

Sparse Representations: An Overview of Theory and Applications

Justin Romberg

Georgia Tech, School of ECE

Tsinghua University

October 14, 2013

Beijing, China

Applied and Computational Harmonic Analysis

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

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

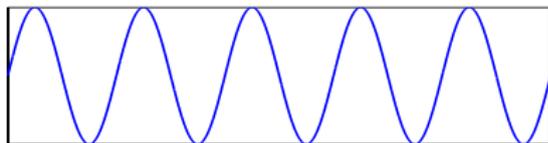
ψ_i = basis functions

α_i = expansion coefficients in ψ -domain

- Classical example: **Fourier series**

ψ_i = complex sinusoids

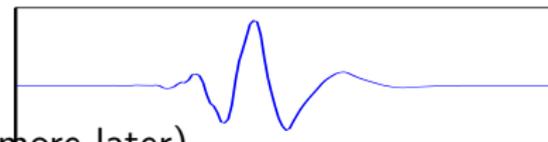
α_i = Fourier coefficients



- Modern example: **wavelets**

ψ_i = "little waves"

α_i = wavelet coefficients



- More exotic example: **curvelets** (more later)

Taking images apart and putting them back together

- Frame operators $\Psi, \tilde{\Psi}$ map images to sequences and back
Two sequences of functions: $\{\psi_i(t)\}, \{\tilde{\psi}_i(t)\}$

Analysis (inner products):

$$\alpha = \tilde{\Psi}^*[f], \quad \alpha_i = \langle \tilde{\psi}_i, f \rangle$$

Synthesis (superposition):

$$f = \Psi[\alpha], \quad f = \sum_i \alpha_i \psi_i(t)$$

- If $\{\psi_i(t)\}$ is an **orthobasis**, then

$$\|\alpha\|_{\ell_2}^2 = \|f\|_{L_2}^2 \quad (\text{Parseval})$$

$$\sum_i \alpha_i \beta_i = \int f(t)g(t) dt \quad (\text{where } \beta = \tilde{\Psi}[g])$$

$$\psi_i(t) = \tilde{\psi}_i(t)$$

i.e. all sizes and angles are preserved

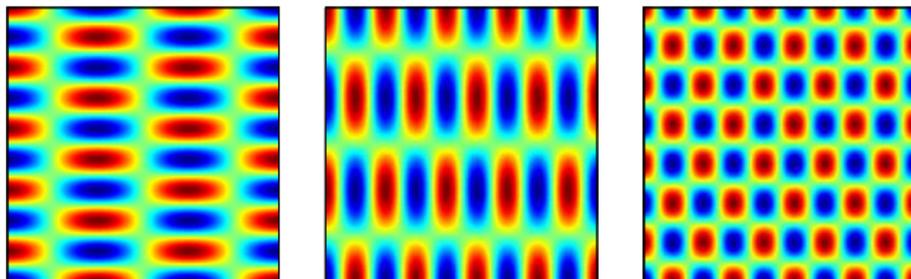
- Overcomplete **tight frames** have similar properties

- 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)
 - ▶ Postmodern: 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)

Example: The discrete cosine transform (DCT)

- For an image $f(t, s)$ on $[0, 1]^2$, we have

$$\psi_{\ell, m}(t, s) = 2\lambda_{\ell}\lambda_m \cdot \cos(\pi\ell t) \cos(\pi m s), \quad \lambda_{\ell} = \begin{cases} 1/\sqrt{2} & \ell = 0 \\ 1 & \text{otherwise} \end{cases}$$



- Closely related to 2D Fourier series/DFT, the DCT is real, and implicitly does symmetric extension
- Can be taken on the whole image, or blockwise (JPEG)

Image approximation using DCT

Take 1% of “low pass” coefficients, set the rest to zero

original



approximated

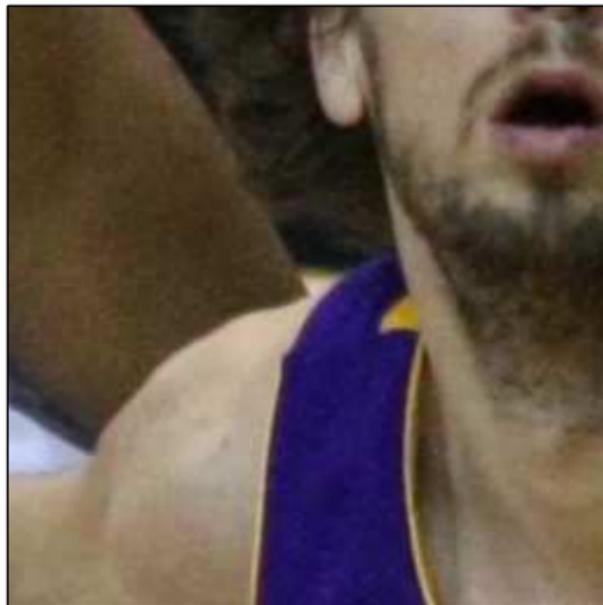


rel. error = 0.075

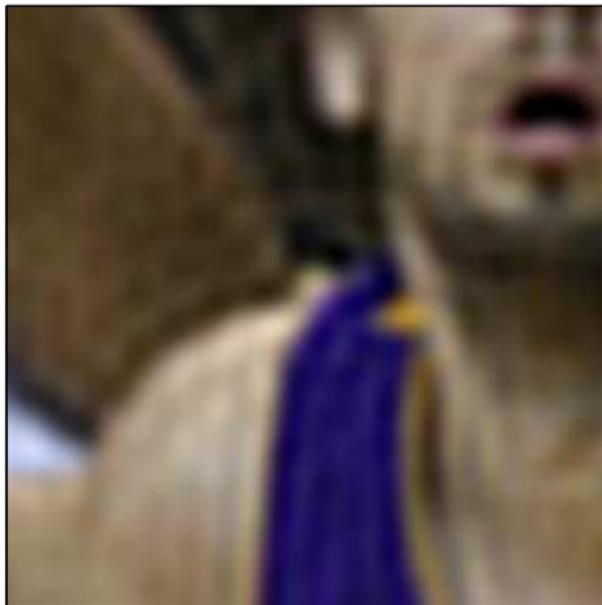
Image approximation using DCT

Take 1% of “low pass” coefficients, set the rest to zero

original



approximated



rel. error = 0.075

Image approximation using DCT

Take 1% of *largest* coefficients, set the rest to zero (adaptive)

original



approximated

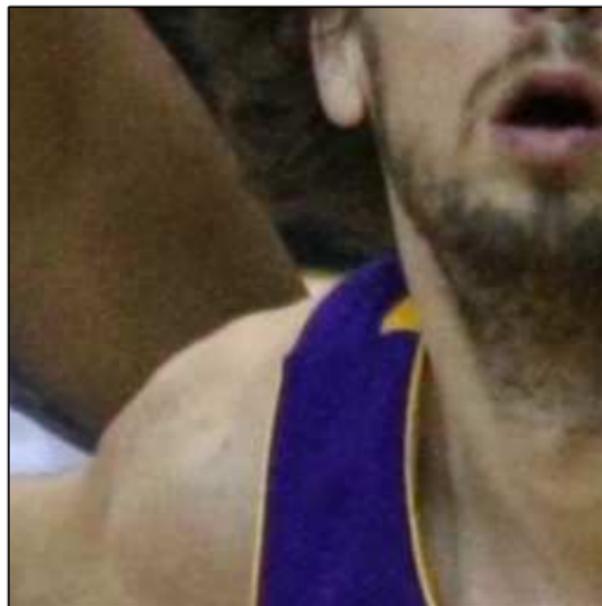


rel. error = 0.057

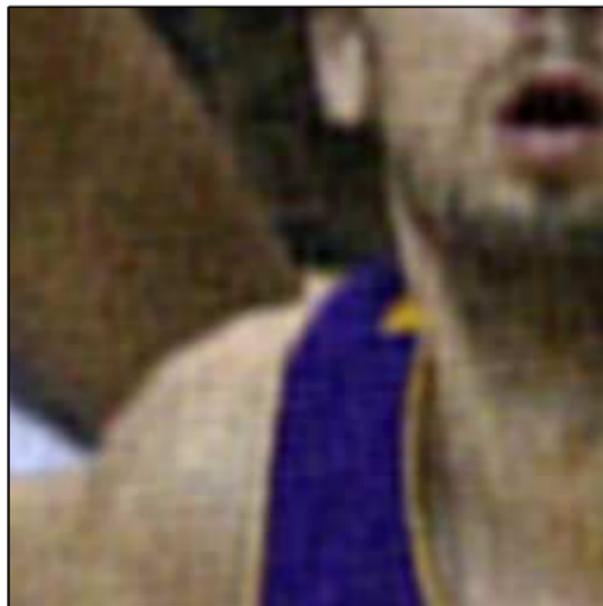
Image approximation using DCT

Take 1% of *largest* coefficients, set the rest to zero (adaptive)

original



approximated

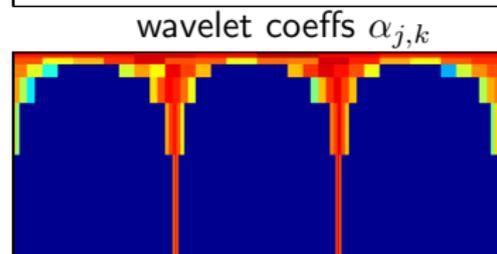
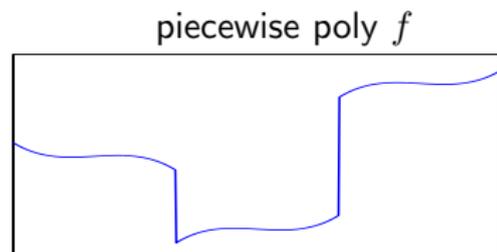
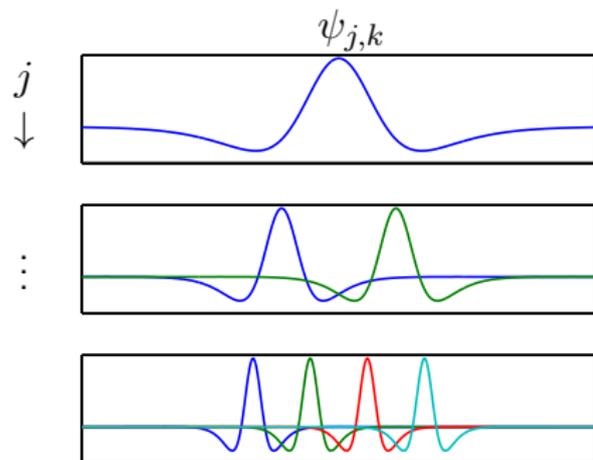


rel. error = 0.057

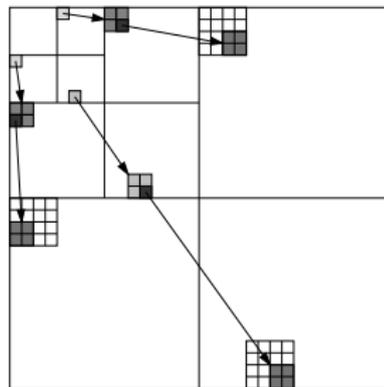
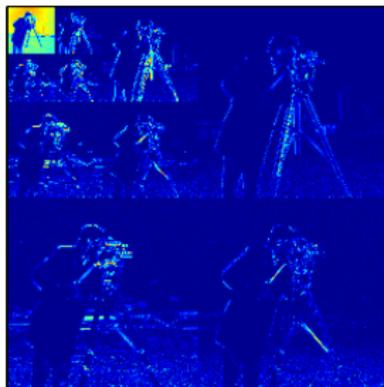
Wavelets

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

- **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, $\alpha_{j,k} = 0$



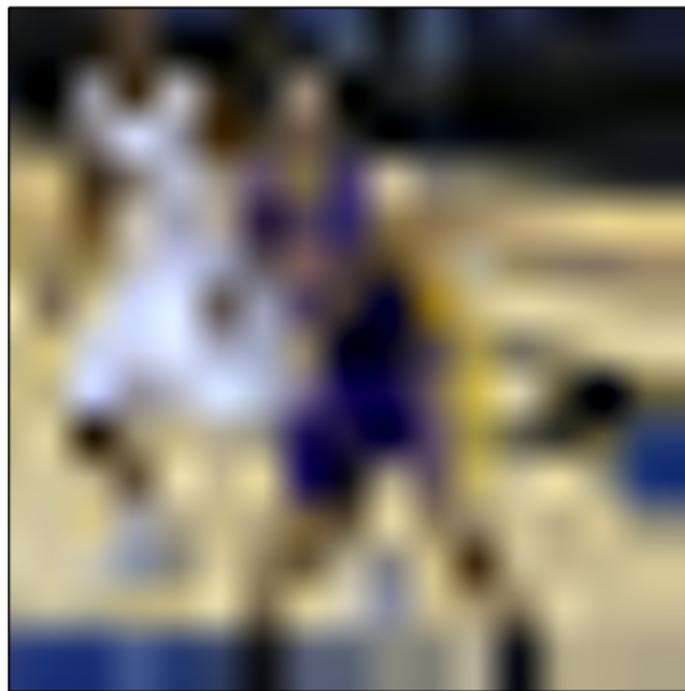
2D wavelet transform



- Sparse: few large coeffs, many small coeffs
- Important wavelets cluster along edges

Multiscale approximations

Scale = 4, 16384:1



rel. error = 0.29

Multiscale approximations

Scale = 5, 4096:1



rel. error = 0.22

Multiscale approximations

Scale = 6, 1024:1



rel. error = 0.16

Multiscale approximations

Scale = 7, 256:1



rel. error = 0.12

Multiscale approximations

Scale = 8, 64:1



rel. error = 0.07

Multiscale approximations

Scale = 9, 16:1



rel. error = 0.04

Multiscale approximations

Scale = 10, 4:1

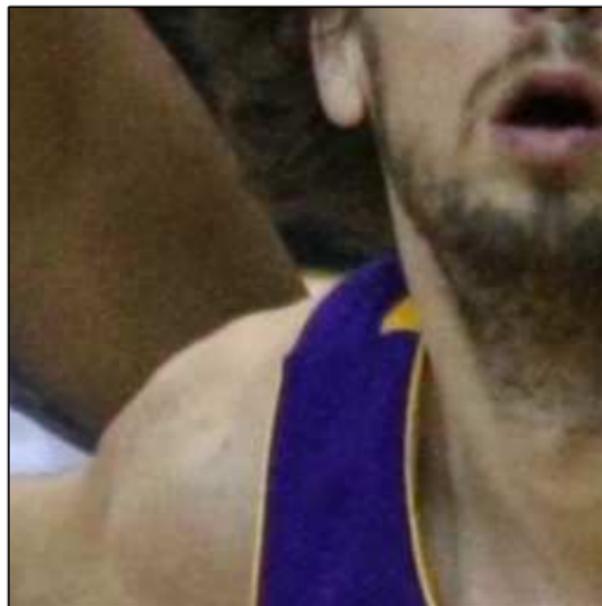


rel. error = 0.02

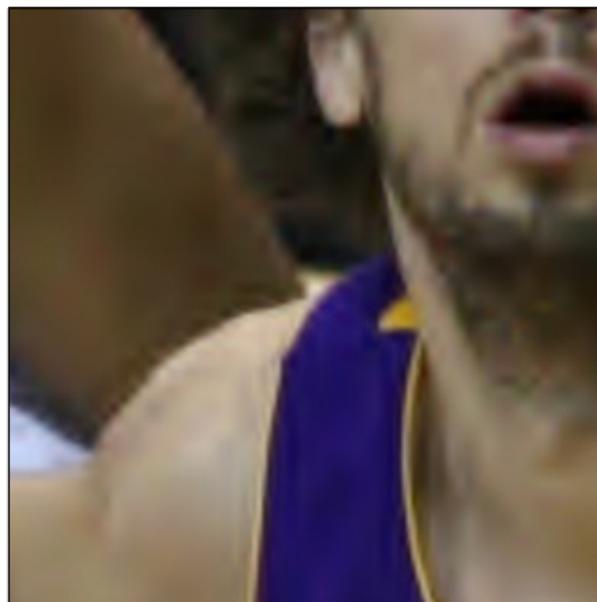
Image approximation using wavelets

Take 1% of *largest* coefficients, set the rest to zero (adaptive)

original



approximated

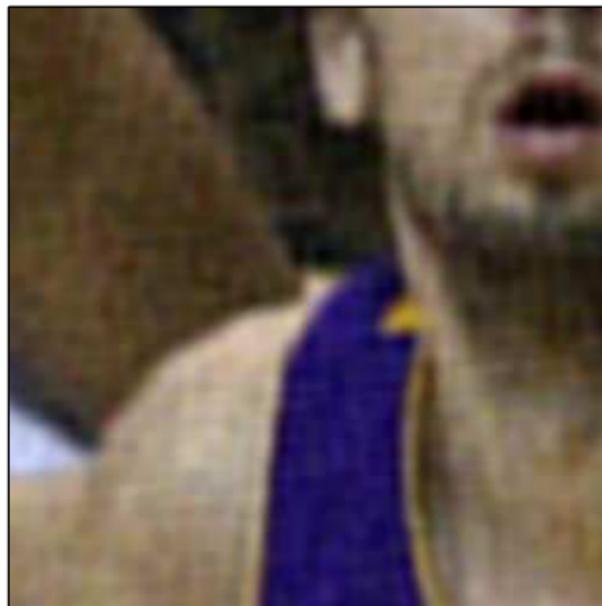


rel. error = 0.031

DCT/wavelets comparison

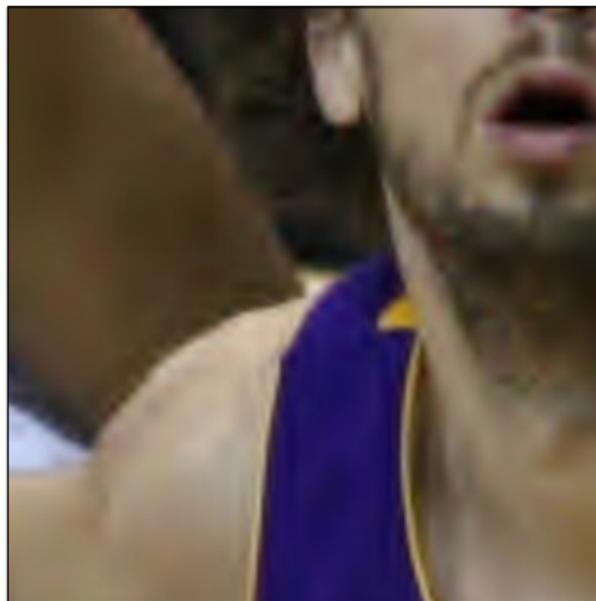
Take 1% of *largest* coefficients, set the rest to zero (adaptive)

DCT



rel. error = 0.057

wavelets



rel. error = 0.031

Linear approximation

- Linear S -term approximation: keep S coefficients in **fixed locations**

$$f_S(t) = \sum_{m=1}^S \alpha_m \psi_m(t)$$

- ▶ projection onto fixed subspace
 - ▶ lowpass filtering, principle components, etc.
- Fast coefficient decay \Rightarrow good approximation

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

- Take $f(t)$ periodic, d -times continuously differentiable,
 $\Psi =$ Fourier series:

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

The smoother the function, the better the approximation

Something similar is true for wavelets ...

Nonlinear approximation

- Nonlinear S -term approximation: keep S *largest* coefficients

$$f_S(t) = \sum_{\gamma \in \Gamma_S} \alpha_\gamma \psi_\gamma(t), \quad \Gamma_S = \text{locations of } S \text{ largest } |\alpha_m|$$

- Fast decay of sorted coefficients \Rightarrow good approximation

$$|\alpha|_{(m)} \lesssim m^{-r} \quad \Rightarrow \quad \|f - f_S\|_2^2 \lesssim S^{-2r+1}$$

$|\alpha|_{(m)}$ = m th largest coefficient

Linear v. nonlinear approximation

- For $f(t)$ *uniformly smooth* with d “derivatives”

S -term approx. error

Fourier, linear	S^{-2d+1}
Fourier, nonlinear	S^{-2d+1}
wavelets, linear	S^{-2d+1}
wavelets, nonlinear	S^{-2d+1}

- For $f(t)$ *piecewise smooth*

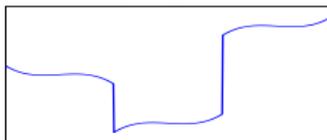
S -term approx. error

Fourier, linear	S^{-1}
Fourier, nonlinear	S^{-1}
wavelets, linear	S^{-1}
wavelets, nonlinear	S^{-2d+1}

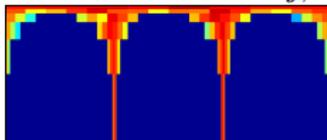
Nonlinear wavelet approximations *adapt* to singularities

Wavelet adaptation

piecewise polynomial $f(t)$

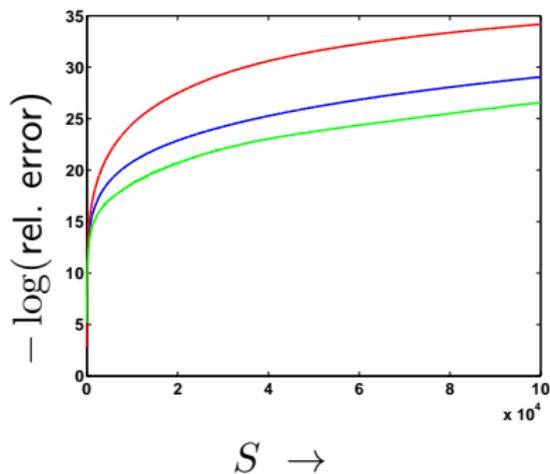


wavelet coeffs $\alpha_{j,k}$



Approximation curves

Approximating Pau with S -terms...



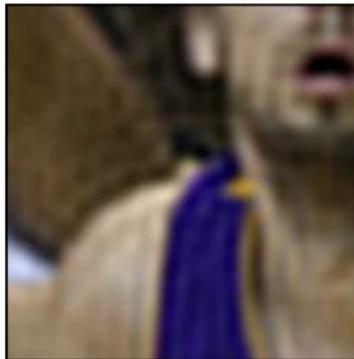
wavelet nonlinear, DCT nonlinear, DCT linear

Approximation comparison

original



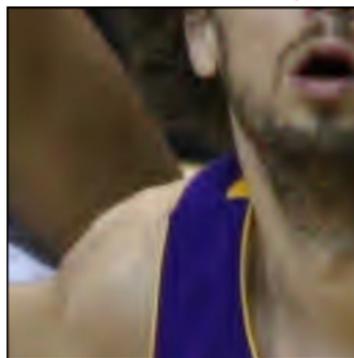
DCT linear (.075)



DCT nonlinear (.057)



wavelet nonlinear (.031)



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
- Basic, “stylized” 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

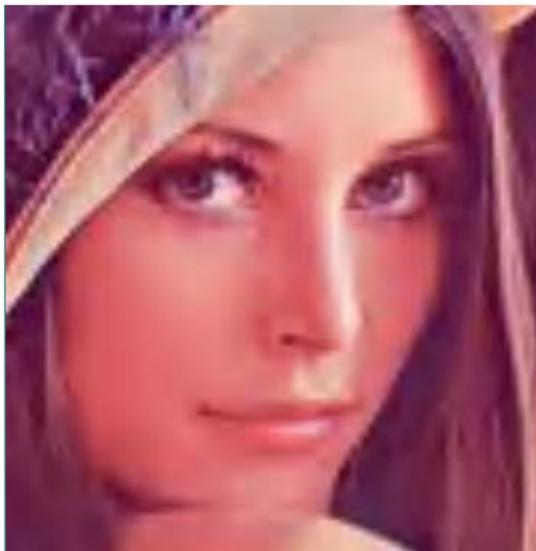
JPEG vs. JPEG2000

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

$$y = f + \sigma z$$

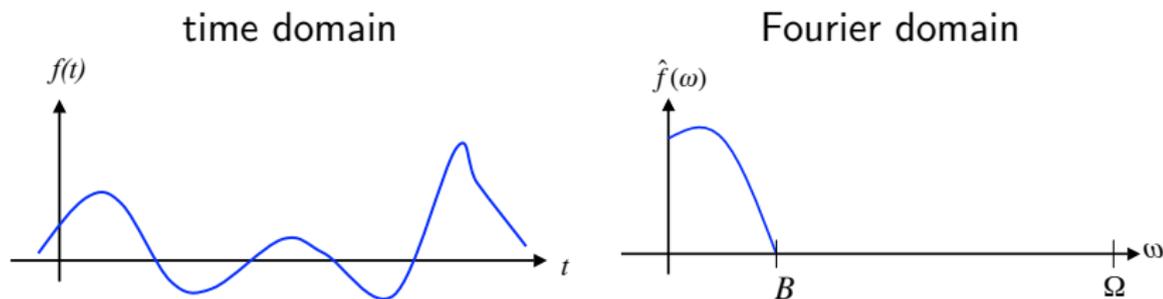
Orthobasis $\{\psi_i\}$:

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

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

Classical estimation example

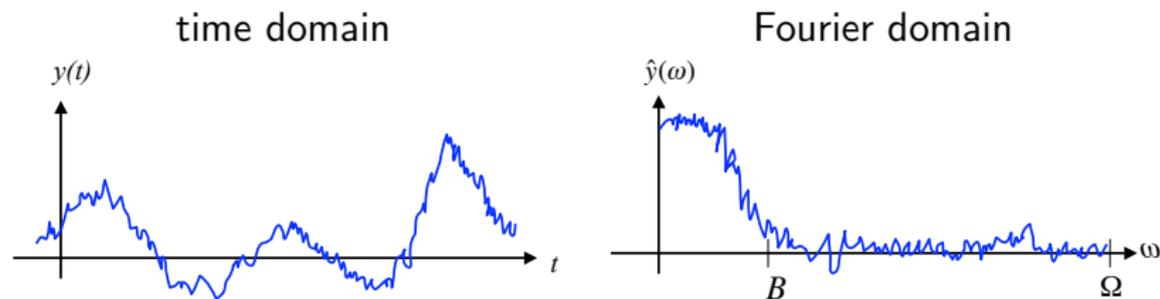
- Classical model: signal of interest f is **lowpass**



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

Classical estimation example

- Add noise: $y = f + z$

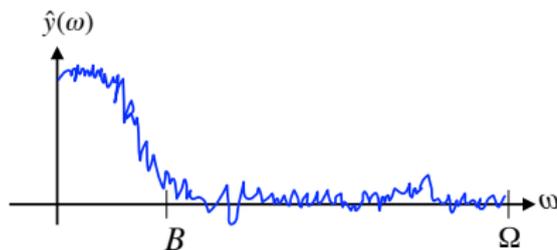


Observation error: $E\|y - f\|_2^2 = E\|\hat{y} - \hat{f}\|_2^2 = \Omega \cdot \sigma^2$

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

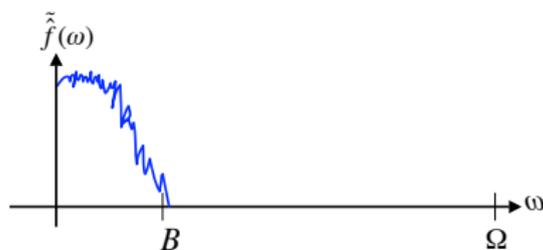
Classical estimation example

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



Original error

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



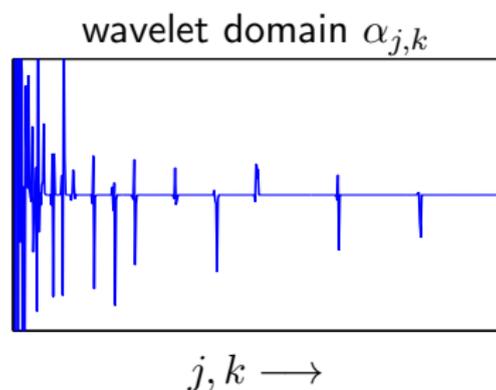
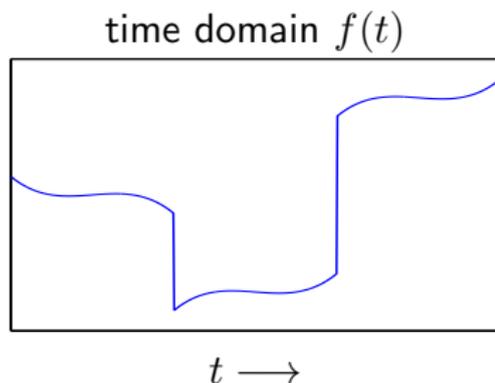
Recovered error

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

- 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
- 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

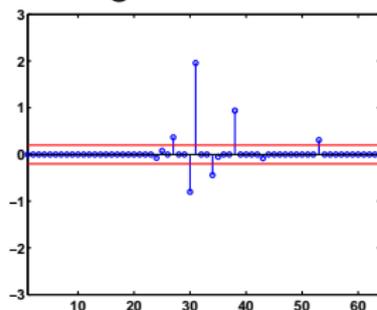
- Oracle Risk:

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

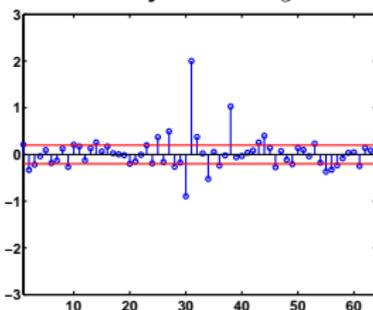
Ideal estimation

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

original coeffs α

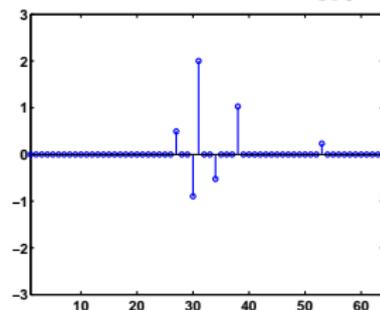


noisy coeffs y



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

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



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

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

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

- 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?*

Denosing 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

Denosing by thresholding

- Thresholding performs (almost) as well as the oracle estimator!
- Donoho and Johnstone:
Form estimate $\tilde{\alpha}^t$ using threshold $\lambda = \sigma\sqrt{2\log N}$,

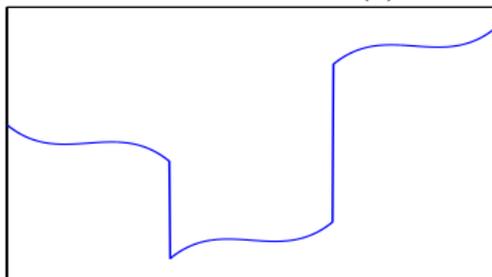
$$\text{MSE}(\tilde{\alpha}^t, \alpha) := E\|\tilde{\alpha}^t - \alpha\|_2^2 \leq (2\log N + 1) \cdot (\sigma^2 + \sum_i \min(\alpha_i^2, \sigma^2))$$

- 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

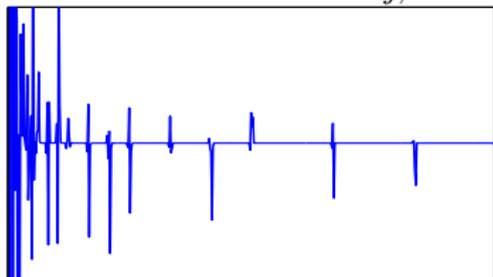
- Signal is **piecewise smooth**, and sparse in the **wavelet domain**

time domain $f(t)$



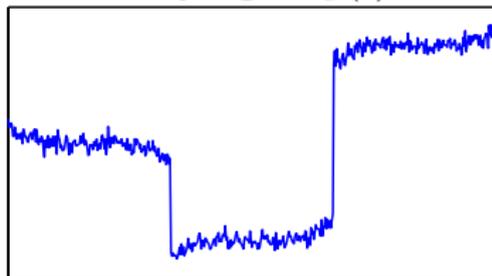
$t \rightarrow$

wavelet domain $\alpha_{j,k}$



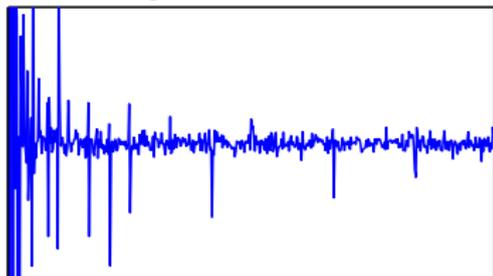
$j, k \rightarrow$

noisy signal $y(t)$



$t \rightarrow$

noisy wavelet coeffs

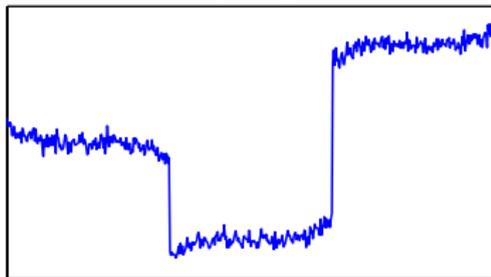


$j, k \rightarrow$

Thresholding wavelets

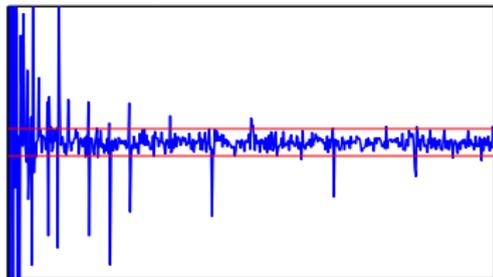
- Denoise (estimate) by soft thresholding

noisy signal



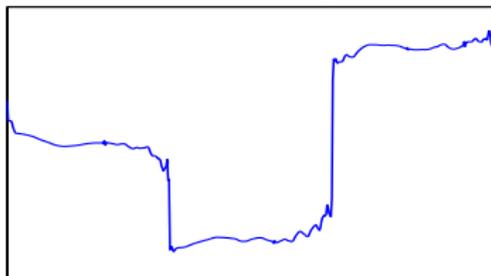
$t \longrightarrow$

noisy wavelet coeffs



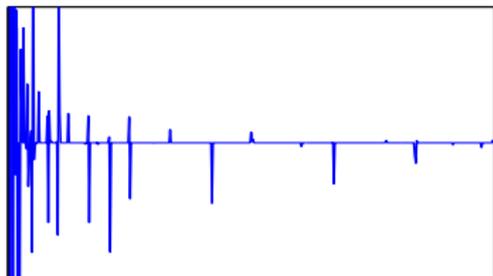
$j, k \longrightarrow$

recovered signal



$t \longrightarrow$

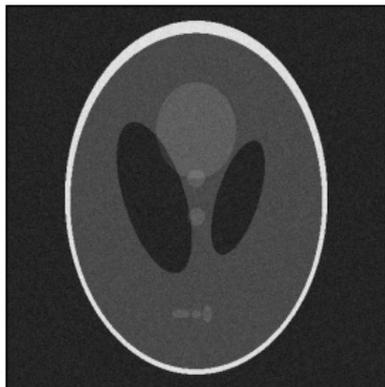
recovered wavelet coeffs



$j, k \longrightarrow$

Denoising the Phantom

noisy



Error = 25.0

lowpass filtered



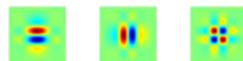
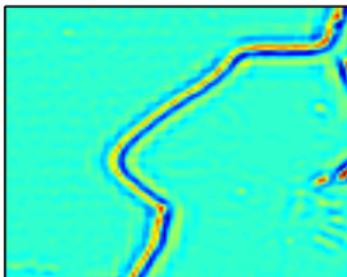
Error = 42.6

wavelet thresholding, $\lambda = 3\sigma$



Error = 11.0

Wavelets and geometry



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

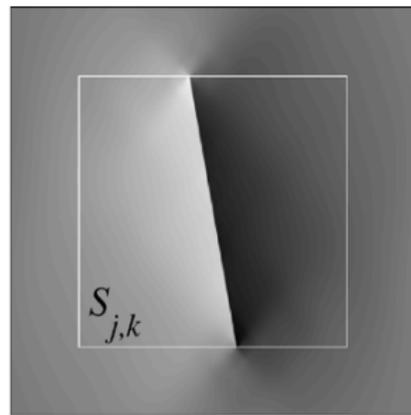
Geometrical transforms



curvelet



bandelet



wedgeprint

These geometrical basis functions are parameterized by *scale*, *location*, and *orientation*

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

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

- Geometrical basis approximation

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

(for some small q ; within log factor of optimal)

Application: Curvelet denoising

Zoom-in on piece of Lena

wavelet thresholding



curvelet thresholding



Summary

- Having a sparse representation plays a fundamental role in how well we can
 - ▶ compress
 - ▶ denoise
 - ▶ restoresignals and images
- The above were accomplished with relatively simple algorithms (in practice, we use similar ideas + a bag a tricks)
- Geometrical representation \longrightarrow better results