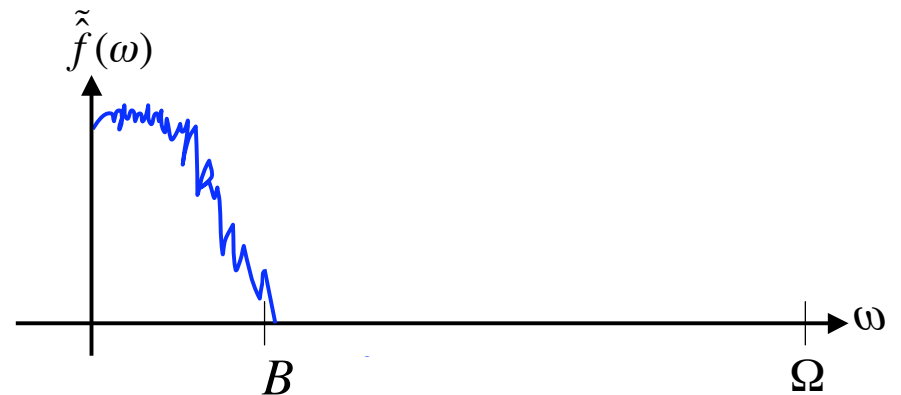
Classical Estimation Example

- Optimal recovery algorithm: lowpass filter (“kill” all $\hat{y}(\omega)$ for $\omega > B$)



Original error

$$E\|\hat{y} - \hat{f}\|_2^2 = \Omega \cdot \sigma^2$$



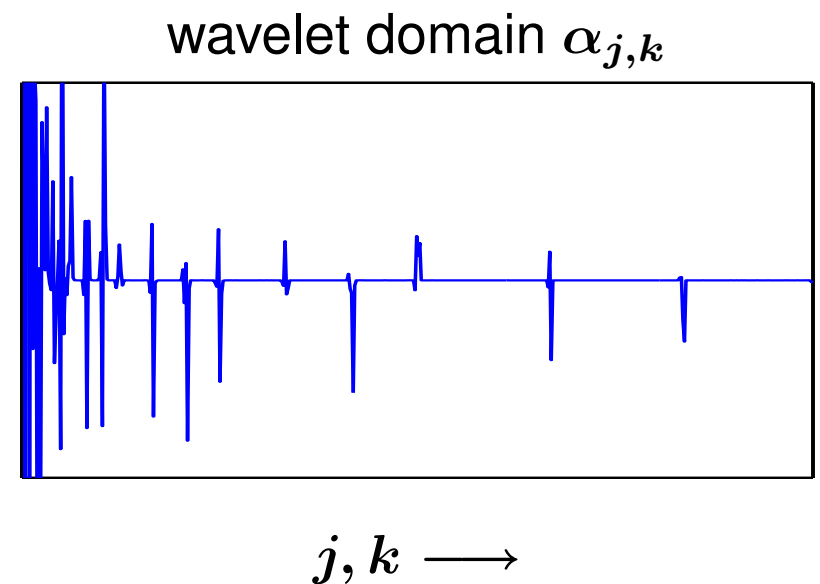
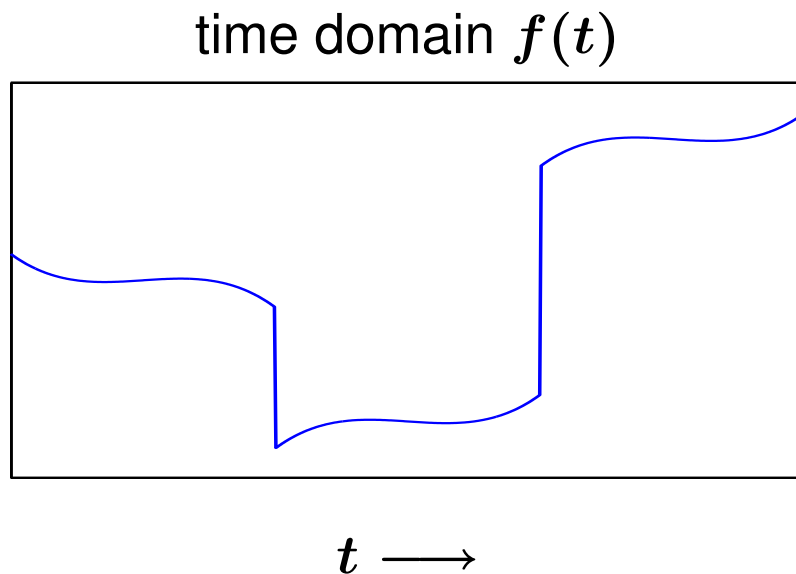
Recovered error

$$E\|\tilde{f} - \hat{f}\|_2^2 = B \cdot \sigma^2$$

- Only the lowpass noise affects the estimate, a savings of $(B/\Omega)^2$

Modern Estimation Example

- Model: signal is **piecewise smooth**
- Signal is sparse in the **wavelet domain**



- Again, the $\alpha_{j,k}$ are concentrated on a small set
- BUT, this set is **signal dependent** (and unknown a priori)
 \Rightarrow we don't know where to "filter"

Ideal Estimation

$$y_i = \alpha_i + \sigma z_i, \quad y \sim \text{Normal}(\alpha, \sigma^2 I)$$

- Suppose an “oracle” tells us which coefficients are above the noise level
- Form the **oracle estimate**

$$\tilde{\alpha}_i^{\text{orc}} = \begin{cases} y_i, & \text{if } |\alpha_i| > \sigma \\ 0, & \text{if } |\alpha_i| \leq \sigma \end{cases}$$

keep the observed coefficients above the noise level, ignore the rest

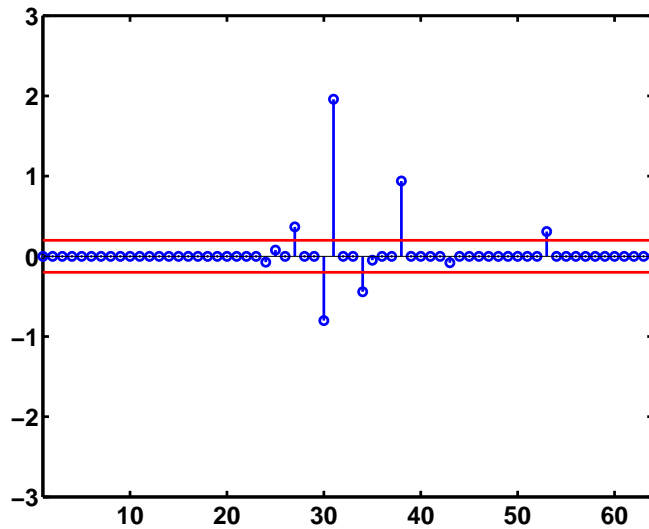
- Oracle Risk:

$$E \|\tilde{\alpha}^{\text{orc}} - \alpha\|_2^2 = \sum_i \min(\alpha_i^2, \sigma^2)$$

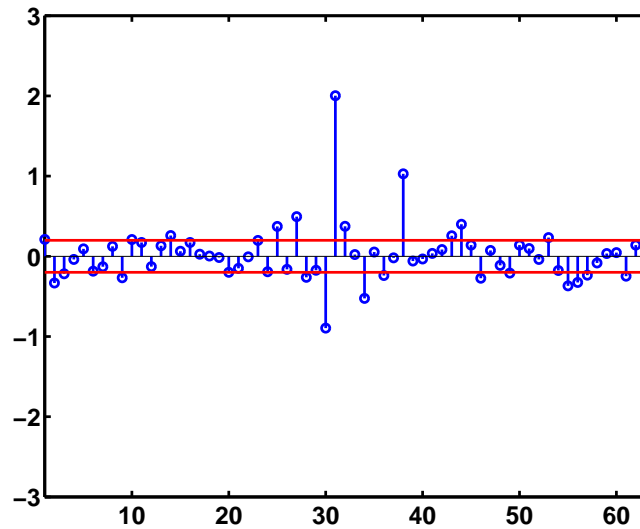
Ideal Estimation

- Transform coefficients α
 - Total length $N = 64$
 - # nonzero components = 10
 - # components above the noise level $S = 6$

original coeffs α

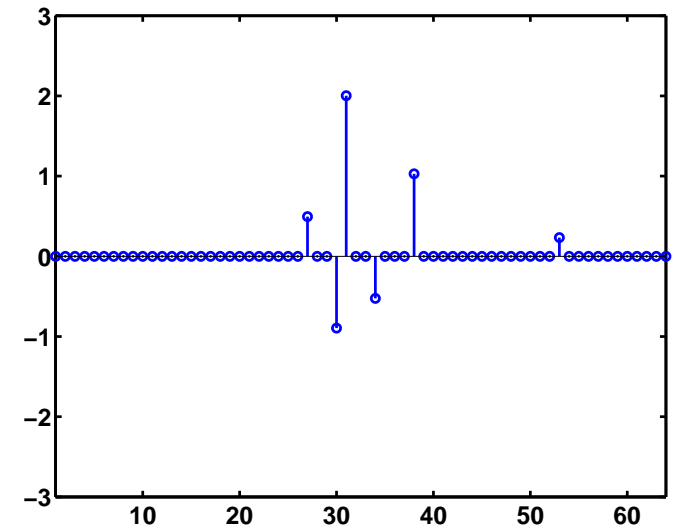


noisy coeffs y



$$E\|y - \alpha\|_2^2 = N \cdot \sigma^2$$

oracle estimate $\tilde{\alpha}_{\text{orc}}$



$$E\|\tilde{\alpha}_{\text{orc}} - \alpha\|_2^2 = S \cdot \sigma^2$$

Interpretation

$$\text{MSE}(\tilde{\alpha}^{\text{orc}}, \alpha) = \sum_i \min(\alpha_i^2, \sigma^2)$$

- Rearrange the coefficients in decreasing order

$$|\alpha|_{(1)}^2 \geq |\alpha|_{(2)}^2 \geq \dots \geq |\alpha|_{(n)}^2$$

- $S(\sigma)$: number of those α_i 's s.t. $\alpha_i^2 \geq \sigma^2$

$$\begin{aligned} \text{MSE}(\tilde{\alpha}^{\text{orc}}, \alpha) &= \sum_{i > N} |\alpha|_{(i)}^2 + S \cdot \sigma^2 \\ &= \|\alpha - \alpha_S\|_2^2 + S \cdot \sigma^2 \\ &= \text{Approx Error} + \text{Number of terms} \times \text{noise level} \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

- The sparser the signal,
 - the better the approximation error (lower bias), and
 - the fewer # terms above the noise level (lower variance)
- *Can we estimate as well without the oracle?*

Denoising by Thresholding

- Hard-thresholding (“keep or kill”)

$$\tilde{\alpha}_i = \begin{cases} y_i, & |y_i| \geq \lambda \\ 0, & |y_i| < \lambda \end{cases}$$

- Soft-thresholding (“shrinkage”)

$$\tilde{\alpha}_i = \begin{cases} y_i - \lambda, & y_i \geq \lambda \\ 0, & -\lambda < y_i < \lambda \\ y_i + \lambda, & y_i \leq -\lambda \end{cases}$$

- Take λ a little bigger than σ
- Working assumption: whatever is above λ is signal, whatever is below is noise

Denoising by Thresholding

- Thresholding performs (almost) as well as the oracle estimator!

- Donoho and Johnstone:

Form estimate $\tilde{\alpha}^t$ using threshold $\lambda = \sigma\sqrt{2\log N}$,

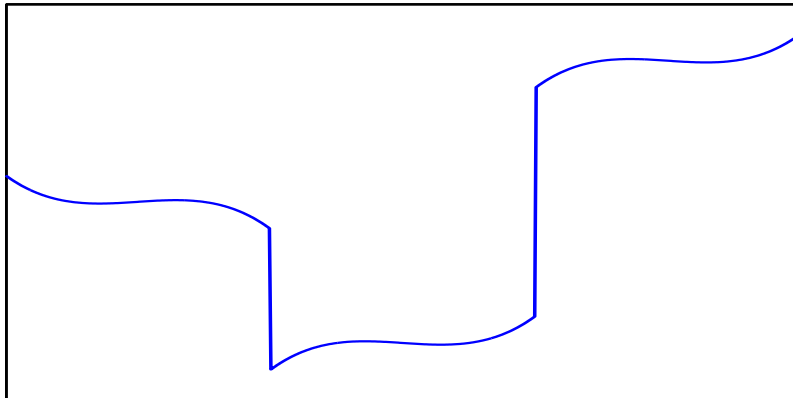
$$\text{MSE}(\tilde{\alpha}^t, \alpha) := E\|\tilde{\alpha}^t - \alpha\|_2^2 \leq (2\log N + 1) \cdot (\sigma^2 + \sum_i \min(\alpha_i^2, \sigma^2))$$

- Thresholding comes within a \log factor of the oracle performance
- The $(2\log N + 1)$ factor is the price we pay for not knowing the locations of the important coeffs
- Thresholding is **simple and effective**
- **Sparsity \Rightarrow good estimation**

Recall: Modern Estimation Example

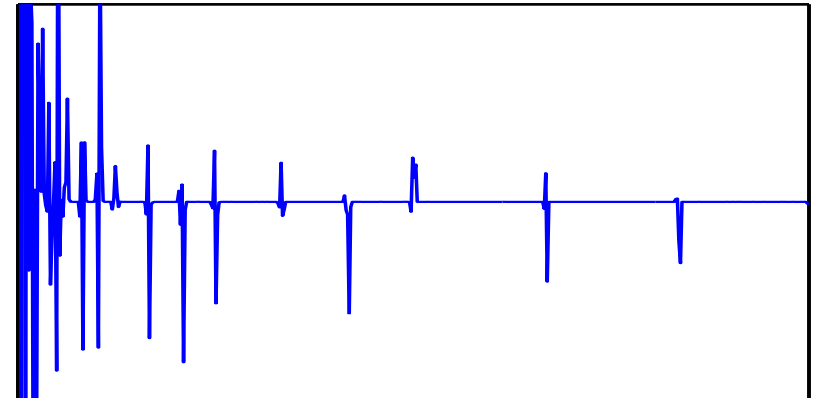
- Signal is **piecewise smooth**, and sparse in the **wavelet domain**

time domain $f(t)$



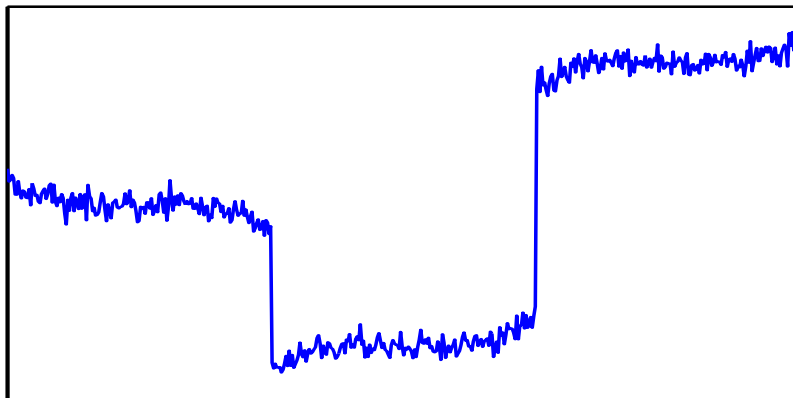
$t \longrightarrow$

wavelet domain $\alpha_{j,k}$



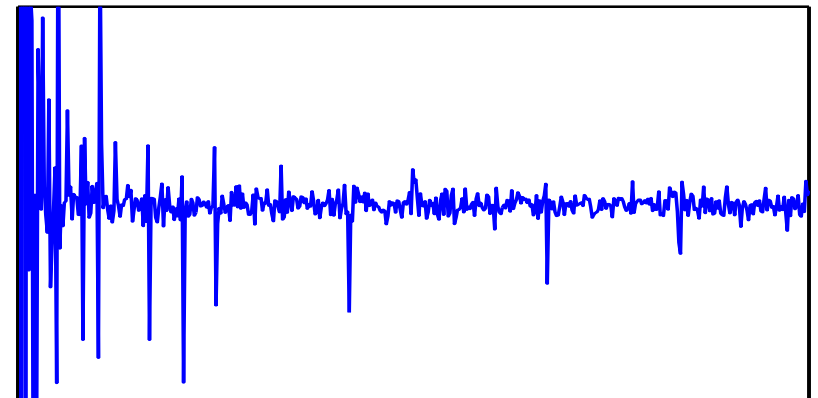
$j, k \longrightarrow$

noisy signal $y(t)$



$t \longrightarrow$

noisy wavelet coeffs

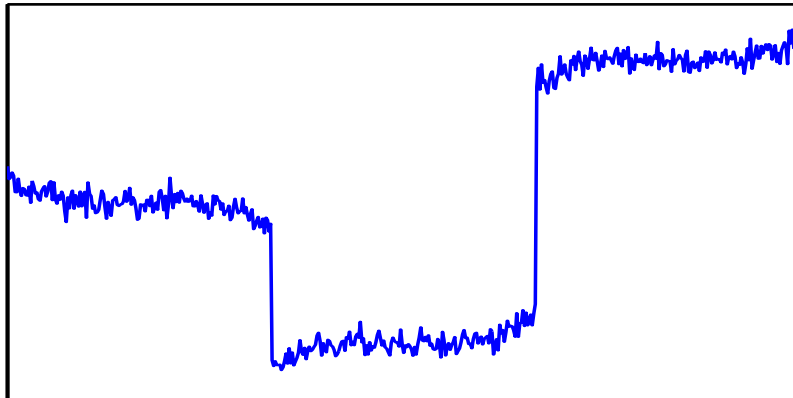


$j, k \longrightarrow$

Thresholding Wavelets

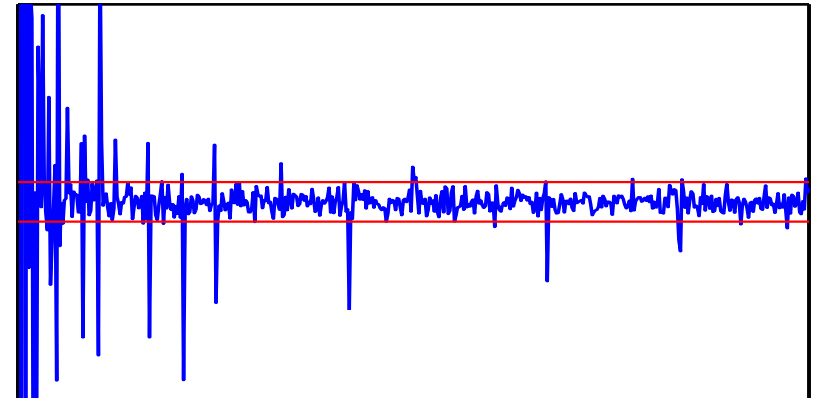
- Denoise (estimate) by soft thresholding

noisy signal



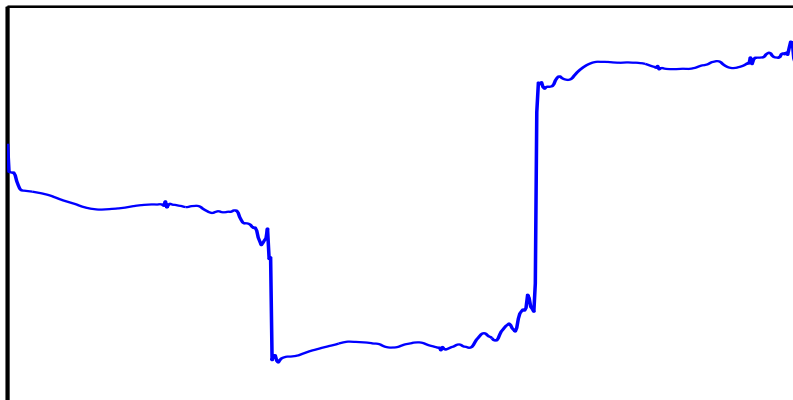
$t \longrightarrow$

noisy wavelet coeffs



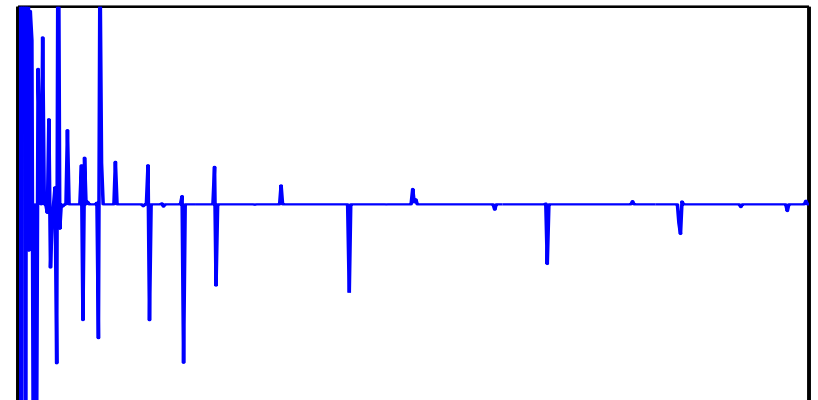
$j, k \longrightarrow$

recovered signal



$t \longrightarrow$

recovered wavelet coeffs



$j, k \longrightarrow$

Denoising the Phantom

noisy



Error = 25.0

lowpass filtered



Error = 42.6

wavelet thresholding, $\lambda = 3\sigma$



Error = 11.0

(Sneak preview: later we will see that we can do even better using *curvelet* thresholding)

Inverse Problems

Linear Inverse Problems

$$y(u) = (Kf)(u) + z(u), \quad u = \text{measurement variable/index}$$

- $f(t)$ object of interest
- K linear operator, indirect measurements

$$(Kf)(u) = \int k(u, t) f(t) dt$$

Examples:

- Convolution (“blurring”)
 - Radon (Tomography)
 - Abel
- $z = \text{noise}$
 - **Ill-posed**: $f = K^{-1}y$ not well defined

Solving Inverse Problems using the SVD

$$K = U\Lambda V^T$$

$$U = \text{col}(u_1, \dots, u_n), \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n), \quad V = \text{col}(v_1, \dots, v_n)$$

- U = orthobasis for the measurement space,
 V = orthobasis for the signal space
- Rewrite action of operator in terms of these bases:

$$y(\nu) = (Kf)(\nu) \Leftrightarrow \langle u_\nu, y \rangle = \lambda_\nu \langle v_\nu, f \rangle$$

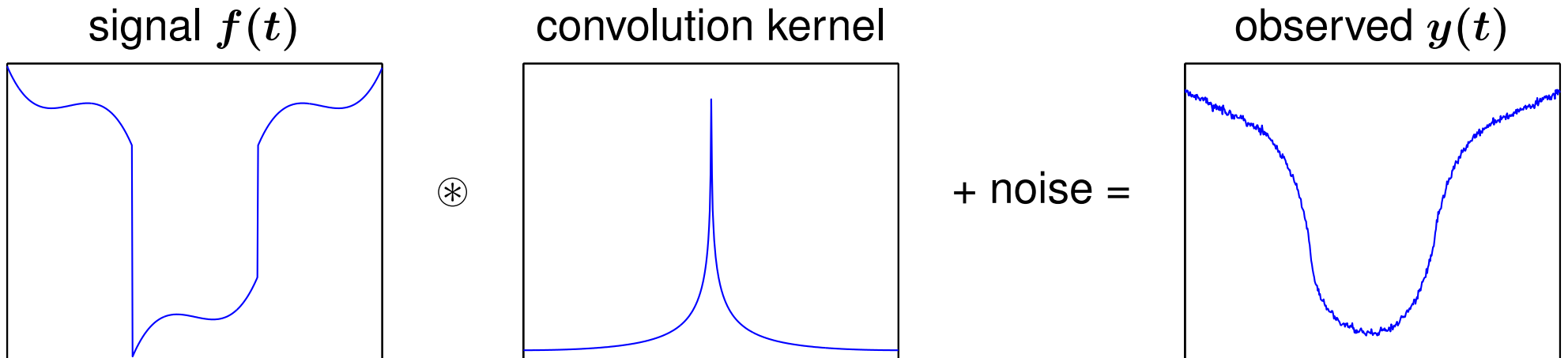
- The inverse operator is also natural:

$$\langle v_\nu, f \rangle = \lambda_\nu^{-1} \langle u_\nu, y \rangle, \quad f = V \begin{pmatrix} \lambda_1^{-1} \langle u_1, y \rangle \\ \lambda_2^{-1} \langle u_2, y \rangle \\ \vdots \end{pmatrix}$$

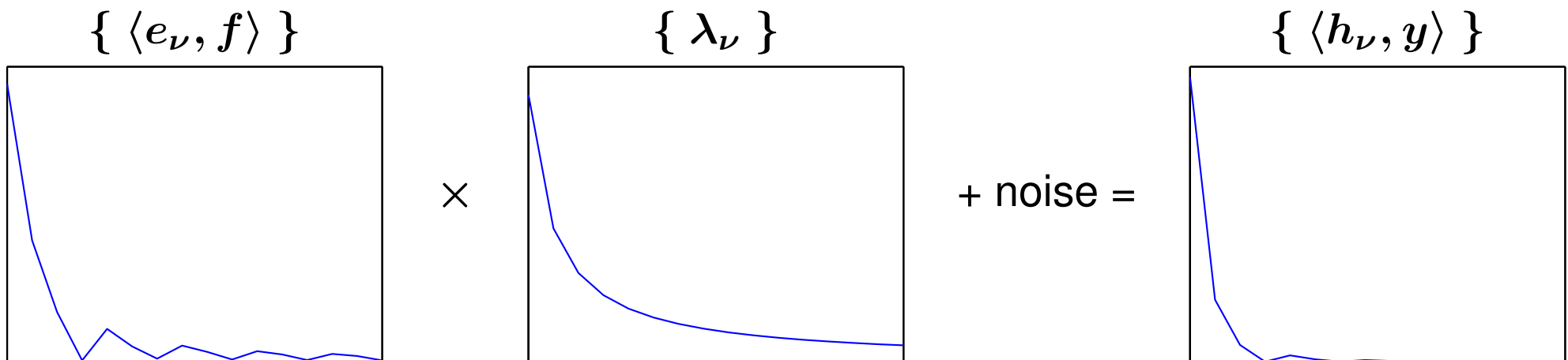
- But in general, $\lambda_\nu \rightarrow 0$, making this unstable

Deconvolution

- Measure $y = Kf + \sigma z$, where K is a convolution operator



- Singular basis: $U = V =$ Fourier transform



Regularization

- Reproducing formula

$$f = \sum_{\nu} \lambda_{\nu}^{-1} \langle u_{\nu}, K f \rangle v_{\nu}$$

- Noisy observations

$$y = K f + \sigma z \quad \Leftrightarrow \quad \langle u_{\nu}, y \rangle = \langle u_{\nu}, K f \rangle + \sigma \hat{z}_{\nu}$$

- Multiply by damping factors w_{ν} to reconstruct from observations y

$$\tilde{f} = \sum_{\nu} w_{\nu} \lambda_{\nu}^{-1} \langle u_{\nu}, y \rangle v_{\nu}$$

want $w_{\nu} \approx 0$ when λ_{ν}^{-1} is large (to keep the noise from exploding)

- If spectral density $\theta_{\nu}^2 = |\langle f, v_{\nu} \rangle|^2$ is known, the MSE optimal weights are

$$w_{\nu} = \frac{\theta_{\nu}^2}{\theta_{\nu}^2 + \sigma^2} = \frac{\text{signal power}}{\text{signal power} + \text{noise power}}$$

This is the **Wiener Filter**

Ideal Damping

- In the SVD domain:

$$y_\nu = \theta_\nu + \sigma_\nu z_\nu$$

$$y_\nu = \langle u_\nu, y \rangle, \quad \theta_\nu = \langle f, v_\nu \rangle, \quad \sigma_\nu = \sigma / \lambda_\nu, \quad z_\nu \sim \text{iid Gaussian}$$

- Again, suppose an oracle tells us which of the θ_ν are above the noise level
- Oracle “keep or kill” window (minimizes MSE)

$$w_\nu = \begin{cases} 1 & |\theta_\nu| > \sigma_\nu \\ 0 & \text{otherwise} \end{cases}$$

Take $\tilde{\theta}_\nu = w_\nu y_\nu$ (thresholding)

- Since V is an isometry, oracle risk is

$$E\|f - \tilde{f}\|_2^2 = E\|\theta - \tilde{\theta}\|_2^2 = \sum_\nu \min(\theta_\nu^2, \sigma_\nu^2)$$

Interpretation

$$\begin{aligned}MSE &= \sum_{\nu} \min(\theta_{\nu}^2, \sigma_{\nu}^2) \\&= \sum_{\nu: |\theta_{\nu}| \lambda_{\nu} \leq \sigma} \theta_{\nu}^2 + \sum_{\nu: |\theta_{\nu}| \lambda_{\nu} > \sigma} \frac{\sigma^2}{\lambda^2} \\&= \text{Bias}^2 + \text{Variance}\end{aligned}$$

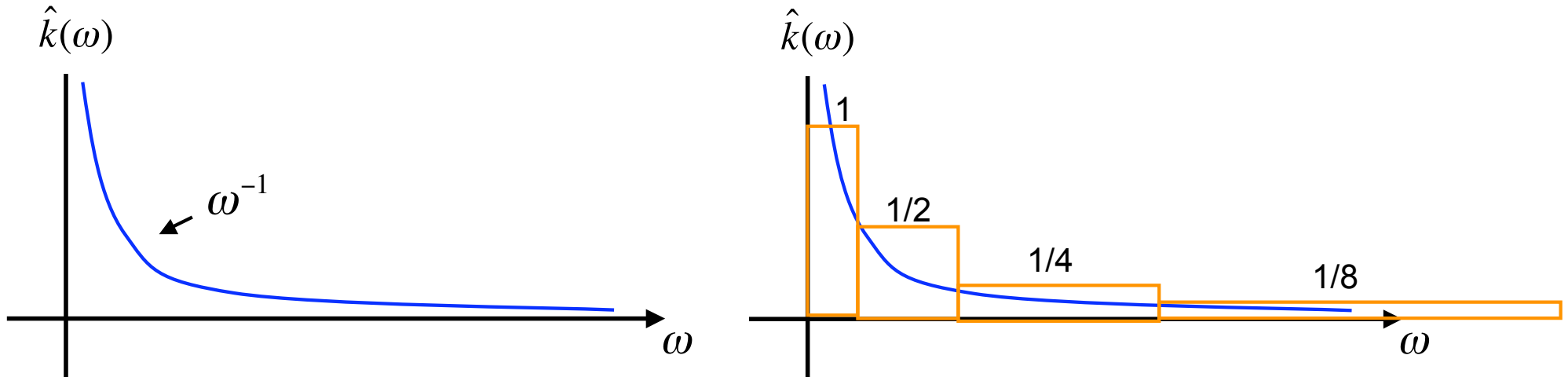
- Again, concentration of the $\theta_{\nu} := \langle f, v_{\nu} \rangle$ on a small set is critical for good performance
- But the v_{ν} are determined only by the operator K !

Typical Situation

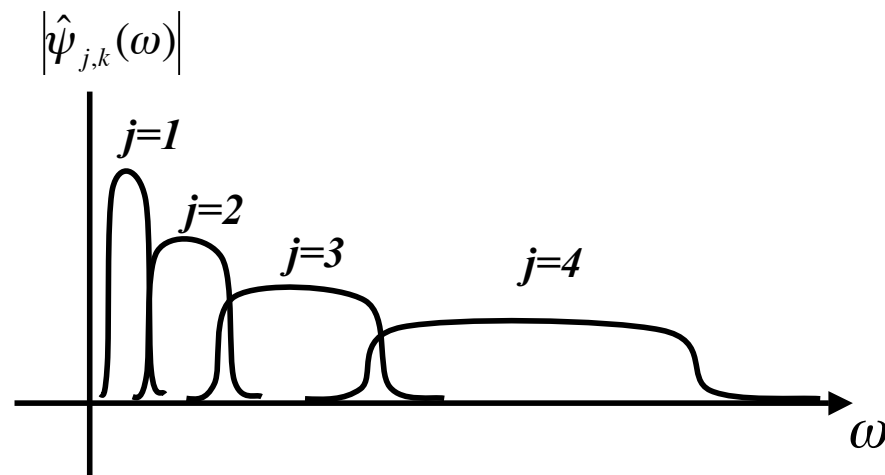
- Convolutions, Radon inversion (tomography)
- $(v_\nu) \sim$ sinusoids
- f has discontinuities (earth, brain, ...)
- SVD basis is *not* a good representation for our signal
- Fortunately, we can find a representation that is simultaneously
 - almost an SVD
 - A sparse decomposition for object we are interested in

Example: Power-law convolution operators

- K = convolution operator with Fourier spectrum $\sim \omega^{-1}$



- Wavelets have dyadic (in scale j) support in Fourier domain



- Spectrum of K is **almost constant** (within a factor of 2) over each subband

The Wavelet-Vaguelette Decomposition (WVD)

Donoho, 1995

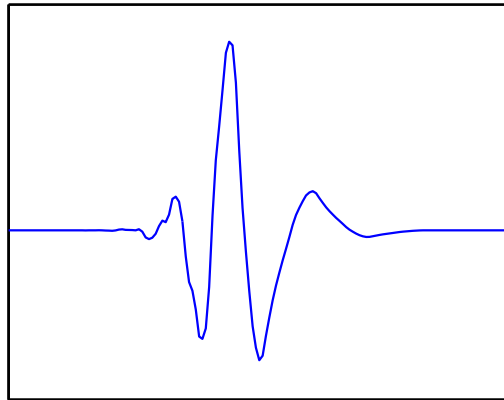
- Wavelet basis $\{\psi_{j,k}\}$ sparsifies piecewise smooth signals
- Vaguelette dual basis $u_{j,k}$ satisfies

$$\langle f, \psi_{j,k} \rangle = 2^{j/2} \langle u_{j,k}, K f \rangle$$

(basis for the measurement space)

- For power-law K , vaguelettes \approx orthogonal, and \approx wavelets

wavelet



vaguelette



- Wavelet-Vaguelette decomposition is **almost an SVD** for Fourier power-law operators

Deconvolution using the WVD

- Observe $y = Kf + \sigma z$,
 $K = 1/|\omega|$ power-law operator, $z = \text{iid Gaussian noise}$
- Expand y in vaguelette basis

$$v_{j,k} = \langle u_{j,k}, y \rangle$$

almost orthonormal, so noise in new basis is \approx independent

- Soft-threshold

$$\tilde{v}_{j,k} = \begin{cases} v_{j,k} - \gamma \text{sign}(v_{j,k}) & |v_{j,k}| > \gamma \\ 0 & |v_{j,k}| \leq \gamma \end{cases}$$

for $\gamma_j \sim 2^{j/2}\sigma$

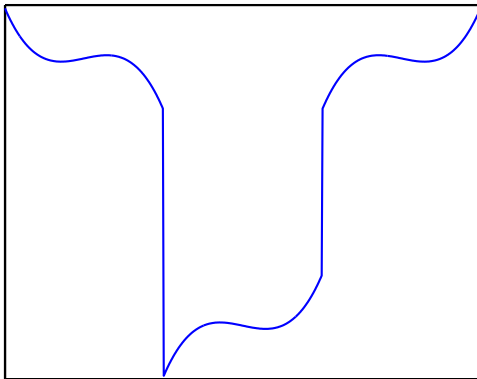
- Weighted reconstruction in the wavelet basis

$$\tilde{f}(t) = \sum_{j,k} 2^{j/2} \tilde{v}_{j,k} \psi_{j,k}(t)$$

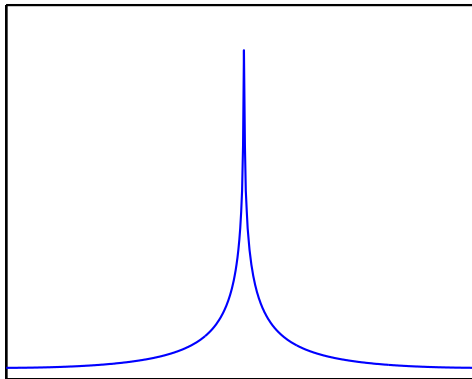
Deconvolution Example

- Measure $y = Kf + \sigma z$, where K is $1/|\omega|$

signal $f(t)$



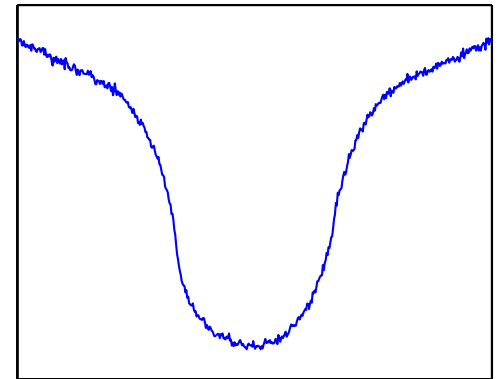
convolution kernel



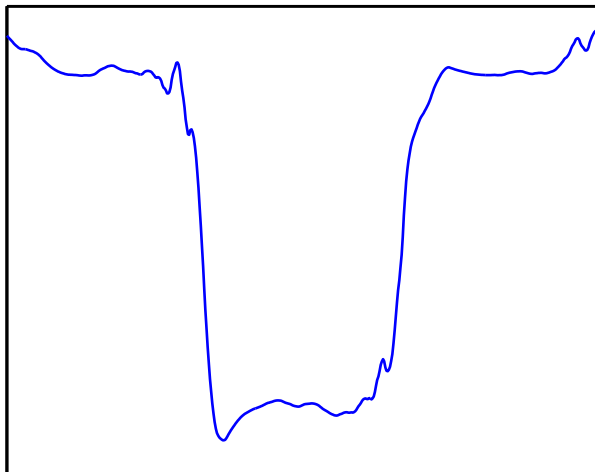
\otimes

+ noise =

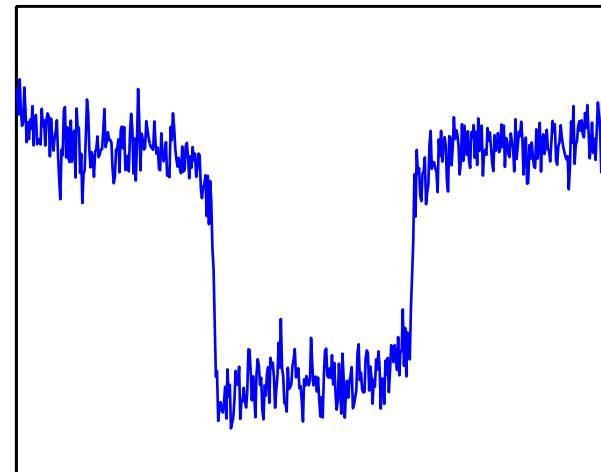
observed $y(t)$



WVD recovery

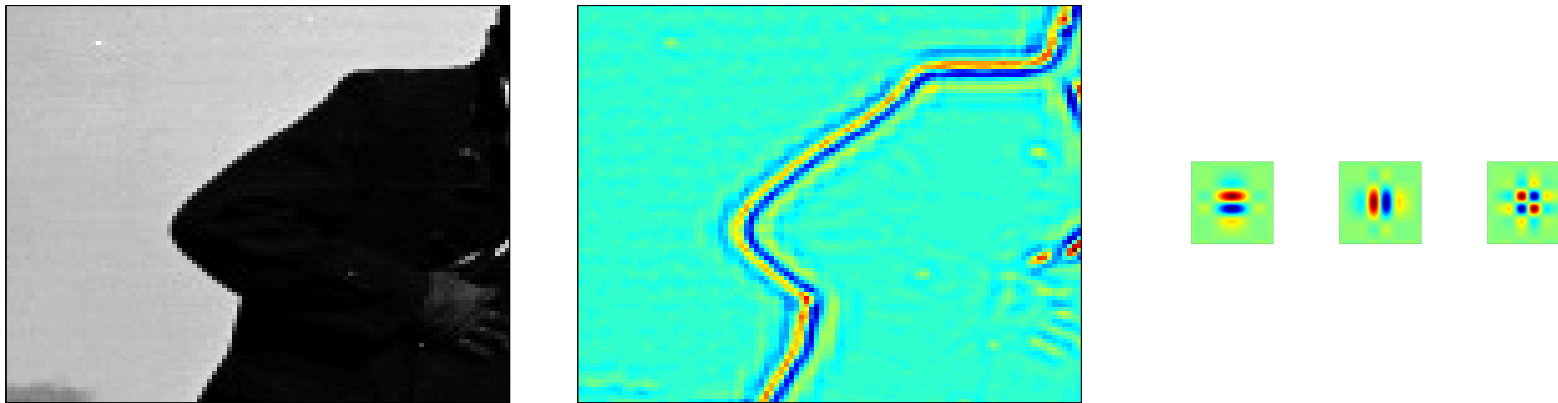


Wiener Filter recovery



Curvelets

Wavelets and Geometry



- Wavelet basis functions are isotropic
⇒ they cannot adapt to *geometrical structure*
- Curvelets offer a more refined scaling concept...

Curvelets

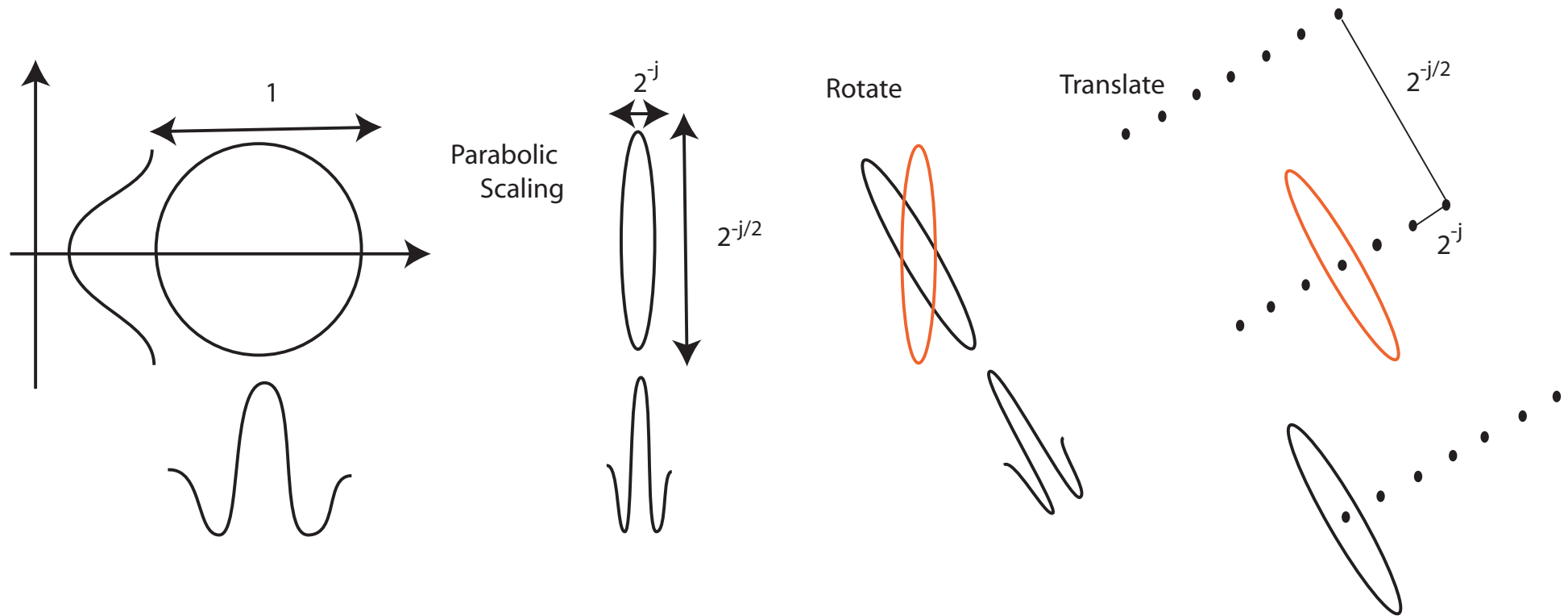
Candes and Donoho, 1999–2004

New multiscale pyramid:

- Multiscale
- Multi-orientations
- *Parabolic (anisotropy) scaling*

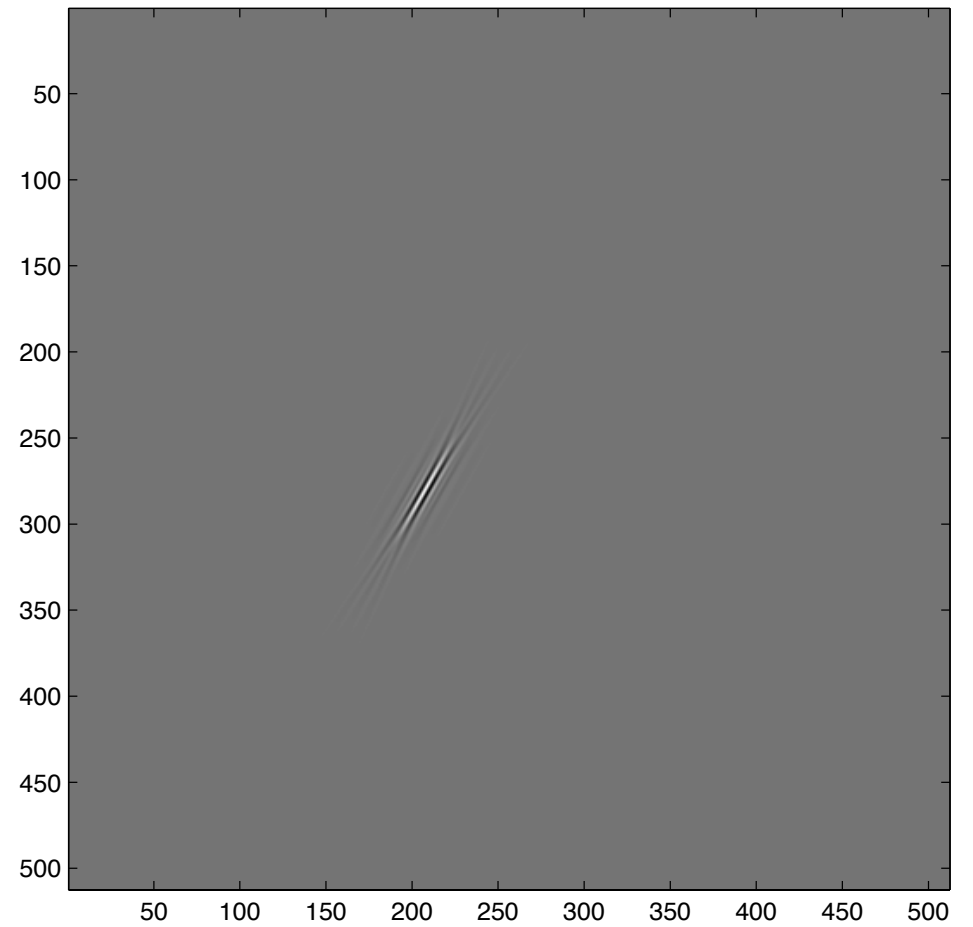
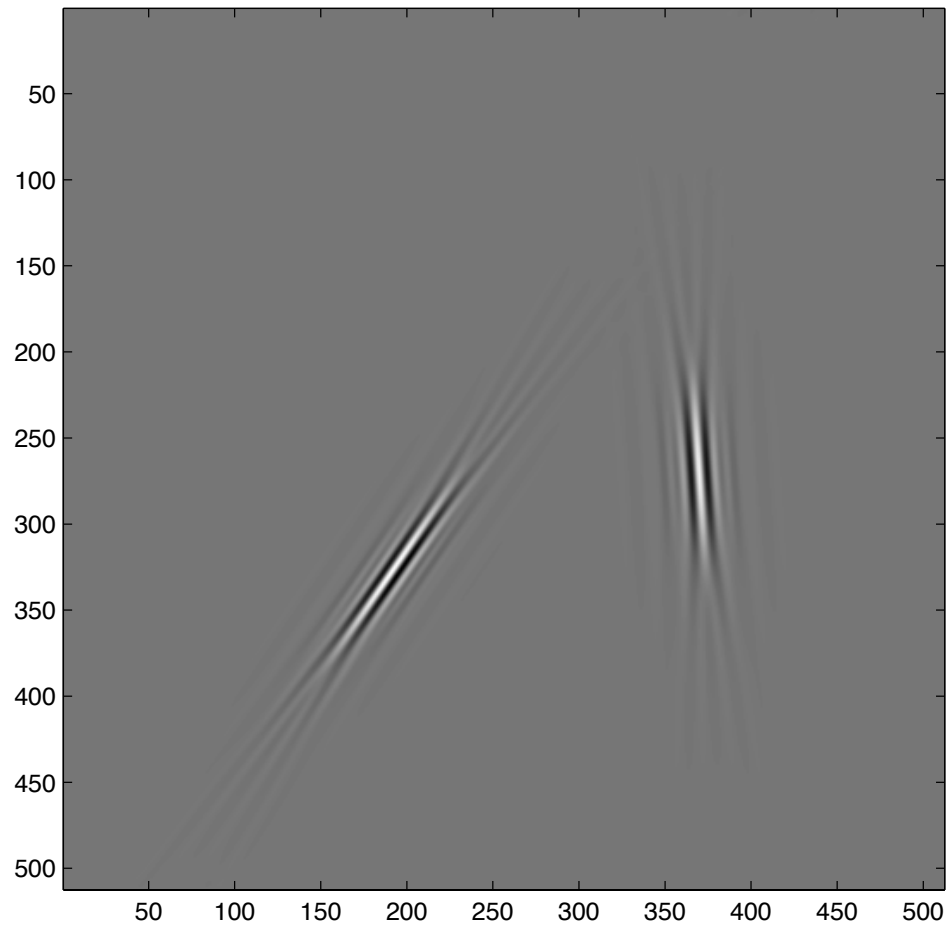
$$width \approx length^2$$

Curvelets in the Spatial Domain



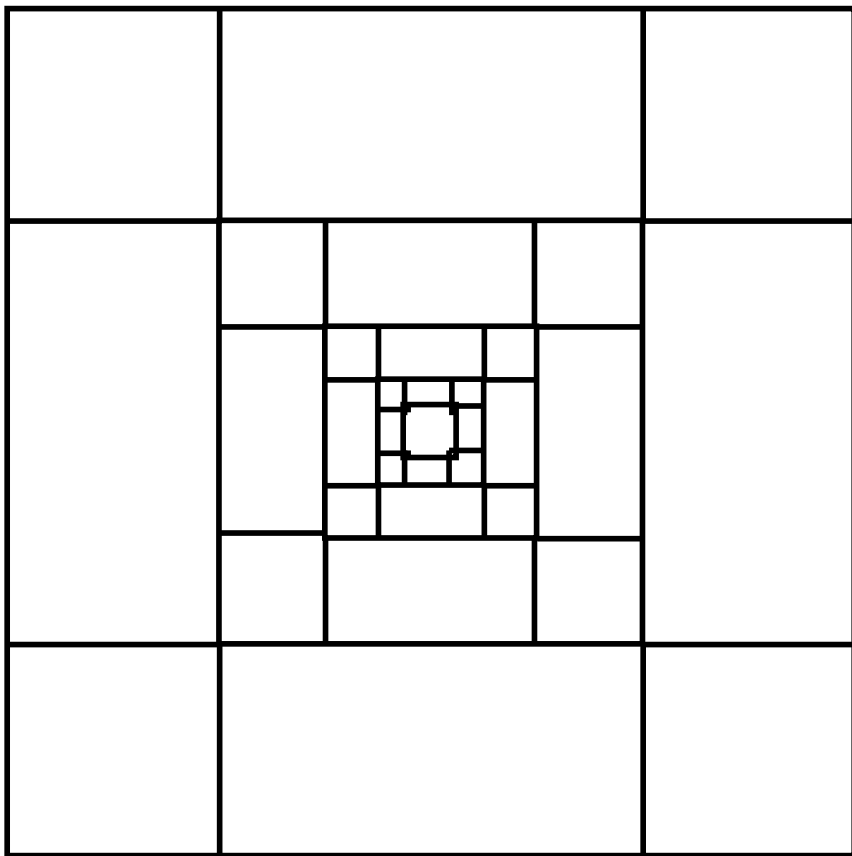
Curvelets parameterized by *scale*, *location*, and *orientation*

Example Curvelets

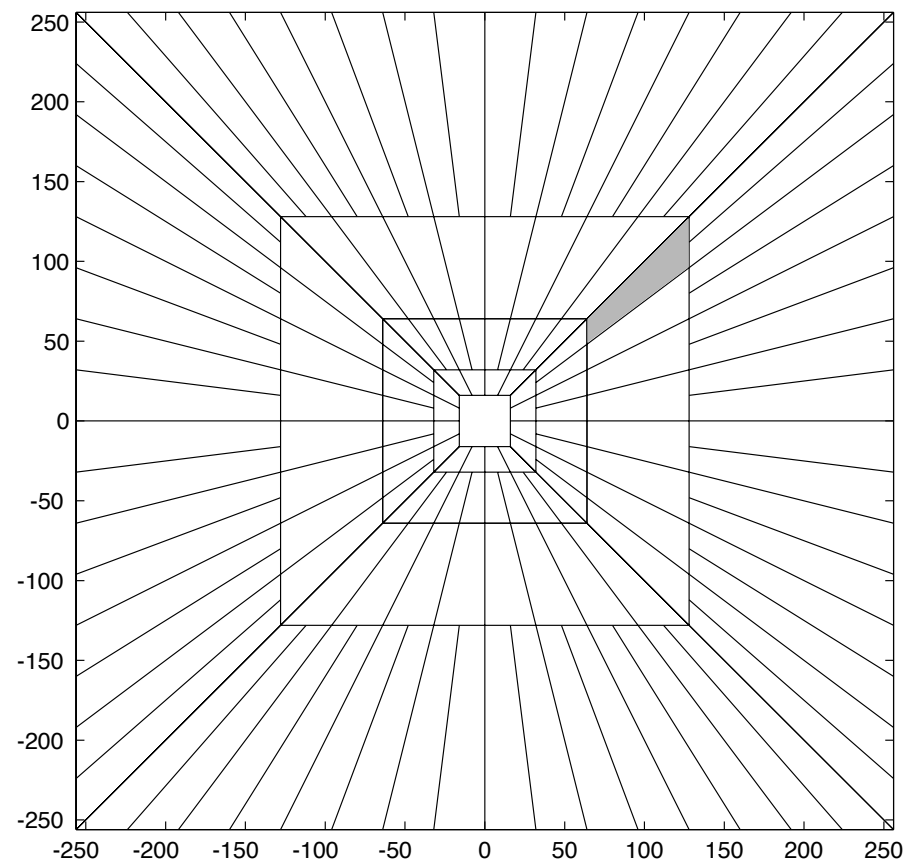


Curvelet Tiling in the Frequency Domain

wavelet



curvelet



Piecewise-smooth Approximation

- Image fragment: C^2 smooth regions separated by C^2 contours
- Fourier approximation

$$\|f - f_S\|_2^2 \lesssim S^{-1/2}$$

- Wavelet approximation

$$\|f - f_S\|_2^2 \lesssim S^{-1}$$

- Curvelet approximation

$$\|f - f_S\|_2^2 \lesssim S^{-2} \log^3 n$$

(within log factor of optimal)

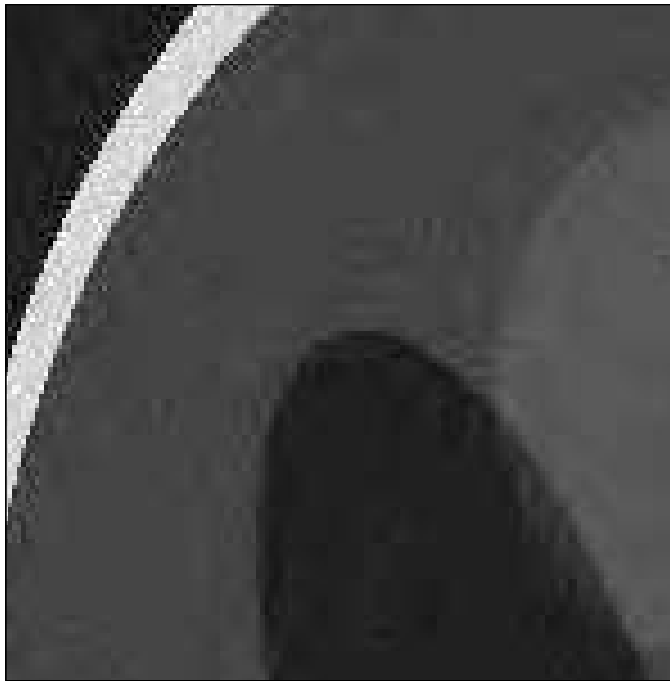
Application: Curvelet Denoising I

Zoom-in on piece of phantom

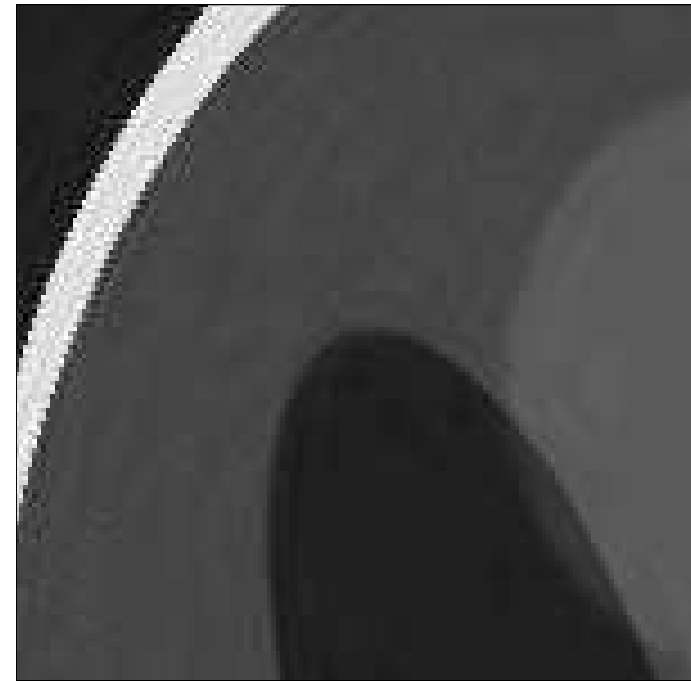
noisy



wavelet thresholding



curvelet thresholding



Application: Curvelet Denoising II

Zoom-in on piece of Lena

wavelet thresholding



curvelet thresholding



Summary (of Part I)

- Having a sparse representation plays a fundamental role in how well we can
 - compress
 - denoise
 - restoreimages
- The above were accomplished with relatively simple algorithms (in practice, we use similar ideas + a bag a tricks)
- Better representation (e.g. curvelets) \longrightarrow better results