

APPLICATION: WAVELET DENOISING

The Discrete Wavelet Transform

The discrete wavelet transform (DWT) is a linear map

$$W^T: \mathbb{R}^N \rightarrow \mathbb{R}^N, \quad \underline{x} \mapsto \underline{y} = W^T \underline{x}$$

satisfying certain special properties.

Although a thorough and rigorous definition and treatment of the DWT is beyond our scope, we can understand it through analogy with the discrete Fourier transform (DFT).

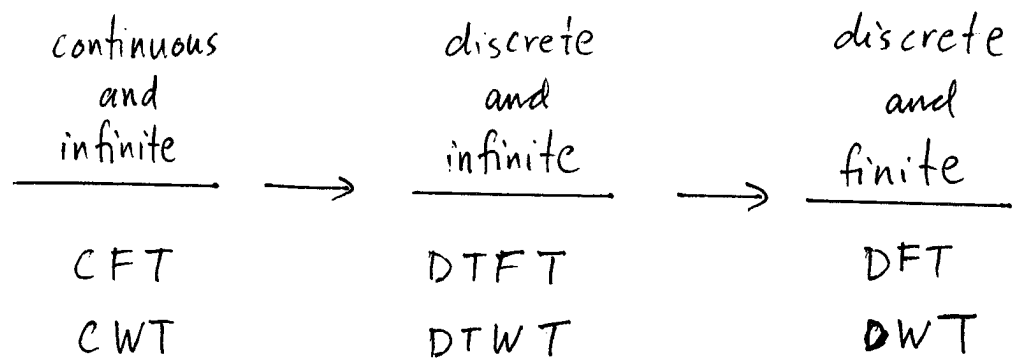
DWT vs. DFT

- Both can be represented by orthogonal matrices
- Both have efficient implementations

$$\text{DFT: } O(N \log N)$$

$$\text{DWT: } O(N)$$

- Both are discretizations of continuous transforms



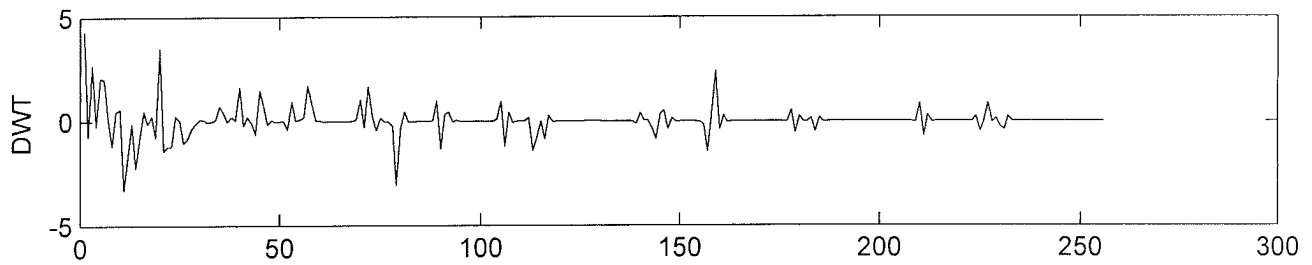
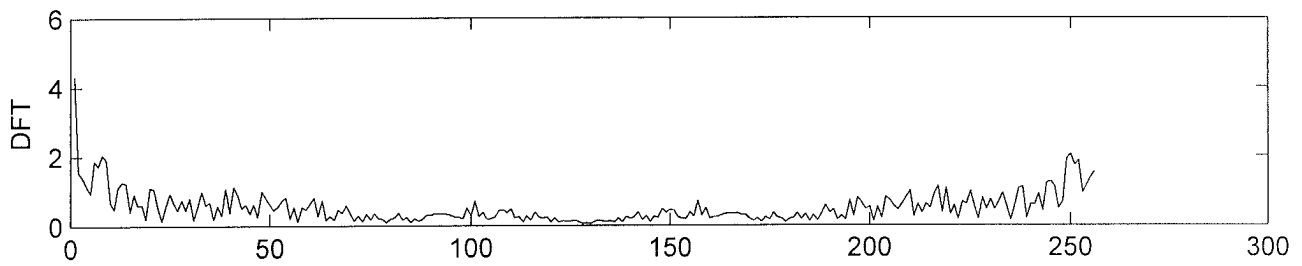
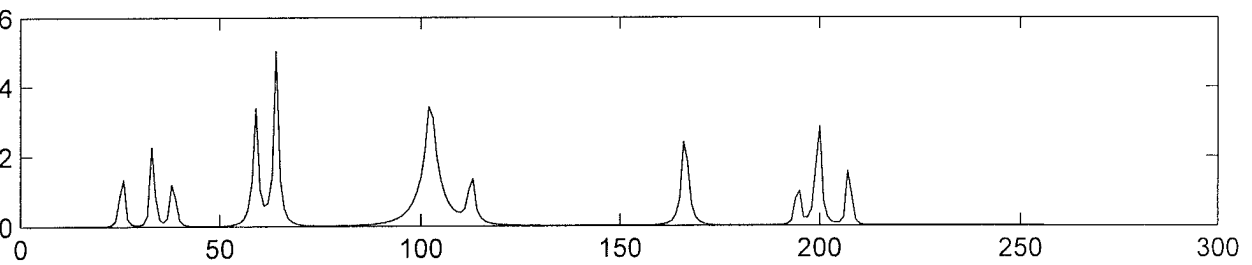
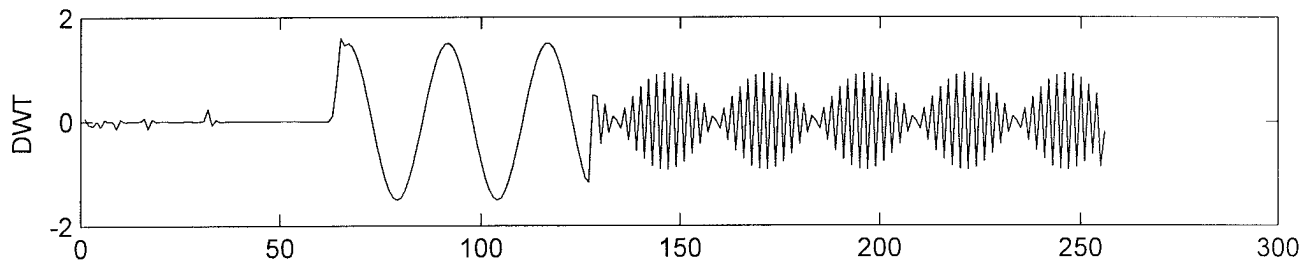
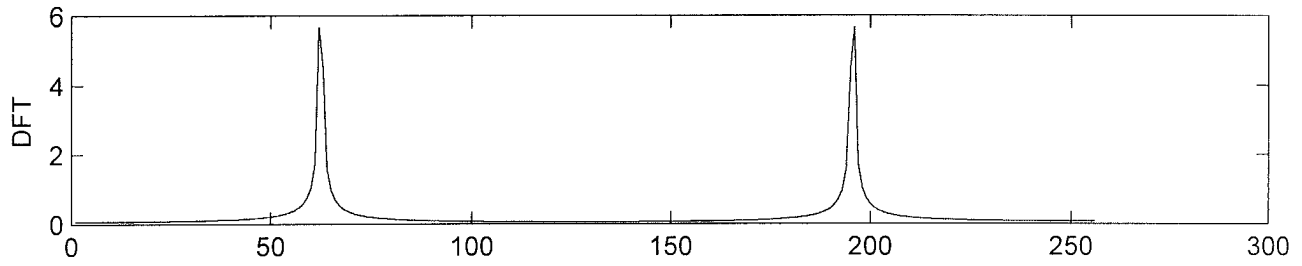
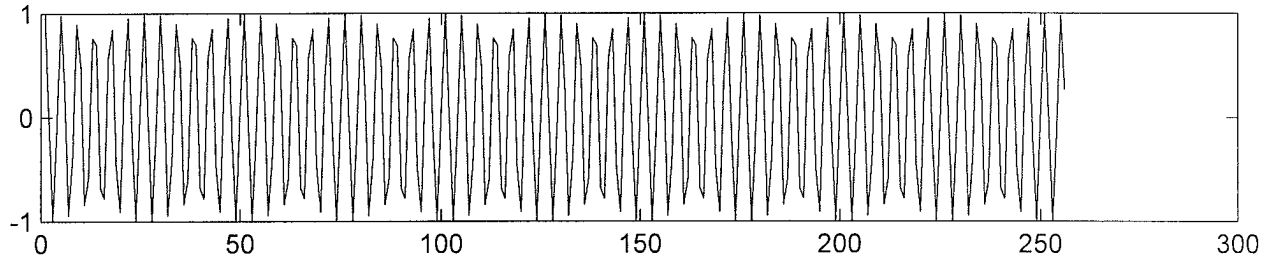
- Both are "change of basis" operators that compute the expansion coefficients of the signal in a different basis

DFT \implies Fourier basis

sparse representation of
sinusoidal signals

DWT \implies Wavelet basis (columns of W)

sparse representation of
piecewise polynomial signals



Haar Wavelet Transform

While there's only one DFT, there are in fact many different DWTs. The simplest is the Haar wavelet transform.

Consider a signal of length $N = 2^L$.

$$\underline{x} = [x(1) \ x(2) \ \dots \ x(N)]^T$$

Define

$$\underline{y}_1 = [c_1(1) \ c_1(2) \ \dots \ c_1(\frac{N}{2}) \ | \ d_1(1) \ \dots \ d_1(\frac{N}{2})]^T$$

where

$$c_1(1) = \frac{x(1) + x(2)}{\sqrt{2}}$$

$$d_1(1) = \frac{x(1) - x(2)}{\sqrt{2}}$$

$$c_1(2) = \frac{x(3) + x(4)}{\sqrt{2}}$$

$$d_1(2) = \frac{x(3) - x(4)}{\sqrt{2}}$$

⋮

Observe

- This transformation is invertible: we can recover \underline{x} from \underline{y}_1
- The transformation is an orthogonal linear map

$$\begin{bmatrix} c_1(1) \\ c_1(2) \\ c_1(3) \\ \vdots \\ d_1(1) \\ d_1(2) \\ d_1(3) \\ \vdots \end{bmatrix} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 & & & & & & \\ & & 1 & 1 & & & & \\ & & & & 1 & 1 & & \\ & & & & & & \ddots & \\ & & & & & & & \ddots \\ 1 & -1 & & & & & & \\ & & 1 & -1 & & & & \\ & & & & 1 & -1 & & \\ & & & & & & \ddots & \\ & & & & & & & \ddots \end{bmatrix} \begin{bmatrix} x(1) \\ x(2) \\ \vdots \\ x(N) \end{bmatrix} = W^T \underline{x}$$

- The coefficients $c_1(n)$ are local averages and represent coarse information about the signal
- The coefficients $d_1(n)$ are local differences and represent detailed information

This operation is called the level-1 Haar wavelet transform. The idea behind the general Haar wavelet transform is to recursively apply this operation to the coarse coefficients.

$$\underline{y}_1 = [c_1(1) \quad \dots \quad c_1(\frac{N}{2}) \mid d_1(1) \quad \dots \quad d_1(\frac{N}{2})]^T$$

$$\underline{y}_2 = [c_2(1) \quad \dots \quad c_2(\frac{N}{4}) \mid d_2(1) \quad \dots \quad d_2(\frac{N}{4}) \mid d_1(1) \quad \dots \quad d_1(\frac{N}{2})]^T$$

$$\underline{y}_3 = [c_3(1) \dots c_3(\frac{N}{8}) \mid d_3(1) \dots d_3(\frac{N}{8}) \mid d_2(1) \dots d_2(\frac{N}{4}) \mid d_1(1) \dots d_1(\frac{N}{2})]^T$$

⋮

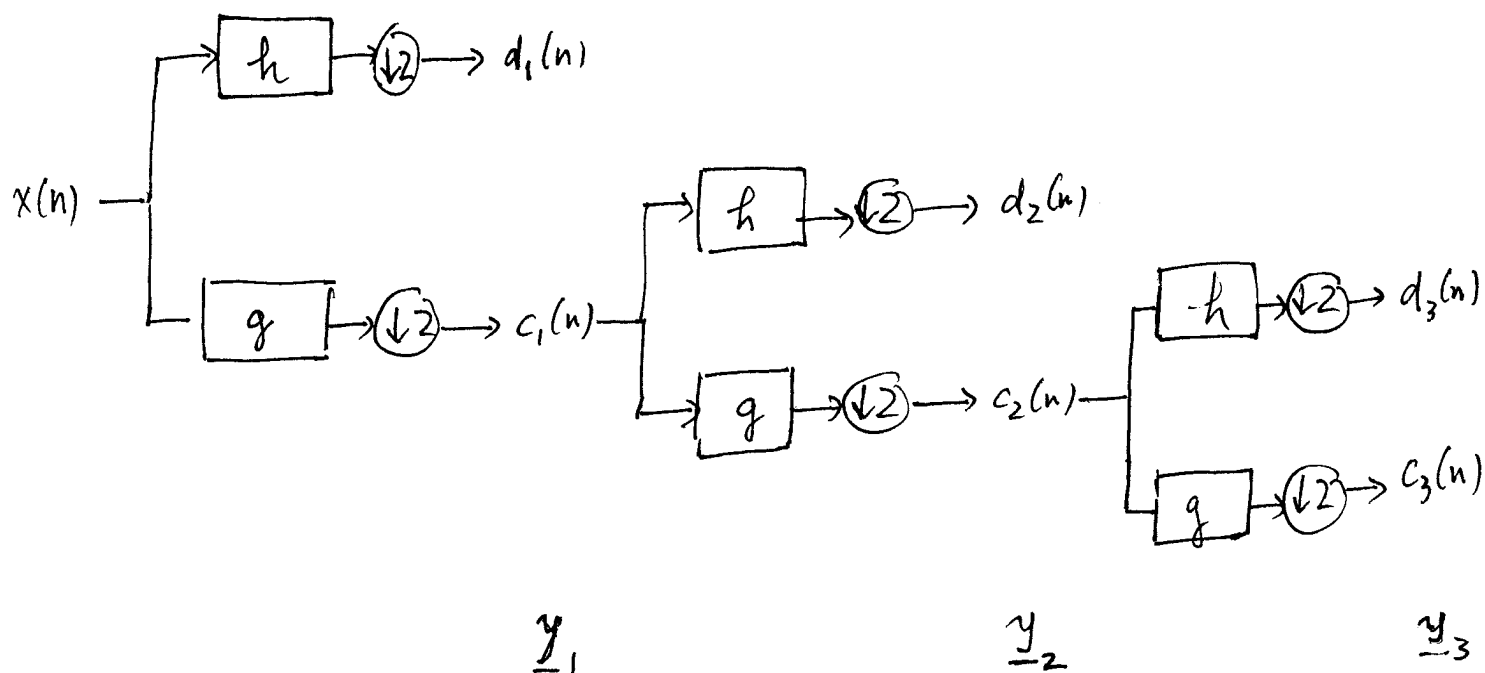
$$\underline{y}_L = [c_L(1) \mid d_L(1) \mid d_{L-1}(1) \quad d_{L-1}(2) \mid d_{L-2}(1) \quad \dots \quad d_{L-2}(4) \mid \dots]^T$$

We call \underline{y}_l the level- l Haar wavelet transform

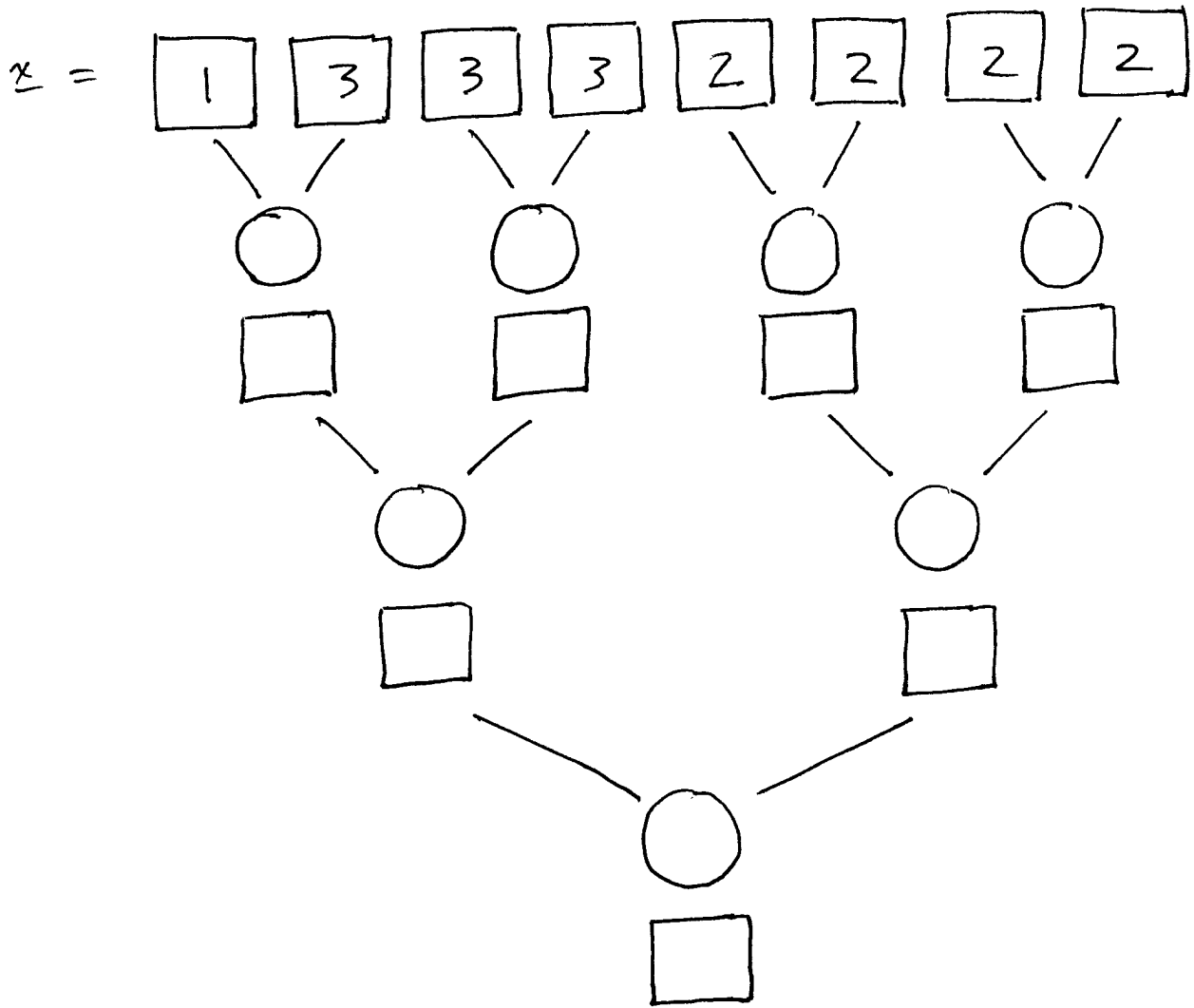
Filter bank implementation:

$$h = \begin{bmatrix} \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix} \quad (\text{high pass})$$

$$g = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix} \quad (\text{low pass})$$



Example



$$\underline{y}_1 = [\quad \quad \quad]^T$$

$$\underline{y}_2 = [\quad \quad \quad]^T$$

$$\underline{y}_3 = [\quad \quad \quad]^T$$

Important things to notice :

- the detail coefficients are zero where the signal is constant. In particular, if $x(n)$ is constant on the interval

$$[k \cdot 2^l + 1, k \cdot 2^l + 2, \dots, (k+1) \cdot 2^l], \text{ then } d_l(k) = 0$$

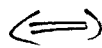
\Rightarrow sparsity

- the detail coefficients have a natural hierarchical (or tree-structured) arrangement; we can say that $d_l(k)$ is the parent of $d_{l-1}(2k-1)$ and $d_{l-1}(2k)$, who are its children
- $c_l(n)$ is a low resolution approximation to $x(n)$; it is the result of averaging and downsampling $x(n)$ l times
- different levels capture different resolutions of detail :

$$\{d_1(k)\}_{k=1}^{2^{L-1}}$$



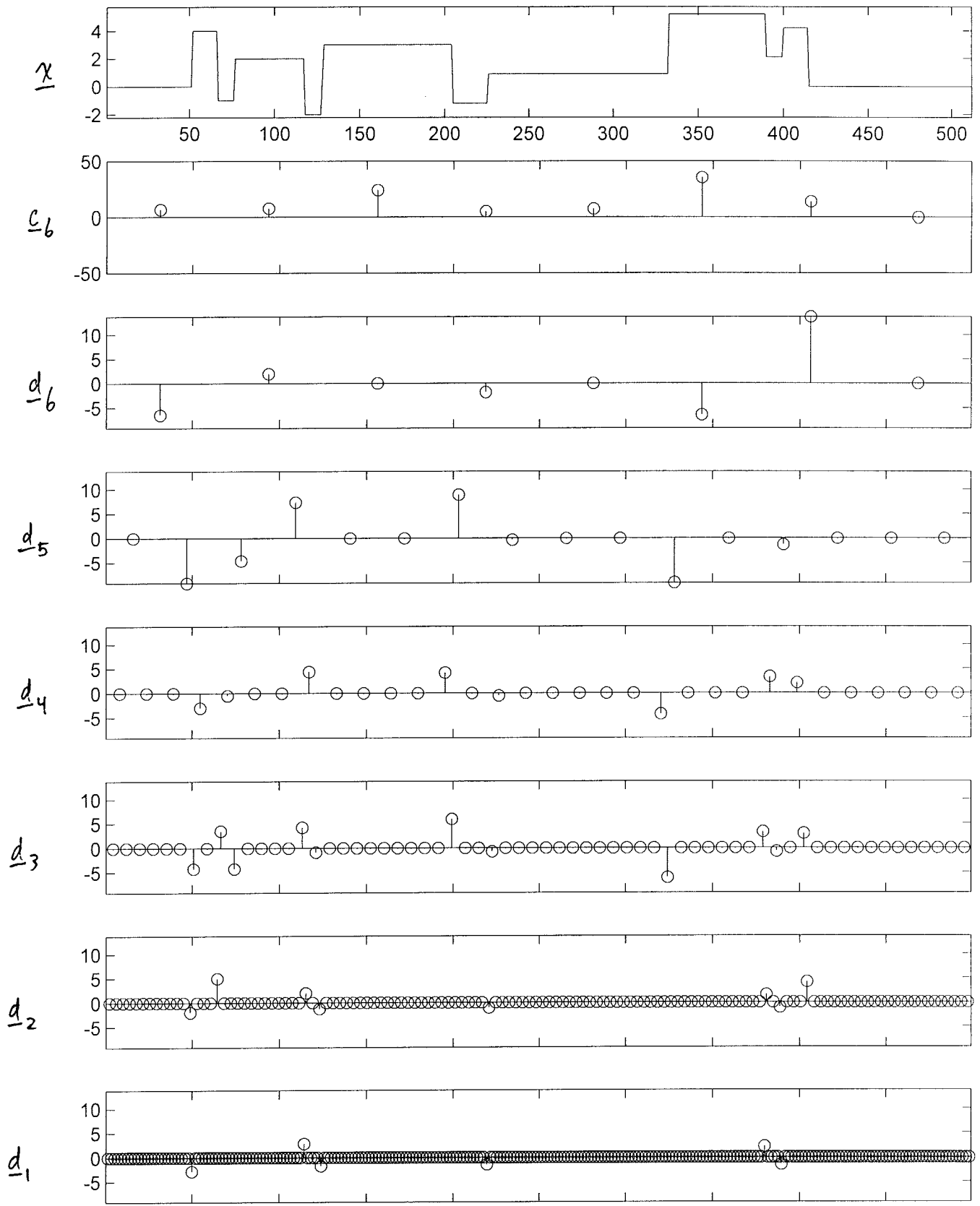
$$\begin{aligned} &\{d_{L-1}(k)\}_{k=1}^2 \\ &\{d_L(k)\}_{k=1}^1 \end{aligned}$$



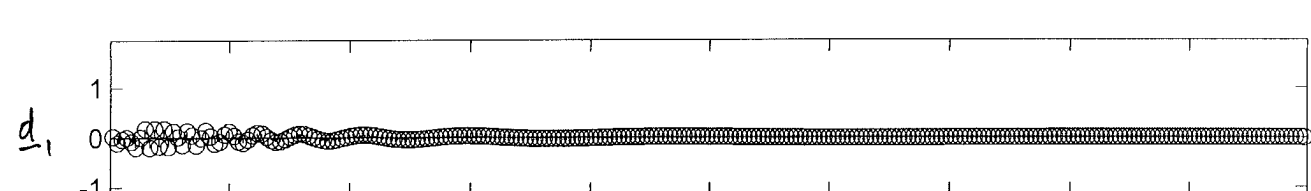
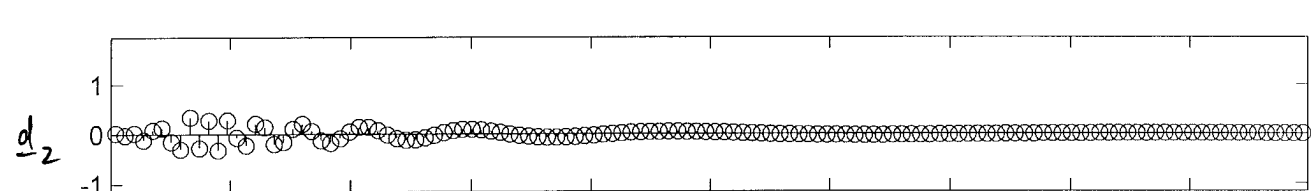
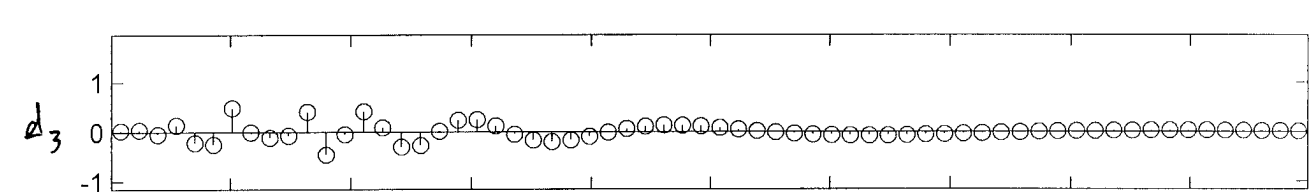
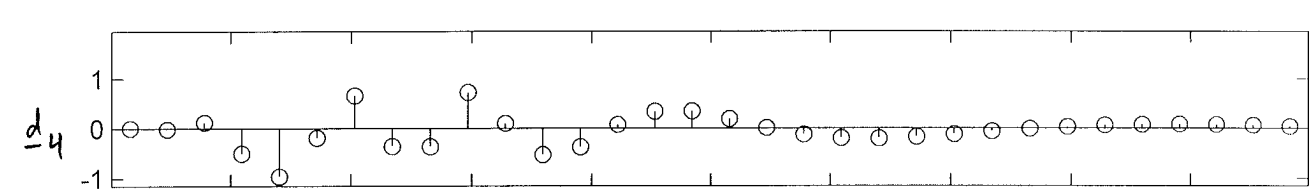
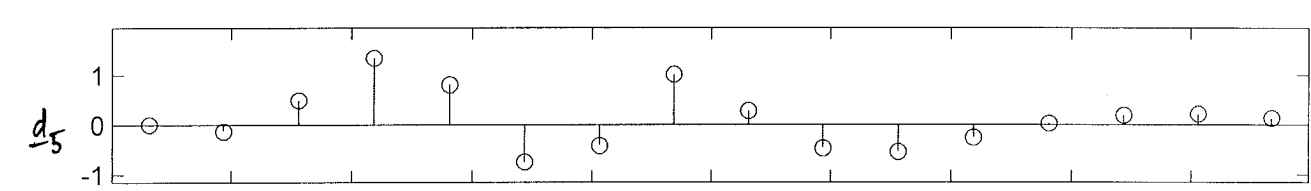
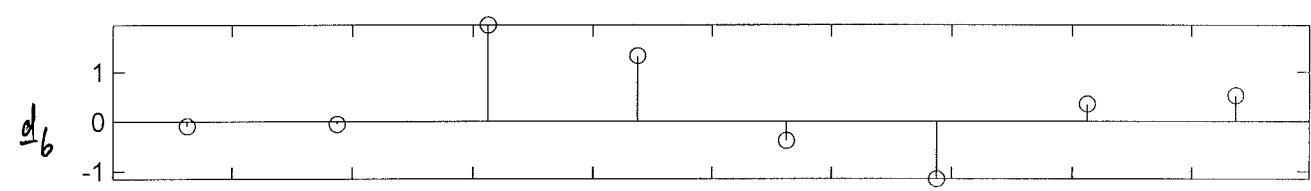
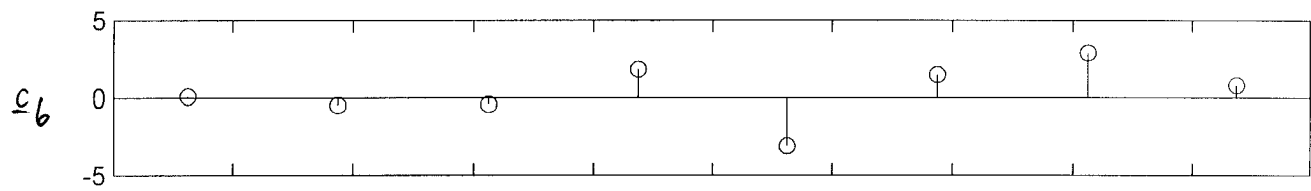
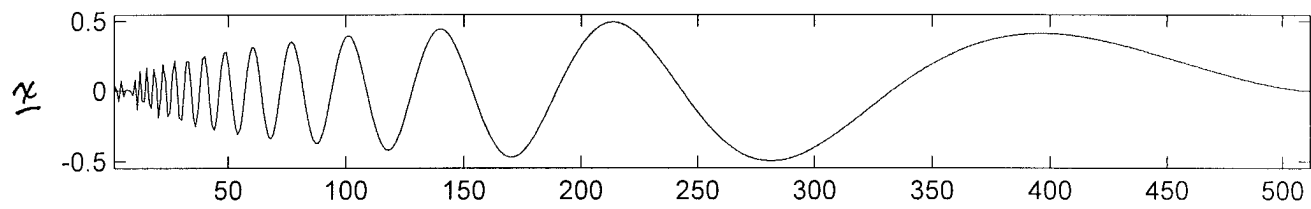
high frequency detail



low frequency detail



\Rightarrow sparse representation



Other Wavelet Transforms

The Haar wavelet transform can be generalized by using different high-pass and low-pass filters h & g . These filters must satisfy certain properties for the resulting transform to be orthogonal (and qualify as a DWT).

The most important generalization are the Daubechies DWT. They are based on certain filters h_p, g_p such that

- the length of h_p and g_p is 2^p .
- if a signal behaves locally like a $(p-1)^{\text{th}}$ order polynomial, the corresponding detail coefficients are zero

Examples

$$p=1 \implies \text{Haar}$$

$$p=2 \implies g = [.4830 \quad .8365 \quad .2241 \quad -.1294]$$

$$h = [.1294 \quad .2241 \quad -.8365 \quad .4830]$$

Wavelet Denoising

Suppose we measure a noisy signal

$$\underline{x} = \underline{s} + \underline{v}$$

(*)

and assume

- $\underline{s} = [s_1, \dots, s_N]^T$ has a sparse representation in a certain wavelet basis, e.g., \underline{s} is piecewise constant / Haar basis
- $\underline{v} \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{I})$

Now take the wavelet transform of (*):

$$\underline{y} = \underline{\theta} + \underline{z}$$

$$\begin{aligned} \underline{y} &= W^T \underline{x} \\ \underline{\theta} &= W^T \underline{s} \\ \underline{z} &= W^T \underline{v} \end{aligned}$$

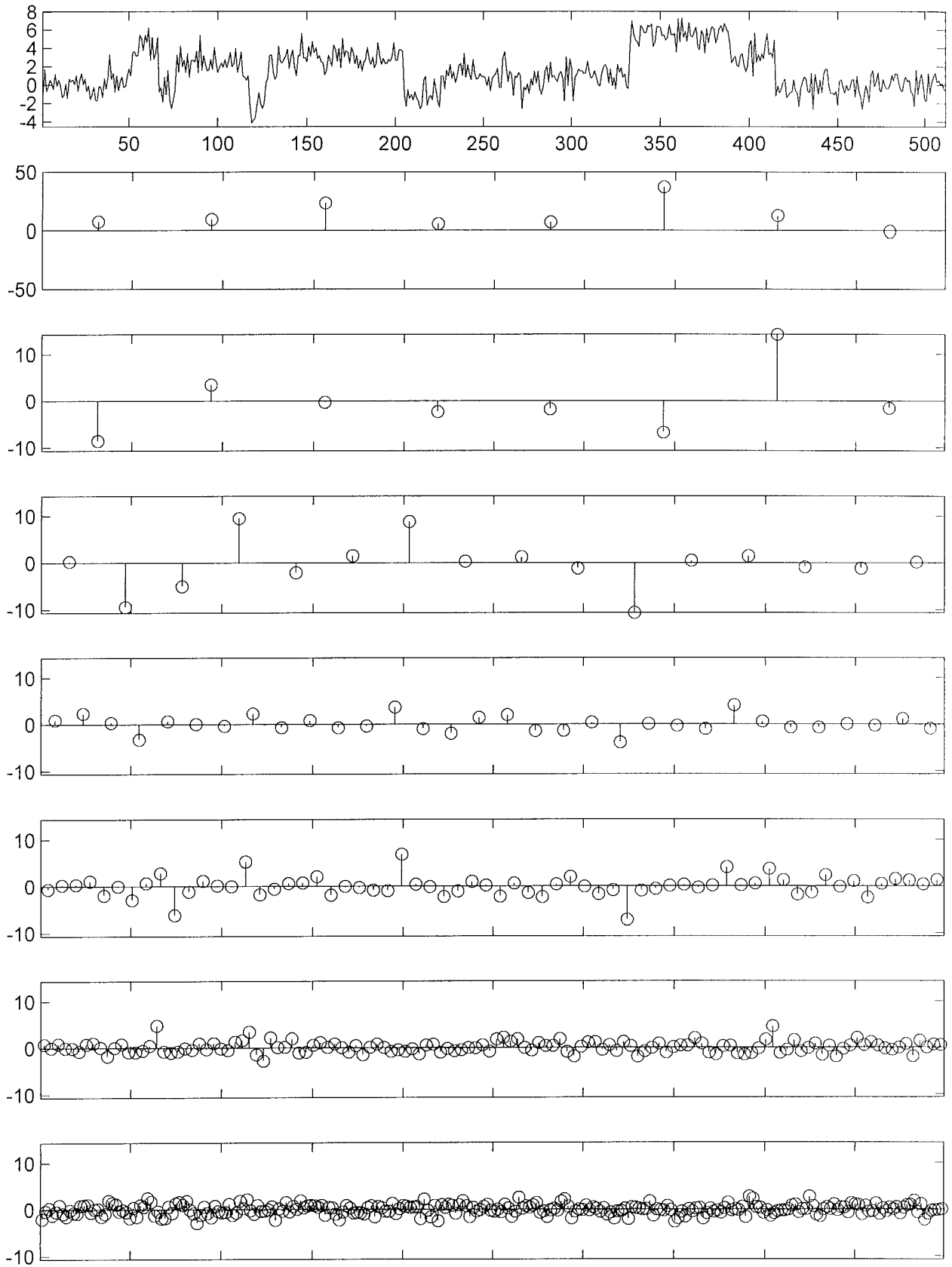
Then we know

- most elements in

$$\underline{\theta} = [\theta_1, \dots, \theta_N]^T$$

are zero or very close to zero

- $\underline{z} \sim \mathcal{N}(\underline{0}, \sigma^2 \mathbf{I})$
because $W^T W = \mathbf{I}$



Since W is orthogonal, the estimation problem amounts to recovery of a signal in iid Gaussian noise, whether we treat the problem in the "time" domain or the "wavelet" domain

From a Bayesian perspective, we should tackle the problem in the domain for which it is easiest to specify a prior.

On one hand, the fact that θ is sparse suggests a subspace model. Unfortunately we don't know a priori which detail coefficients will be zero.

Q: How can we take advantage of the prior knowledge that θ is sparse? What statistical model captures this information?

A: One solution is to employ a _____

Mixture Modeling

View the detail coefficients $\theta_2, \dots, \theta_N$
as realizations of a single random variable θ .

We know

- Most θ_i are small (sparsity assumption)
- Some θ_i are large (W is orthogonal, so energy must be preserved)
- θ is zero mean, since θ_i are local differences
- θ_i are "approximately" independent, since the θ_i are local differences

This suggests the following prior:

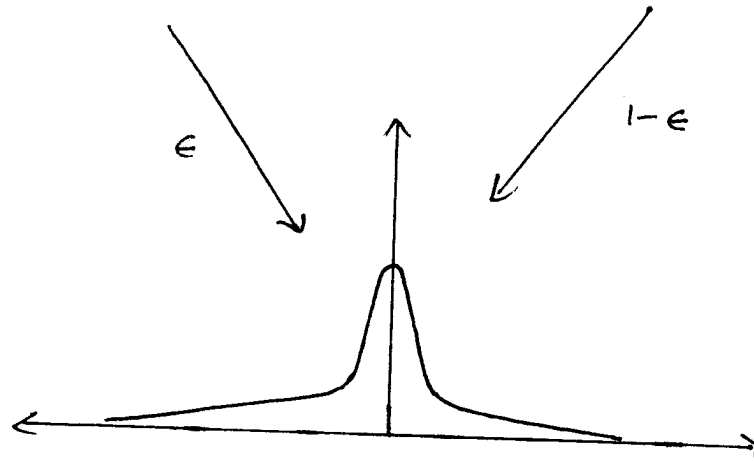
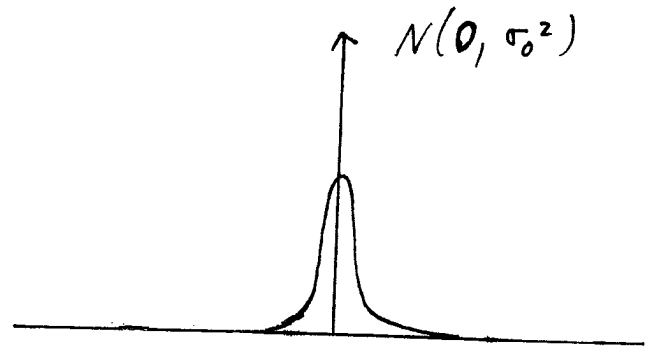
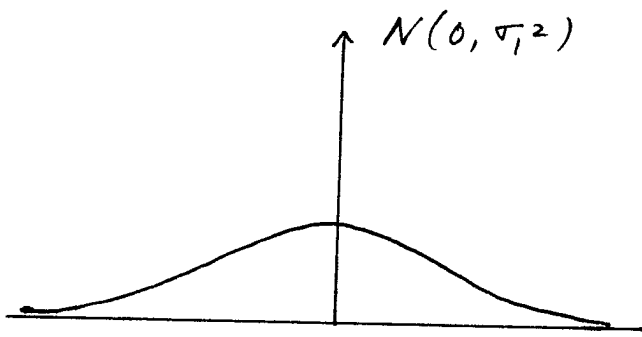
$$\theta_i \stackrel{\text{iid}}{\sim} \epsilon \mathcal{N}(0, \sigma_1^2) + (1-\epsilon) \mathcal{N}(0, \sigma_0^2).$$

where $\sigma_0^2 \ll \sigma_1^2$ and

ϵ = proportion of "significant" coefficients

σ_1^2 = variance of significant "

σ_0^2 = " " insignificant "



According to this prior, a detail coefficient θ is generated according to the following algorithm:

1. Flip an "e-coin"
2. If heads,
 $\theta \sim N(0, \sigma_1^2)$
- Else
 $\theta \sim N(0, \sigma_0^2)$

The Big Picture

- We observe $\underline{x} = \underline{s} + \underline{v}$, $\underline{v} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$
- Compute $\underline{y} = \underline{\theta} + \underline{w}$ by taking wavelet transform
- View

$$y_i = \theta_i + w_i , \quad w_i \sim \mathcal{N}(0, \sigma^2)$$

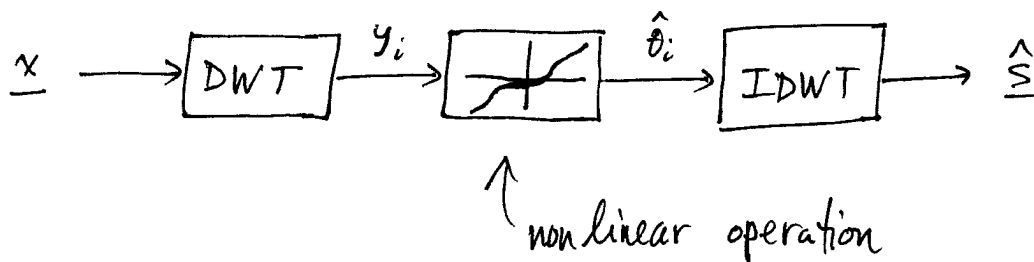
as independent estimation problems.

- Leave "coarse" coefficients unprocessed: noise will be "averaged out" since these are local averages
- Assume a mixture model prior on the detail coefficients and estimate

$$\hat{\theta}_i = \mathbb{E}[\theta_i | y_i]$$

- Apply inverse wavelet transform to obtain

$$\hat{\underline{s}} = \mathbf{W} \hat{\underline{\theta}}$$



Setting Parameters

Before this method is practical, we need ways to set σ^2 , ϵ , σ_1^2 , and σ_0^2 .

Donoho and Johnstone suggested the estimate

$$\hat{\sigma} = \frac{\text{MAD}(y_i)_{i > N/2}}{.6745},$$

which takes the "median absolute deviation" of the wavelet coefficients at the "finest" level of detail, and .6745 makes the estimate unbiased if all of the θ_i are in fact 0.

The mixture model parameters ϵ , σ_1^2 , σ_0^2 may be estimated via maximum likelihood:

$$(\hat{\epsilon}, \hat{\sigma}_1^2, \hat{\sigma}_0^2) = \arg \max_{(\epsilon, \sigma_1^2, \sigma_0^2)} l(\epsilon, \sigma_1^2, \sigma_0^2; \underline{y})$$

Exercise | Determine a formula for the likelihood

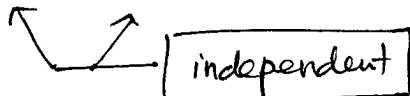
of ϵ , σ_1^2 , σ_0^2 given the detail coefficients $\underline{y}_{\text{detail}}$

$= [y_2, \dots, y_N]^T$ (assuming a max-level wavelet transform)

Solution | Denote

$$\phi(y; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y-\mu)^2}{2\sigma^2}\right\}.$$

Since

$$y_i = \theta_i + w_i$$


A diagram with two arrows pointing from the terms θ_i and w_i in the equation above to a rectangular box containing the word "independent".

where

$$f(\theta_i) = \epsilon \phi(\theta_i; 0, \sigma_1^2) + (1-\epsilon) \phi(\theta_i; 0, \sigma_0^2)$$

$$f(w_i) = \phi(w_i; 0, \sigma^2)$$

it follows that

$$f(y_i) = f(\theta_i) * f(w_i)$$

$$= \epsilon \phi(y_i; 0, \sigma^2 + \sigma_1^2) + (1-\epsilon) \phi(y_i; 0, \sigma^2 + \sigma_0^2)$$

Hence the likelihood of $\epsilon, \sigma_1^2, \sigma_0^2$ is

$$l(\epsilon, \sigma_1^2, \sigma_0^2; \underline{y}_{\text{detail}}) = \prod_{i=2}^N f(y_i; \epsilon, \sigma_1^2, \sigma_0^2)$$

$$= \prod_{i=2}^N \left[\epsilon \phi(y_i; 0, \sigma^2 + \sigma_1^2) + (1-\epsilon) \phi(y_i; 0, \sigma^2 + \sigma_0^2) \right]$$

Typically one uses an iterative algorithm such as an EM algorithm to fit mixture models. However, because there are only 3 unknowns, we could also maximize the likelihood by an exhaustive grid search.

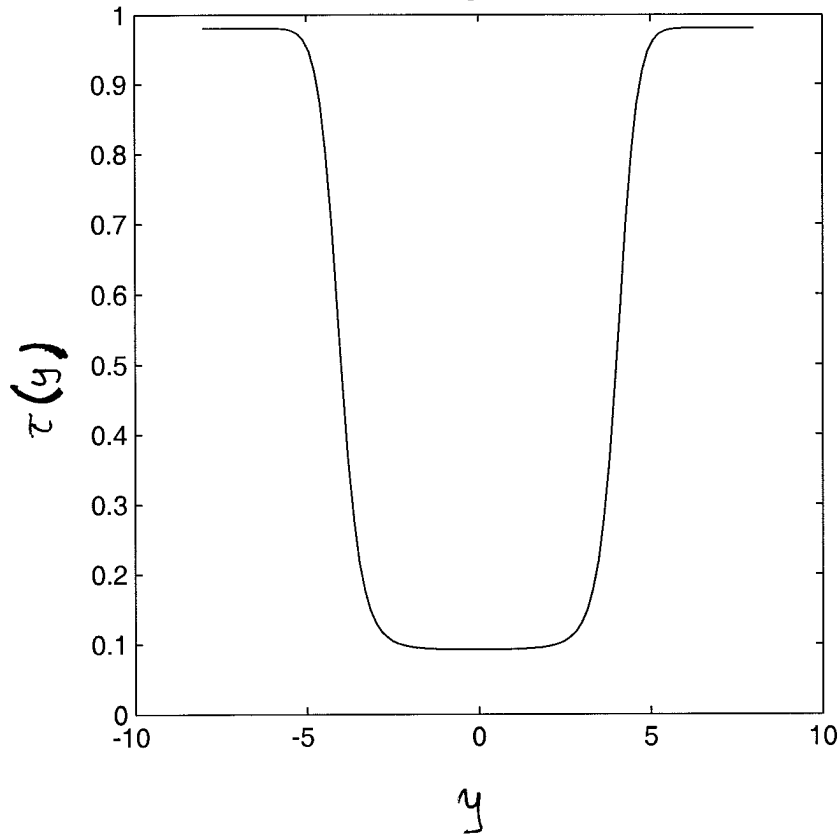
Note that we are using the data to determine our prior - hence we are not strictly adhering to the Bayesian philosophy. This kind of procedure is called an empirical Bayesian method.

You will show on the homework that

$$\hat{\theta} = E[\theta | y] = \tau(y) \cdot y$$

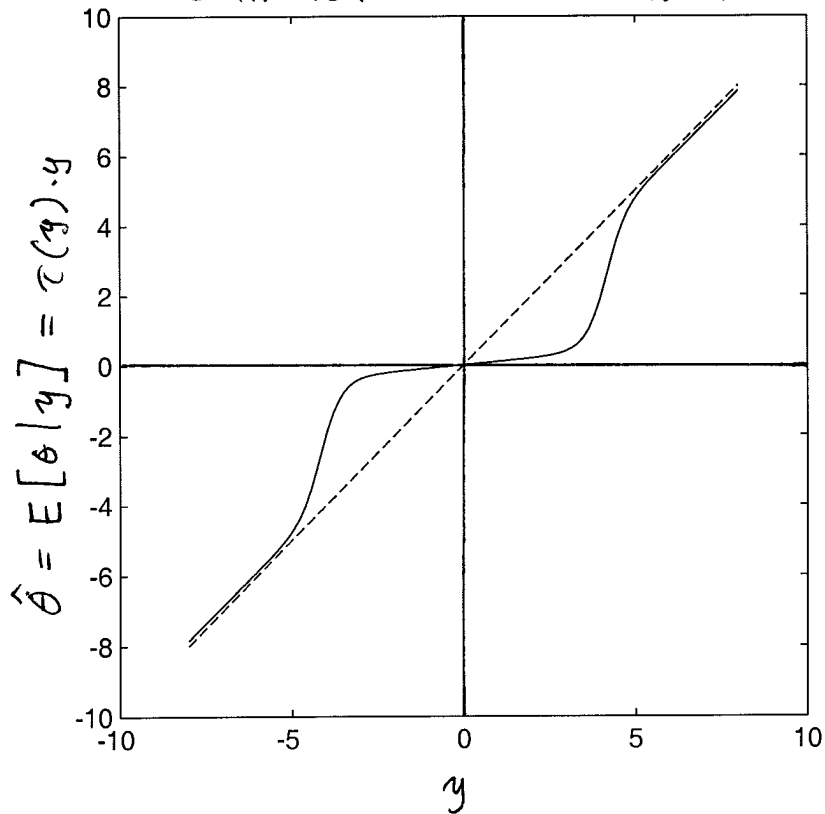
where $0 < \tau(y) < 1$ is called the shrinkage factor.

the shrinkage factor



$\epsilon = 0.05$
 $\sigma^2 = 1$
 $\sigma_0^2 = .1$
 $\sigma_1^2 = 50$

estimated wavelet coefficient



The effect of the shrinkage factor is that

- small coefficients are set nearly to zero
- large coefficients are virtually unaltered.

This property is consistent with our understanding

- small coefficients are mostly noise
- large coefficients contain actual signal

Extensions

- Image denoising
- More sophisticated priors

Summary

Gaussian noise
Mixture of Gaussians prior } \Rightarrow Wavelet shrinkage