

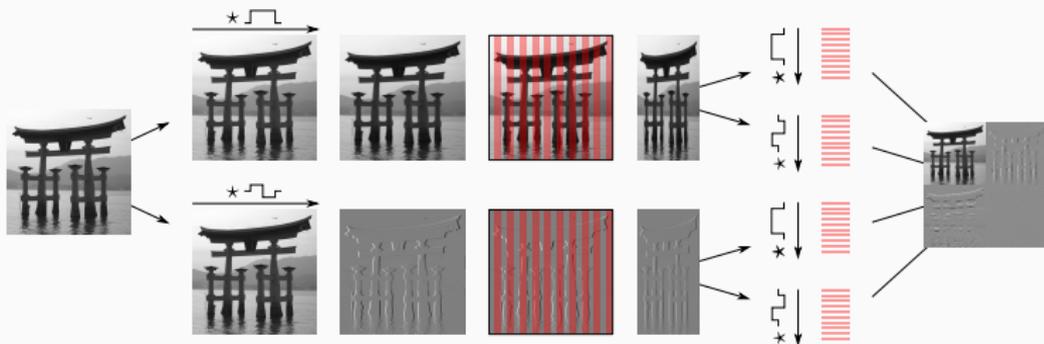
ECE 285

Image and video restoration

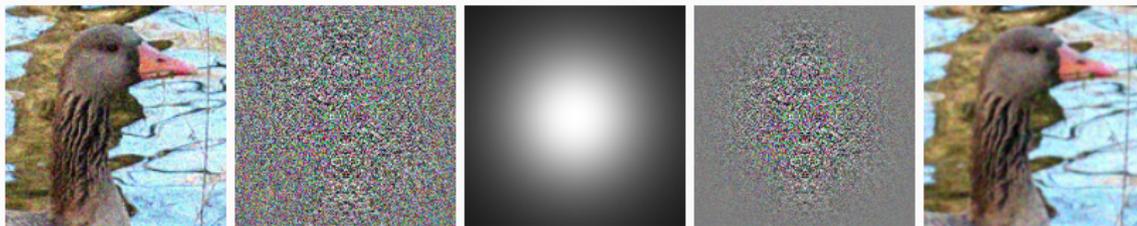
Chapter VI – Sparsity, shrinkage and wavelets

Charles Deledalle

June 7, 2019



Motivations



(a) $y = x + w$ (b) $z = \mathbf{F}y$ (c) $\frac{\lambda_i^2}{\lambda_i^2 + \sigma^2}$ (d) $\hat{z}_i = \frac{\lambda_i^2}{\lambda_i^2 + \sigma^2} z_i$ (e) $\hat{x} = \mathbf{F}^{-1} \hat{z}$

Wiener filter (LMMSE in the Fourier domain)

- Assume Fourier coefficients to be decorrelated (white),
- Modulate frequencies based on the mean power spectral density λ_i^2 .

Limits

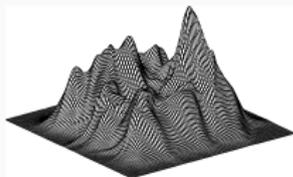
- Linear: no adaptation to the content \Rightarrow **Unable to preserve edges, Blurry solutions.**

Facts and consequences

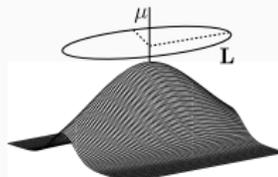
- Assume Fourier coefficients to be decorrelated (white)
- Removing Gaussian noise \Rightarrow need to be adaptive \Rightarrow Non linear
- Assuming Gaussian noise + Gaussian prior \Rightarrow Linear

Deductive reasoning

Fourier coefficients of clean images are not Gaussian distributed



Underlying prior $x \mapsto p(x)$



Gaussian prior $x \sim \mathcal{N}(\mu; \mathbf{L})$

How are Fourier coefficients distributed?

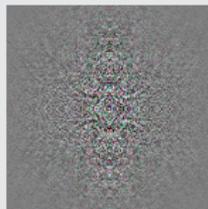
How are Fourier coefficients distributed?

1. Perform whitening with DFT

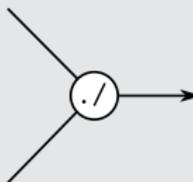


A clean image x

E^T
 $\propto \text{fft2}$



Standard deviation λ_i
for each frequency



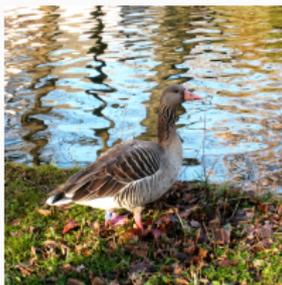
Whitening $\eta = \mathbf{L}^{-1/2}x$
close to white noise

$$\text{Var}[x] = \mathbf{L} = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^* \quad \text{with} \quad \mathbf{E} = \frac{1}{\sqrt{n}}\mathbf{F}$$
$$\text{diag}(\mathbf{\Lambda}) = (\lambda_1^2, \dots, \lambda_n^2) = n^{-1}\text{MPSD}$$

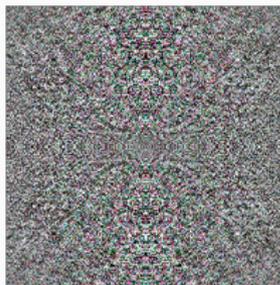
How are Fourier coefficients distributed?

2. Look at the histogram

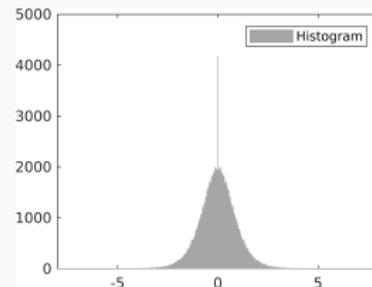
- The histogram of η has a symmetric bell shape around 0.
- It has a peak at 0 (a large number of Fourier coefficients are zero).
- It has large/heavy tails (many coefficients are “outliers”/abnormal).



(a) x



(b) Whitening η of x

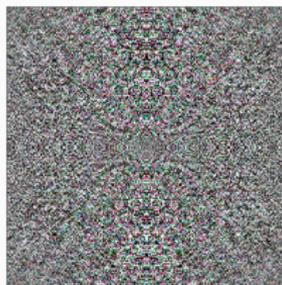


(c) Histogram of η

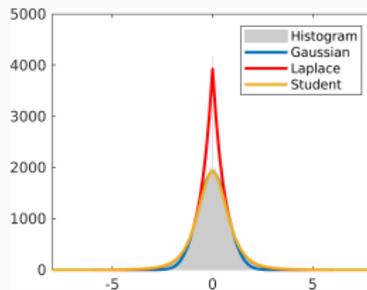
How are Fourier coefficients distributed?

3. Look for the distribution that best fits (in log scale)

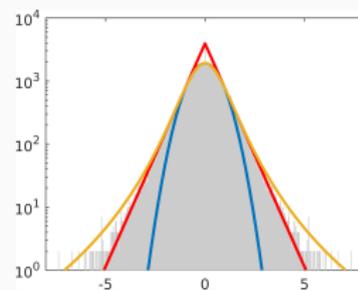
- Gaussian: bell shape ✓, peak ✗, tail ✗
 - Laplacian: bell shape ✗, peak ✓, tail ✓
 - Student: bell shape ✓, peak ✗, tail ✓ (heavier)
-
- Others: alpha stables and generalized Gaussian distributions



(a) Whitening η of x



(b) Histogram of η



(c) Log-histogram of η

Motivations – Distribution of Fourier coefficients

Model expression (zero mean, variance = 1)

- Gaussian: bell shape ✓, peak ✗, tail ✗

$$p(\eta_i) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\eta_i^2}{2}\right)$$

- Laplacian: bell shape ✗, peak ✓, tail ✓

$$p(\eta_i) = \frac{1}{\sqrt{2}} \exp\left(-\sqrt{2}|\eta_i|\right)$$

- Student: bell shape ✓, peak ✗, tail ✓ (heavier)

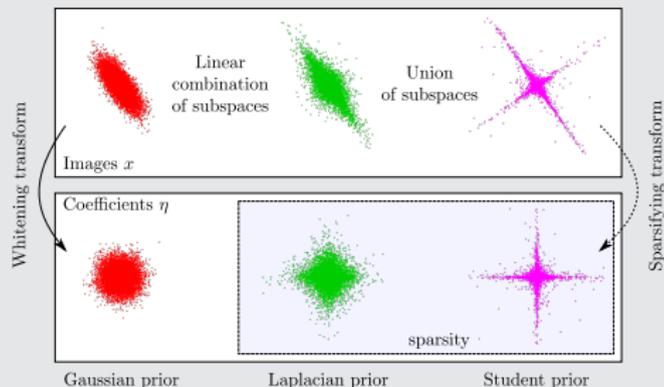
$$p(\eta_i) = \frac{1}{Z} \left(\frac{1}{(2r-2) + \eta_i^2} \right)^{r+1/2}$$

(Z normalization constant, $r > 1$ controls the tails)

How do they look in multiple-dimensions?

Motivations – Distribution of Fourier coefficients

- Gaussian prior $\left\{ \begin{array}{l} \bullet \text{ images are concentrated in an elliptical cluster,} \\ \bullet \text{ outliers are rare (images outside the cluster).} \end{array} \right.$
- Peaky & heavy tailed priors: shape **between a diamond and a star**.



- **union of subspaces:** most images lie in one of the branches of the star,
- **sparsity:** most of their coefficients η_i are zeros,
- **robustness:** outlier coefficients are frequent.

Shrinkage functions

Consider the following Gaussian denoising problem

- Let $y \in \mathbb{R}^n$ and $x \in \mathbb{R}^P$ be two random vectors such that

$$y | x \sim \mathcal{N}(x, \sigma^2 \text{Id}_n)$$

$$\mathbb{E}[x] = 0 \quad \text{and} \quad \text{Var}[x] = L = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^*$$

- Let $\eta = \mathbf{\Lambda}^{-1/2} \mathbf{E}^* x$ (whitening / decorrelation of x)

Goal: estimate x from y
assuming a non-Gaussian prior p_η for η .
(such as Laplacian or Student)

Shrinkage functions

Bayesian shrinkage functions

- Assume η_i are also **independent and identically distributed (iid)**.
- Then, the MMSE and MAP estimators both read as

$$\hat{x}^* = \underbrace{\mathbf{E}\hat{z}}_{\text{Come back}} \quad \text{where} \quad \underbrace{\hat{z}_i = s(z_i; \lambda_i, \sigma)}_{\text{shrinkage}} \quad \text{and} \quad z = \underbrace{\mathbf{E}^*y}_{\text{Change of basis}}$$

- The function $z_i \mapsto s(z_i; \lambda_i, \sigma)$ is called **shrinkage function**.
- Unlike the LMMSE, s will depend on the prior distribution of η_i .
- As for the LMMSE, the solution can be computed in the eigenspace.
- We say that the estimator is **separable** in the eigenspace (ex: Fourier).

Remark

independence \Rightarrow uncorrelation
 \neg uncorrelation \Rightarrow \neg independence
correlation \Rightarrow dependence

\Rightarrow

Whitening is a necessarily step
for independence but not a
sufficient one.
(Except in the Gaussian case)

How are the shrinkage functions defined for the MMSE and MAP?

Shrinkage functions

- Recall that the MMSE is the posterior mean

$$\hat{x}^* = \int_{\mathbb{R}^n} xp(x|y) dx = \frac{\int_{\mathbb{R}^n} xp(y|x)p(x) dx}{\int_{\mathbb{R}^n} p(y|x)p(x) dx}$$

MMSE Shrinkage functions

- Under the previous assumptions

$$\hat{x}^* = \underbrace{\mathbf{E} \hat{z}}_{\text{Come back}} \quad \text{where} \quad \underbrace{\hat{z}_i = s(z_i; \lambda_i, \sigma)}_{\text{shrinkage}} \quad \text{and} \quad z = \underbrace{\mathbf{E}^* y}_{\text{Change of basis}}$$

$$\text{with} \quad s(z; \lambda, \sigma) = \frac{\int_{\mathbb{R}} \tilde{z} \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2}\right) p_{\eta}\left(\frac{\tilde{z}}{\lambda}\right) d\tilde{z}}{\int_{\mathbb{R}} \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2}\right) p_{\eta}\left(\frac{\tilde{z}}{\lambda}\right) d\tilde{z}}$$

where p_{η} is the prior distribution on the entries of η .

- Separability: n dimensional optimization $\longrightarrow n \times 1d$ integrations.

Shrinkage functions

- Recall that the MAP is the optimization problem

$$\hat{x}^* \in \operatorname{argmax}_{x \in \mathbb{R}^n} p(x|y) = \operatorname{argmin}_{x \in \mathbb{R}^n} [-\log p(y|x) - \log p(x)]$$

MAP Shrinkage functions

- Under the previous assumptions

$$\hat{x}^* = \underbrace{\mathbf{E} \hat{z}}_{\text{Come back}} \quad \text{where} \quad \underbrace{\hat{z}_i = s(z_i; \lambda_i, \sigma)}_{\text{shrinkage}} \quad \text{and} \quad z = \underbrace{\mathbf{E}^* y}_{\text{Change of basis}}$$

$$\text{with} \quad s(z; \lambda, \sigma) = \operatorname{argmin}_{\tilde{z} \in \mathbb{R}} \left[\frac{(z - \tilde{z})^2}{2\sigma^2} - \log p_\eta \left(\frac{\tilde{z}}{\lambda} \right) \right]$$

where p_η is the prior distribution on the entries of η .

- Separability: n dimensional integration $\rightarrow n \times 1d$ optimisations.

Example (Gaussian noise + Gaussian prior)

- MMSE Shrinkage

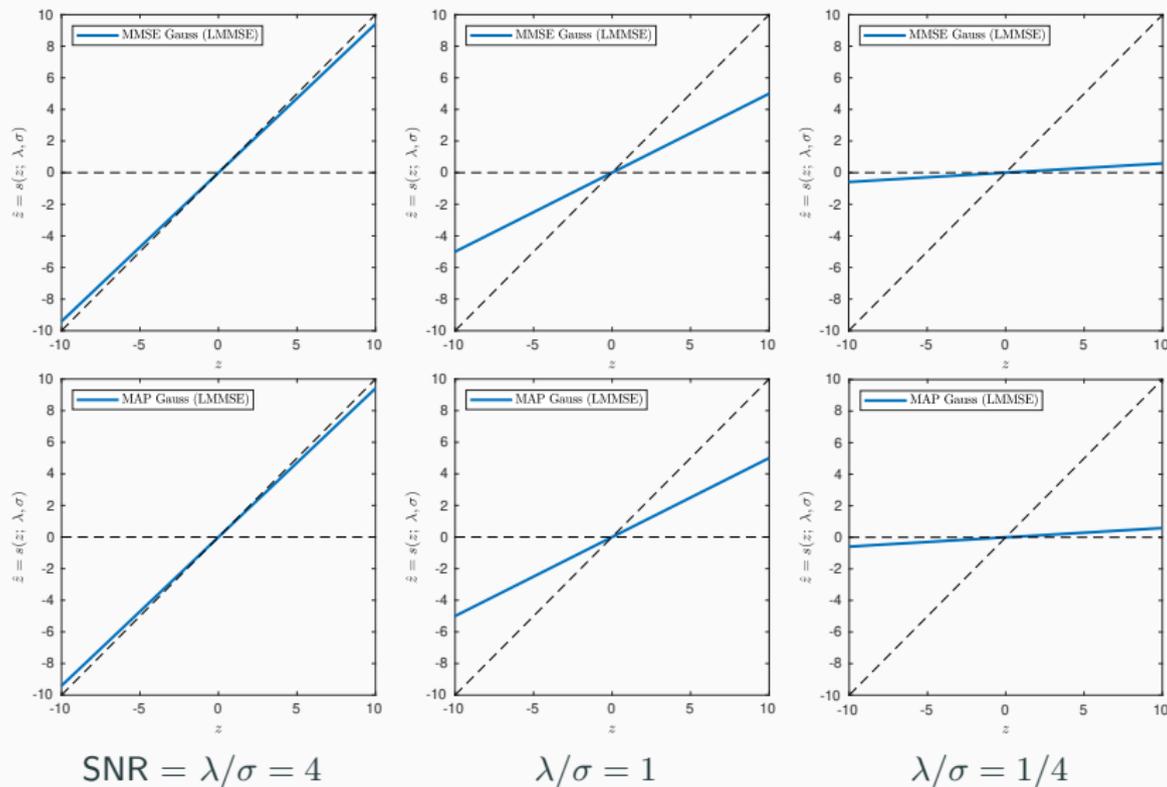
$$s(z; \lambda, \sigma) = \frac{\int_{\mathbb{R}} \tilde{z} \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2} - \frac{\tilde{z}^2}{2\lambda^2}\right) d\tilde{z}}{\int_{\mathbb{R}} \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2} - \frac{\tilde{z}^2}{2\lambda^2}\right) d\tilde{z}} = \frac{\lambda^2}{\lambda^2 + \sigma^2} z$$

- MAP Shrinkage

$$s(z; \lambda, \sigma) = \operatorname{argmin}_{\tilde{z} \in \mathbb{R}} \left[\frac{(z - \tilde{z})^2}{2\sigma^2} + \frac{\tilde{z}^2}{2\lambda^2} \right] = \frac{\lambda^2}{\lambda^2 + \sigma^2} z$$

- Gaussian prior: MAP = MMSE = Linear shrinkage.
- We retrieve the LMMSE as expected.

Gaussian noise + Gaussian prior



Example (Gaussian noise + Laplacian prior)

- MMSE Shrinkage

$$s(z; \lambda, \sigma) = \frac{\int \tilde{z} \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2} - \frac{\sqrt{2}|\tilde{z}|}{\lambda}\right) d\tilde{z}}{\int \exp\left(-\frac{(z-\tilde{z})^2}{2\sigma^2} - \frac{\sqrt{2}|\tilde{z}|}{\lambda}\right) d\tilde{z}}$$

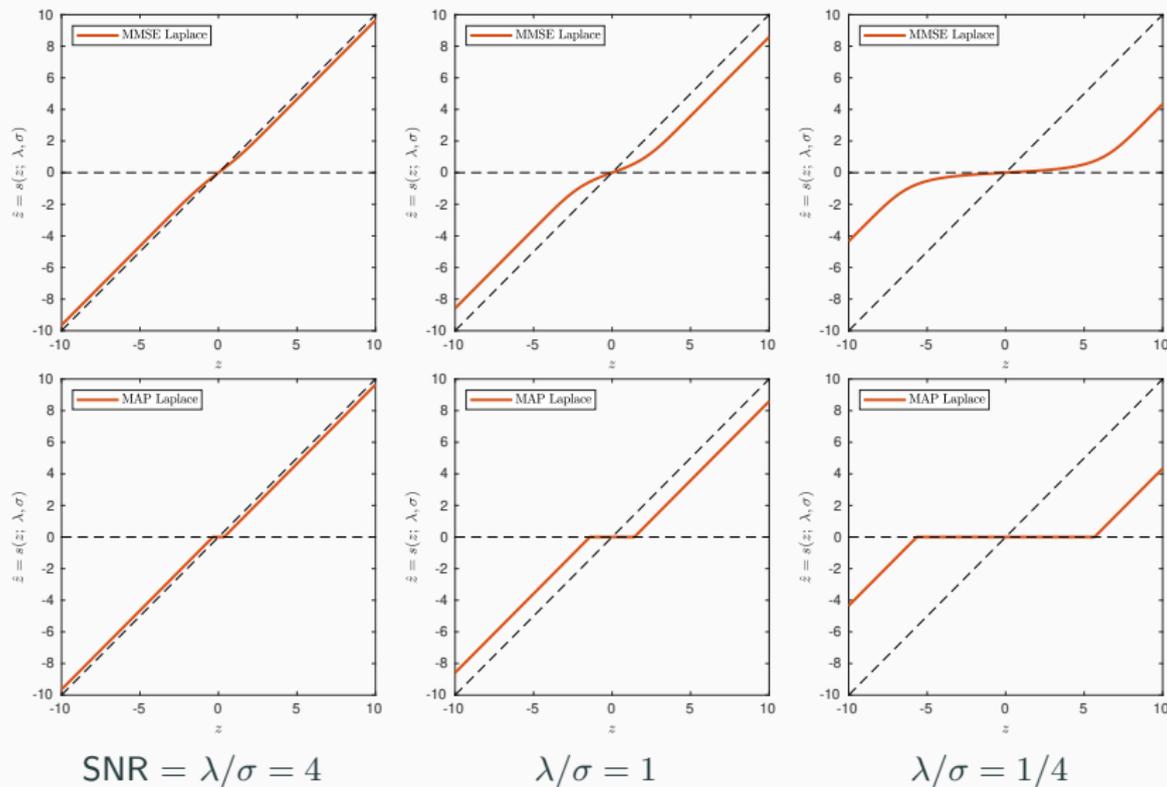
$$= z - \frac{\gamma \left(\operatorname{erf}\left(\frac{z-\gamma}{\sqrt{2}\sigma}\right) - \exp\left(\frac{2\gamma z}{\sigma^2}\right) \operatorname{erfc}\left(\frac{\gamma+z}{\sqrt{2}\sigma}\right) + 1 \right)}{\operatorname{erf}\left(\frac{z-\gamma}{\sqrt{2}\sigma}\right) + \exp\left(\frac{2\gamma z}{\sigma^2}\right) \operatorname{erfc}\left(\frac{\gamma+z}{\sqrt{2}\sigma}\right) + 1}, \quad \gamma = \frac{\sqrt{2}\sigma^2}{\lambda}$$

- MAP Shrinkage (soft-thresholding)

$$s(z; \lambda, \sigma) = \operatorname{argmin}_{\tilde{z} \in \mathbb{R}} \left[\frac{(z - \tilde{z})^2}{2\sigma^2} + \frac{\sqrt{2}|\tilde{z}|}{\lambda} \right] = \underbrace{\begin{cases} 0 & \text{if } |z| < \gamma \\ z - \gamma & \text{if } z > \gamma \\ z + \gamma & \text{if } z < -\gamma \end{cases}}_{\text{Soft-T}(z, \gamma)}$$

Non-gaussian prior: MAP \neq MMSE \rightarrow **Non-linear shrinkage.**

Gaussian noise + Laplacian prior



Example (Gaussian noise + Student prior)

- MMSE Shrinkage

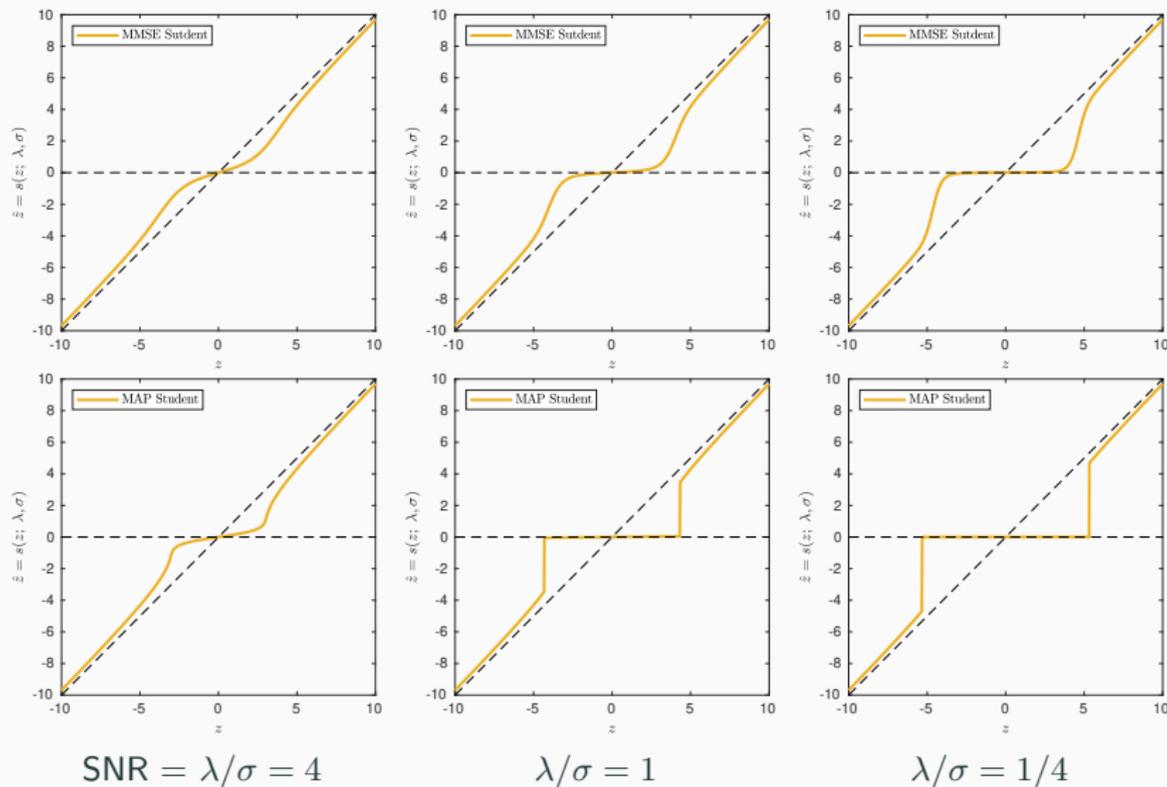
No simple expression, requires 1d numerical integration

- MAP Shrinkage

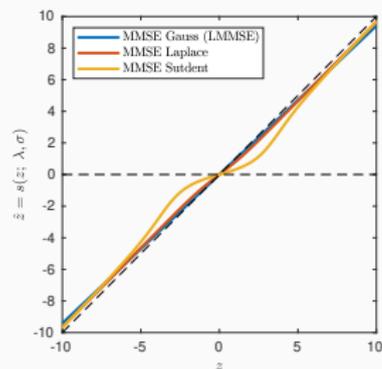
No simple expression, requires 1d numerical optimization

**For efficiency, the 1d functions
can be evaluated offline and stored in a look-up-table.**

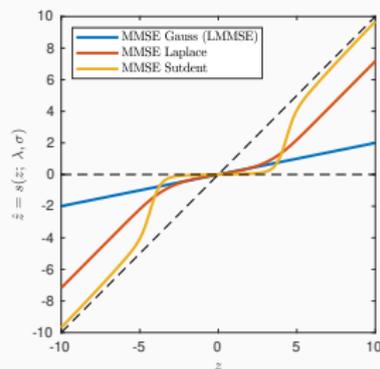
Gaussian noise + Student prior



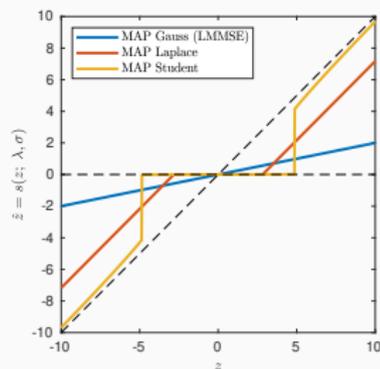
Posterior mean – Shrinkage functions – Examples



SNR = $\lambda/\sigma = 4$



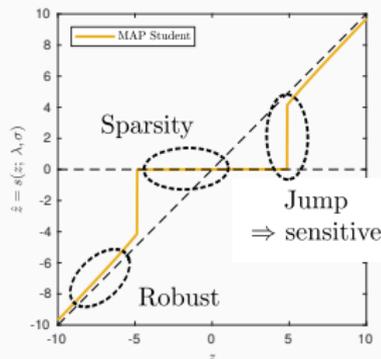
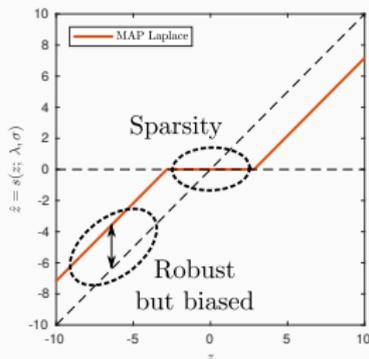
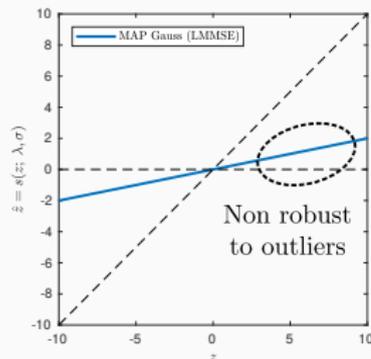
$\lambda/\sigma = 1/2$



$\lambda/\sigma = 1/2$

- Coefficients are **shrunk** towards zero
- Signs are preserved
- Non-Gaussian priors leads to non-linear filtering:
 - **sparsity**: small coefficients are shrunk (likely due to noise)
 - **robustness**: large coefficients are preserved (likely encoding signal)
- Larger SNR = $\frac{\lambda}{\sigma} \Rightarrow$ shrinkage becomes close to identity.

Interpretation



Sparsity: zero for small values.

Robustness: remain close to the identity for large values.

Transition: bias/variance tradeoff.

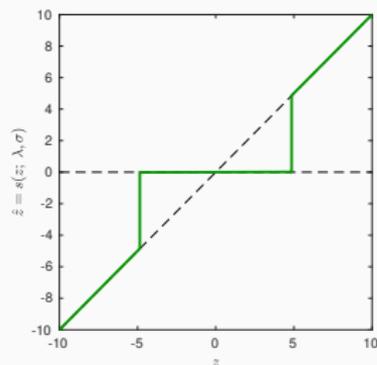
Can we design our own shrinkage according to what we want?

Shrinkage functions (*a.k.a.*, thresholding functions)

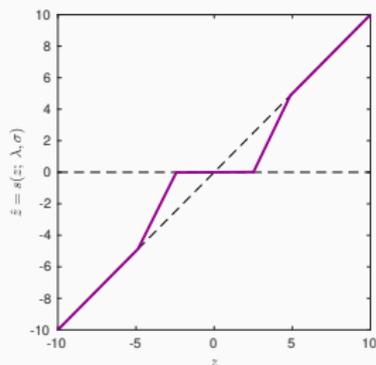
- Pick a shrinkage function s satisfying
 - **Shrink:** $|s(z)| \leq |z|$ (non-expansive)
 - **Preserve sign:** $z \cdot s(z) \geq 0$
 - **Kill low SNR:** $\lim_{\frac{\lambda}{\sigma} \rightarrow 0} s(z; \lambda, \sigma) = 0$
 - **Keep high SNR:** $\lim_{\frac{\lambda}{\sigma} \rightarrow \infty} s(z; \lambda, \sigma) = z$
 - **Increasing:** $z_1 \leq z_2 \Leftrightarrow s(z_1) \leq s(z_2)$
- **Beyond Bayesian:** No need to relate s to a prior distribution p_η .

Shrinkage functions

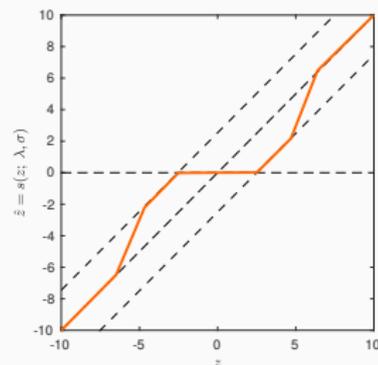
A few examples (among many others)



Hard-thresholding



MCP



SCAD

- Though not necessarily related to a prior distribution,
- Often related to a **penalized least square problem**, ex:

$$\text{Hard-T}(z) = \underset{\tilde{z} \in \mathbb{R}}{\text{argmin}} \left[(z - \tilde{z})^2 + \tau^2 \mathbf{1}_{\{\tilde{z} \neq 0\}} \right] = \begin{cases} 0 & \text{if } |z| < \tau \\ z & \text{otherwise} \end{cases}$$

- Hard-thresholding: similar behavior to Student's shrinkage.

Link with penalized least square (1/2)

- $D = L^{1/2} = E\Lambda^{1/2}$ is an **orthogonal dictionary** of n atoms/words

$$D = (d_1, d_2, \dots, d_n) \quad \text{with} \quad \|d_i\| = \lambda_i \quad \text{and} \quad \langle d_i, d_j \rangle = 0 \quad (\text{for } i \neq j)$$

- Goal: Look for the n coefficients η_i , such that \hat{x} close to y

$$\hat{x} = D\eta = \sum_{i=1}^n \eta_i d_i = \text{“linear comb. of the orthogonal atoms } d_i \text{ of } D\text{”}$$

- Choosing $\eta_i = \left\langle \frac{d_i}{\|d_i\|^2}, y \right\rangle$, i.e., $\eta = \Lambda^{-1/2} E^* y$, is optimal:

$$\hat{x} = y$$

but, it also reconstructs the noise component.

- Idea: **penalize the coeffs to prevent from reconstructing the noise.**

Link with penalized least square (2/2)

- Penalization on the coefficients controls shrinkage and sparsity:

- $\frac{1}{2}\|y - D\eta\|_2^2 + \frac{\tau^2}{2}\|\eta\|_2^2 \Rightarrow \hat{z}_i = \frac{\lambda_i^2}{\lambda_i^2 + \tau^2} z_i$

- $\frac{1}{2}\|y - D\eta\|_2^2 + \tau\|\eta\|_1 \Rightarrow \hat{z}_i = \text{Soft-T}(z_i, \gamma_i) \quad \text{with} \quad \gamma_i = \frac{\tau}{\lambda_i}$

- $\frac{1}{2}\|y - D\eta\|_2^2 + \frac{\tau^2}{2}\|\eta\|_0 \Rightarrow \hat{z}_i = \text{Hard-T}(z_i, \gamma_i) \quad \text{with} \quad \gamma_i = \frac{\tau}{\lambda_i}$

ℓ_0 pseudo-norm: $\|\eta\|_0 = \lim_{p \rightarrow 0} \left(\sum_{i=1}^n \eta_i^p \right)^{1/p} = \text{"\# of non-zero coefficients"}$

Sparsity: $\|\eta\|_0$ small compared to n

Shrinkage in the discrete Fourier domain



(a) x

(b) $y = x + w$

```
sig      = 20
► y      = x + sig * np.random.randn(x.shape)
n1, n2  = y.shape[:2]
n        = n1 * n2
lbd      = np.sqrt(prior_mpsd(n1, n2) / n)

z        = nf.fft2(y, axes=(0, 1)) / np.sqrt(n)
zhat     = shrink(z, lbd, sig)
xhat     = np.real(nf.ifft2(zhat, axes=(0, 1))) * np.sqrt(n)
```

$$z = \mathbf{F}y / \sqrt{n}$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

$$\hat{x} = \sqrt{n} \mathbf{F}^{-1} \hat{z}$$

Shrinkage in the discrete Fourier domain



(a) x



(b) λ

```
sig      = 20
y        = x + sig * np.random.randn(x.shape)
n1, n2  = y.shape[:2]
n        = n1 * n2
▶ lbd    = np.sqrt(prior_mpsd(n1, n2) / n)

z        = nf.fft2(y, axes=(0, 1)) / np.sqrt(n)
zhat     = shrink(z, lbd, sig)
xhat     = np.real(nf.ifft2(zhat, axes=(0, 1))) * np.sqrt(n)
```

$$z = \mathbf{F}y / \sqrt{n}$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

$$\hat{x} = \sqrt{n} \mathbf{F}^{-1} \hat{z}$$

Shrinkage in the discrete Fourier domain



(a) x

(b) z

```
sig    = 20
y      = x + sig * np.random.randn(x.shape)
n1, n2 = y.shape[:2]
n      = n1 * n2
lbd    = np.sqrt(prior_mpsd(n1, n2) / n)
```

```
► z     = nf.fft2(y, axes=(0, 1)) / np.sqrt(n)
zhat   = shrink(z, lbd, sig)
xhat   = np.real(nf.ifft2(zhat, axes=(0, 1))) * np.sqrt(n)
```

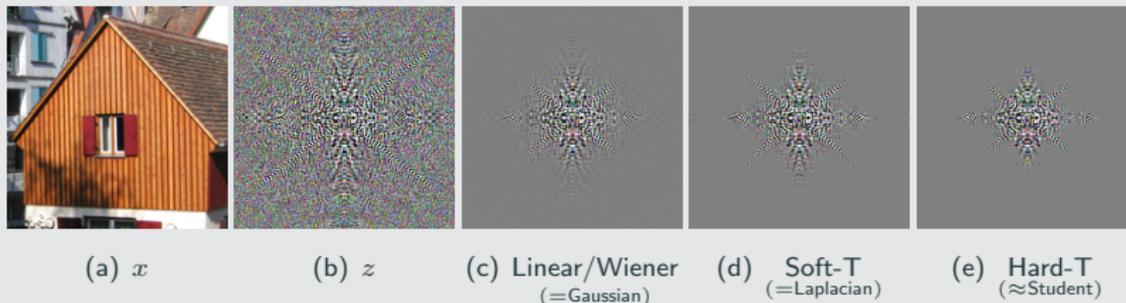
$$z = \mathbf{F}y/\sqrt{n}$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

$$\hat{x} = \sqrt{n}\mathbf{F}^{-1}\hat{z}$$

Posterior mean – Shrinkage in the Fourier domain

Shrinkage in the discrete Fourier domain



```
sig    = 20
y      = x + sig * np.random.randn(x.shape)
n1, n2 = y.shape[:2]
n      = n1 * n2
lbd    = np.sqrt(prior_mpsd(n1, n2) / n)

z      = nf.fft2(y, axes=(0, 1)) / np.sqrt(n)
▶ zhat = shrink(z, lbd, sig)
xhat   = np.real(nf.ifft2(zhat, axes=(0, 1))) * np.sqrt(n)
```

$$z = \mathbf{F}y/\sqrt{n}$$
$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$
$$\hat{x} = \sqrt{n}\mathbf{F}^{-1}\hat{z}$$

Posterior mean – Shrinkage in the Fourier domain

Shrinkage in the discrete Fourier domain

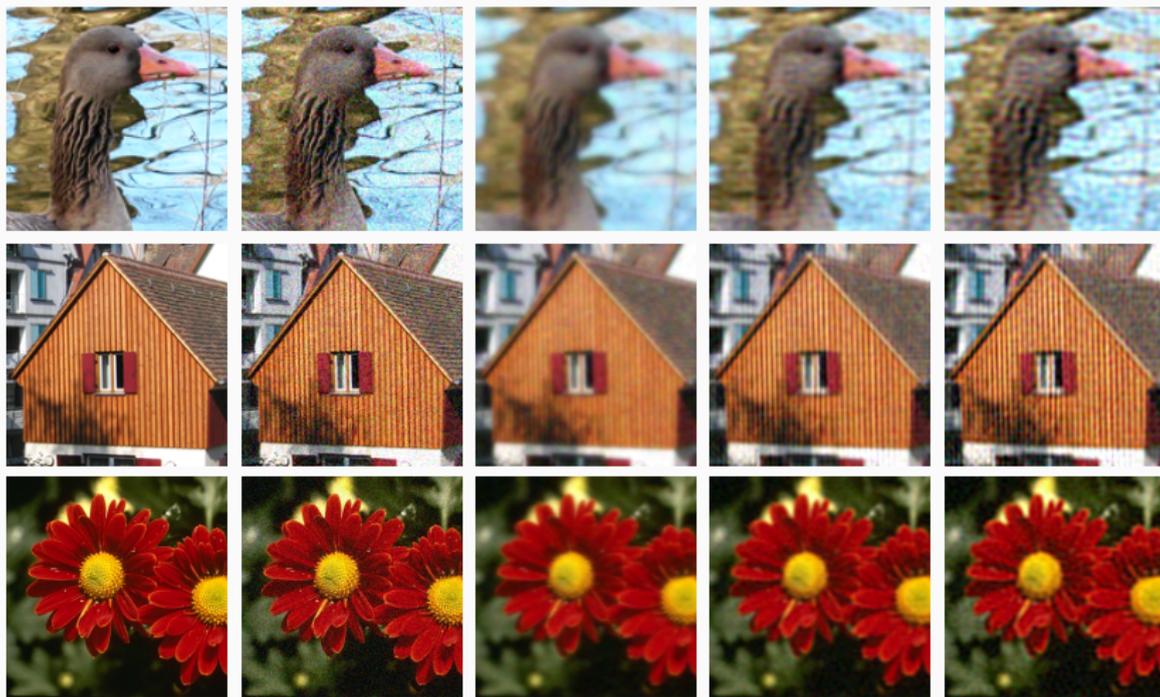


```
sig      = 20
y        = x + sig * np.random.randn(x.shape)
n1, n2  = y.shape[:2]
n        = n1 * n2
lbd      = np.sqrt(prior_mpsd(n1, n2) / n)

z        = nf.fft2(y, axes=(0, 1)) / np.sqrt(n)
zhat     = shrink(z, lbd, sig)
▶ xhat   = np.real(nf.ifft2(zhat, axes=(0, 1))) * np.sqrt(n)
```

$$z = \mathbf{F}y/\sqrt{n}$$
$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$
$$\hat{x} = \sqrt{n}\mathbf{F}^{-1}\hat{z}$$

Shrinkage functions – Fourier domain – Results



(a) x

(b) y ($\sigma = 20$)

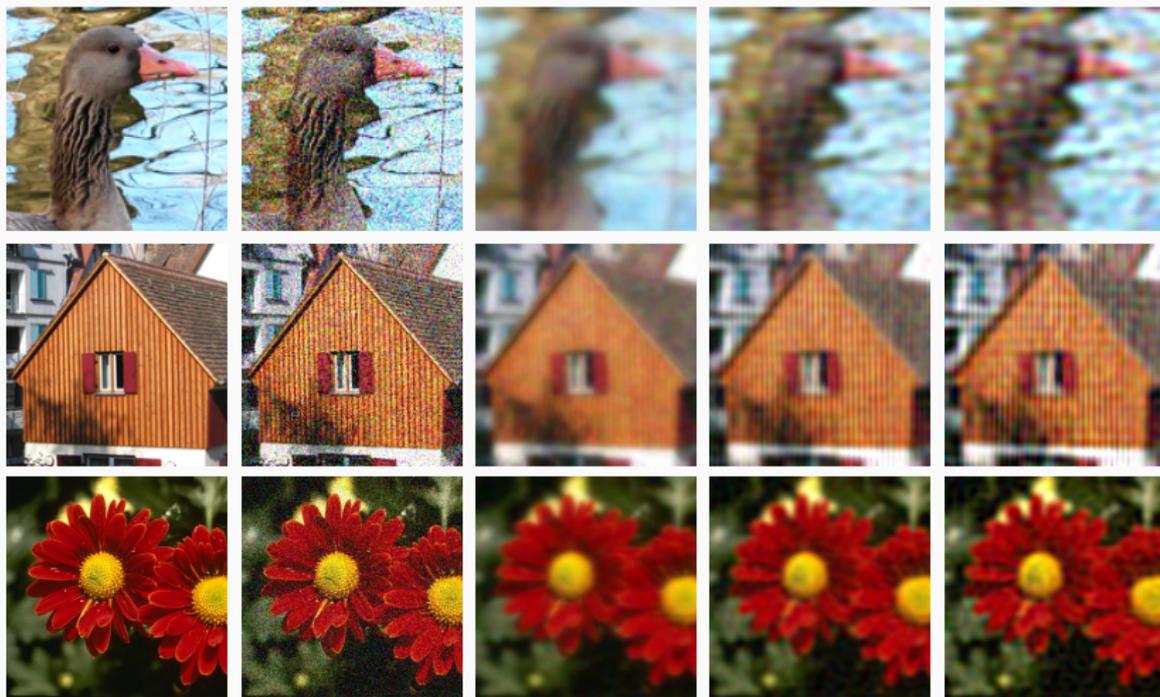
(c) Linear/Wiener
(=Gaussian)

(d) Soft-T
(=Laplacian)

(e) Hard-T
(\approx Student)

Bias \longleftrightarrow Variance

Shrinkage functions – Fourier domain – Results



(a) x

(b) y ($\sigma = 40$)

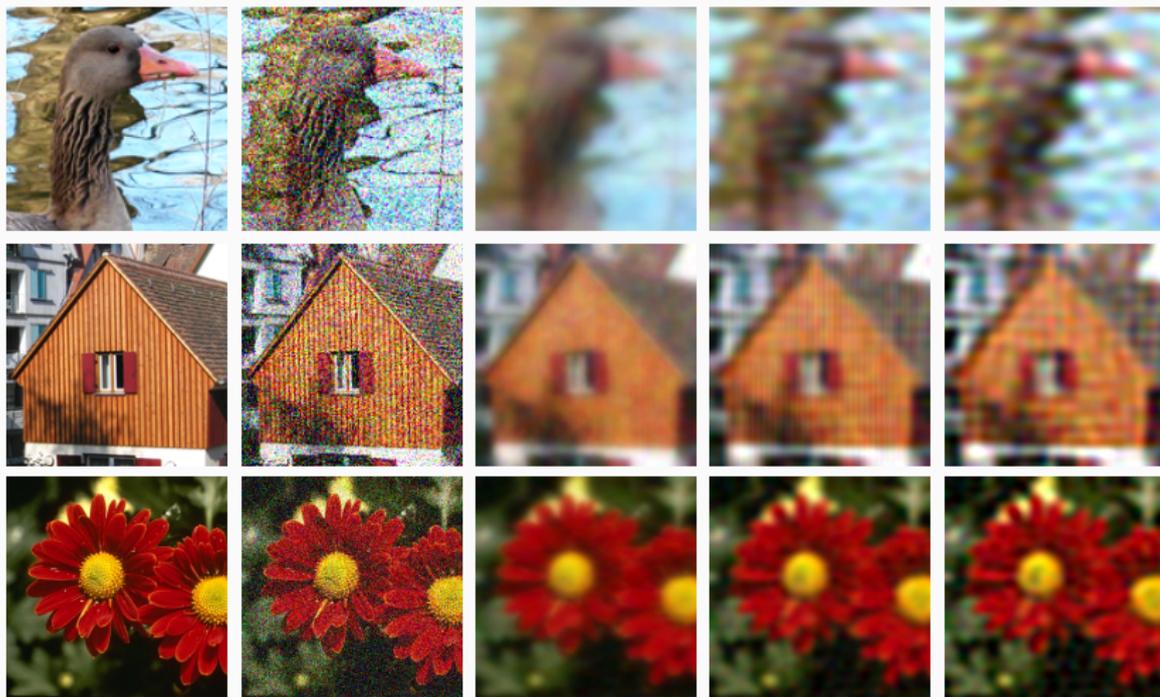
(c) Linear/Wiener
(=Gaussian)

(d) Soft-T
(=Laplacian)

(e) Hard-T
(\approx Student)

Bias \longleftrightarrow Variance

Shrinkage functions – Fourier domain – Results



(a) x

(b) y ($\sigma = 60$)

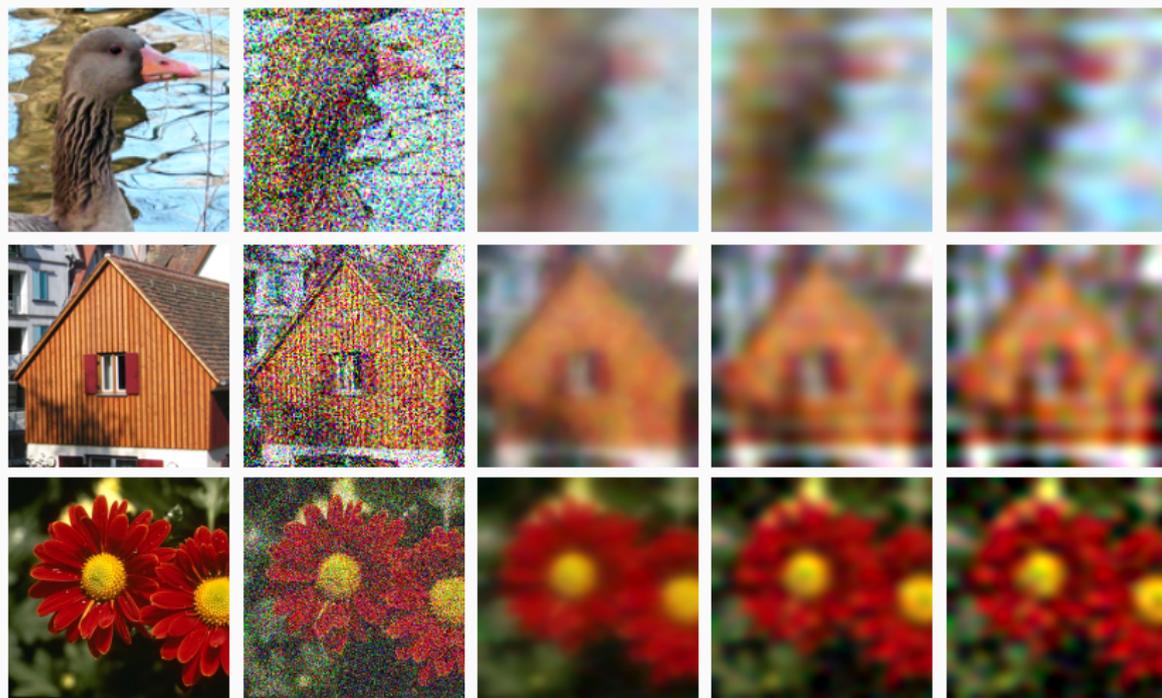
(c) Linear/Wiener
(=Gaussian)

(d) Soft-T
(=Laplacian)

(e) Hard-T
(\approx Student)

Bias \longleftrightarrow Variance

Shrinkage functions – Fourier domain – Results



(a) x

(b) y ($\sigma = 120$)

(c) Linear/Wiener
(=Gaussian)

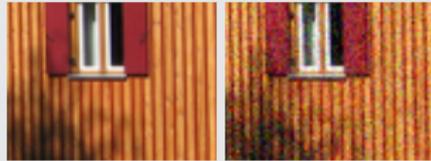
(d) Soft-T
(=Laplacian)

(e) Hard-T
(\approx Student)

Bias \longleftrightarrow Variance

Posterior mean – Limits of shrinkage in the Fourier domain

Limits of shrinkage in the discrete Fourier domain

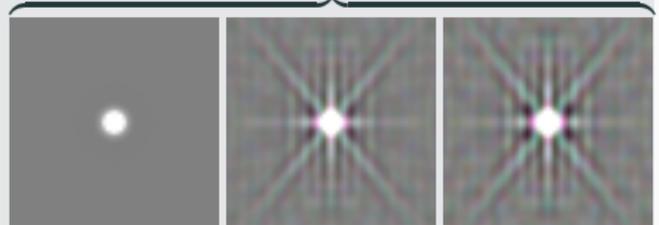


(a) x

(b) y



convolution kernels



(c) Linear
(=Gaussian)

(d) Soft-T
(=Laplacian)

(e) Hard-T
(\approx Student)

- Linear shrinkage (Wiener)
⇒ Non-adaptive,
- Non-linear shrinkage
⇒ **Adaptive convolution**,
- Adapts to the frequency content,
- but not to the spatial content.

$$\hat{z}_i = s(z_i; \tau, \sigma) = \underbrace{\frac{s(z_i; \tau, \sigma)}{z_i}}_{\text{element-wise product}} \times z_i$$

\Leftrightarrow

$$\hat{x} = \underbrace{\nu(y)}_{\substack{\text{spatial average} \\ \text{adapted to the spectrum of } y}} * y$$

Motivations

Consequences

- Modulating Fourier coefficients \Rightarrow Non spatially adaptive
- Assuming Fourier coefficients to be white+sparse \Rightarrow Shrinkage in Fourier

Deductive reasoning

Need another representation for sparsifying clean images

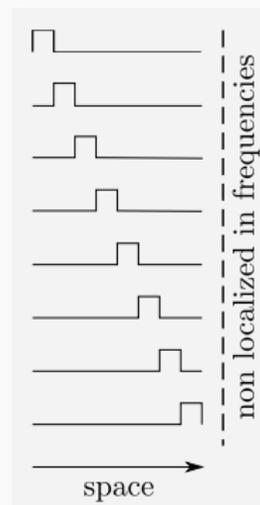
$$\mathbf{E} \neq \frac{1}{\sqrt{n}} \left(\begin{array}{c} \text{[Grid of patterns]} \\ \times n \end{array} \right) \leftarrow \text{Columns were the Fourier atoms}$$

DFT: \mathbf{F}

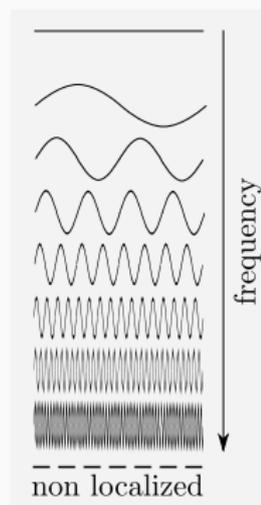
What transform can make signal white and sparse and captures both spatial and spectral contents?

Wavelet transforms

Canonical basis

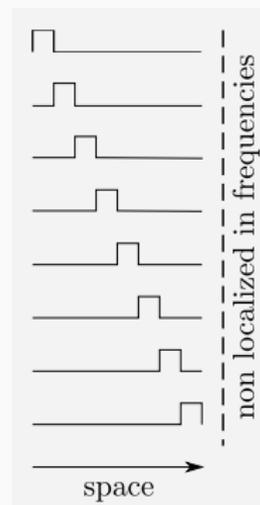


Fourier basis

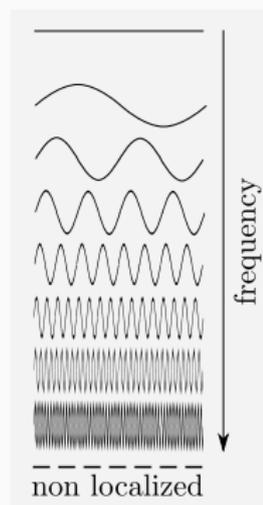


$$\text{Id} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{F} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & e^{-2\pi i/4} & e^{-2\pi i 2/4} & e^{-2\pi i 3/4} \\ 1 & e^{-2\pi i 2/4} & e^{-2\pi i 4/4} & e^{-2\pi i 6/4} \\ 1 & e^{-2\pi i 3/4} & e^{-2\pi i 6/4} & e^{-2\pi i 9/4} \end{pmatrix}$$

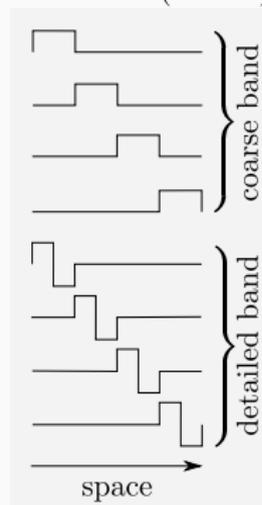
Canonical basis



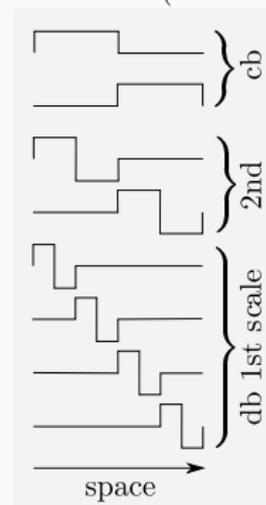
Fourier basis



Haar basis (1 scale)



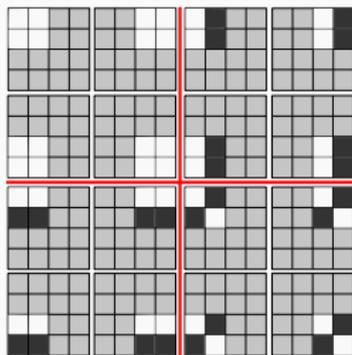
Haar basis (2 scales)



$$\mathcal{H}^{1st} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}$$

$$\text{and } \mathcal{H}^{2nd} = \begin{matrix} 1/2 \\ 1/2 \\ 1/\sqrt{2} \\ 1/\sqrt{2} \end{matrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}$$

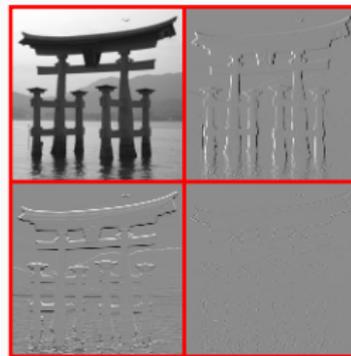
Introduction to wavelets – Haar (2d case)



(a) \mathcal{H}^{1st} (4×4 image)



(b) x

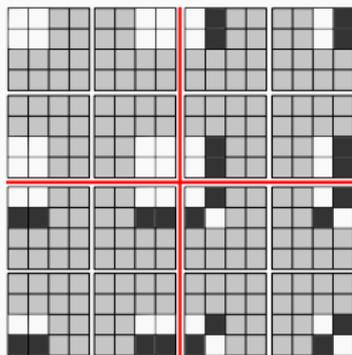


(c) $\mathcal{H}^{1st}x$

2d Haar representation

- 4 sub-bands
- Coarse sub-band
 - Vertical detailed sub-band
 - Horizontal detailed sub-band
 - Diagonal detailed sub-band

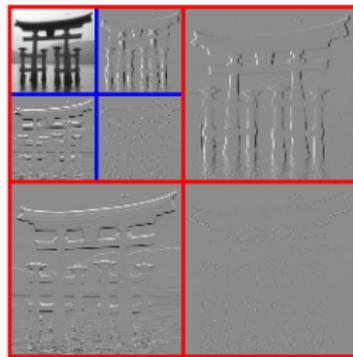
Introduction to wavelets – Haar (2d case)



(a) \mathcal{H}^{1st} (4×4 image)



(b) x

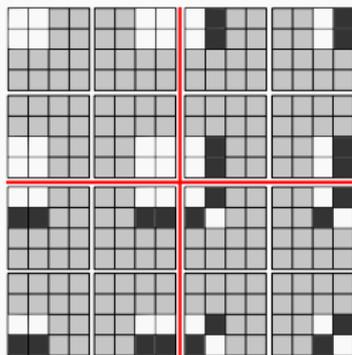


(c) $\mathcal{H}^{2nd} x$

Multi-scale 2d Haar representation

- Repeat recursively J times
- Dyadic decomposition
- Multi-scale representation
- Related to scale spaces

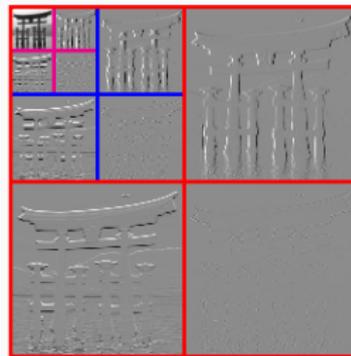
Introduction to wavelets – Haar (2d case)



(a) \mathcal{H}^{1st} (4×4 image)



(b) x

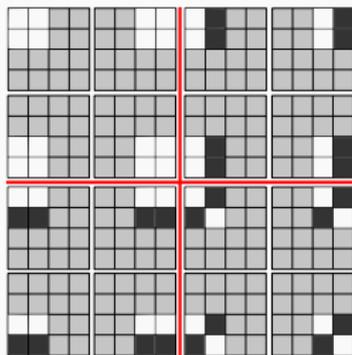


(c) $\mathcal{H}^{3rd} x$

Multi-scale 2d Haar representation

- Repeat recursively J times
- Dyadic decomposition
- Multi-scale representation
- Related to scale spaces

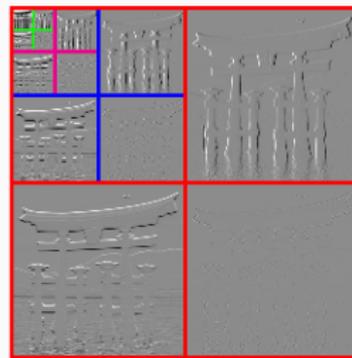
Introduction to wavelets – Haar (2d case)



(a) \mathcal{H}^{1st} (4×4 image)



(b) x

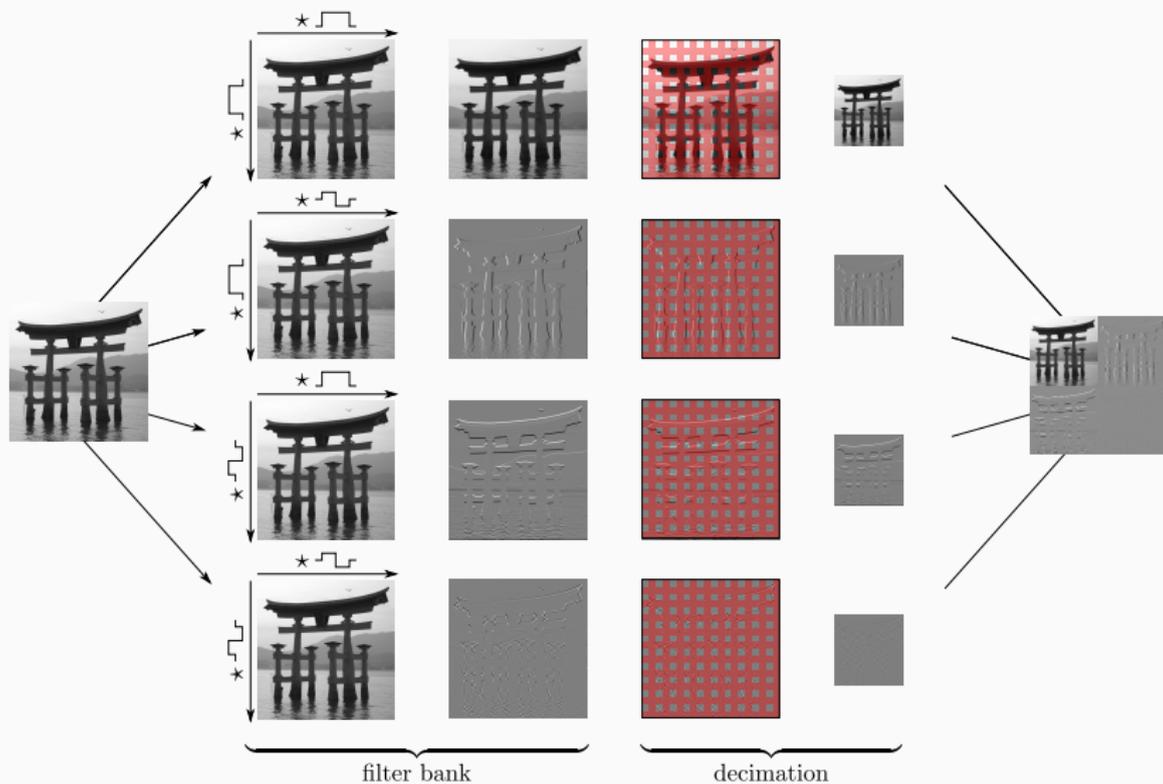


(c) $\mathcal{H}^{4th} x$

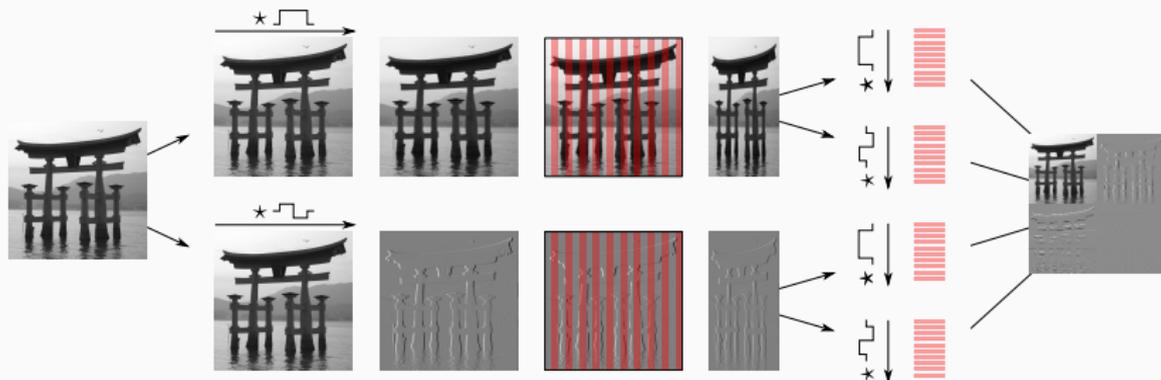
Multi-scale 2d Haar representation

- Repeat recursively J times
- Dyadic decomposition
- Multi-scale representation
- Related to scale spaces

Introduction to wavelets – Haar transform - Filter bank



Introduction to wavelets – Haar transform - Separability



Properties of the 2d Haar transform

- Separable: 1d Haar transforms in horizontal and next vertical direction
- First: perform a **low pass** and **high pass** filtering
- Next: perform decimation by a factor of 2

Can we choose other low and high pass filters
to get a better transform?

Discrete wavelet transform (DWT) (1/3)

(1d and n even)

- Let $h \in \mathbb{R}^n$ (with periodical boundary conditions) satisfying

$$\sum_{i=0}^{n-1} h_i = 0$$

$$\sum_{i=0}^{n-1} h_i^2 = 1$$

and $\sum_{i=0}^{n-1} h_i h_{i+2k} = 0$ for all integer $k \neq 0$

Example (Haar as a particular case)

$$h = \frac{1}{\sqrt{2}}(0 \dots 0 \quad -1 \quad +1 \quad 0 \dots 0)$$

Discrete wavelet transform (DWT) (2/3)

(1d and n even)

- Define the high and low pass filters $H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $G : \mathbb{R}^n \rightarrow \mathbb{R}^n$ as

$$(Hx)_k = (h * x)_k = \sum_{i=0}^{n-1} h_i x_{k-i}$$

$$(Gx)_k = (g * x)_k = \sum_{i=0}^{n-1} g_i x_{k-i} \quad \text{where } g_i = (-1)^i h_{n-1-i}$$

- Note: necessarily $\sum_{i=0}^{n-1} g_i = \sqrt{2}$

Example (Haar as a particular case)

$$h = \frac{1}{\sqrt{2}}(0 \dots 0 \quad -1 \quad +1 \quad 0 \dots 0)$$

$$g = \frac{1}{\sqrt{2}}(0 \dots 0 \quad +1 \quad +1 \quad 0 \dots 0)$$

- Define the decimation by 2 of a matrix $M \in \mathbb{R}^{n \times n}$ as

$$M \downarrow_2 = "M[:, :2, :]" \in \mathbb{R}^{n/2 \times n}$$

i.e., the matrix obtained by removing every two rows.

- $M \downarrow_2 x$: apply M to x and next remove every two entries.

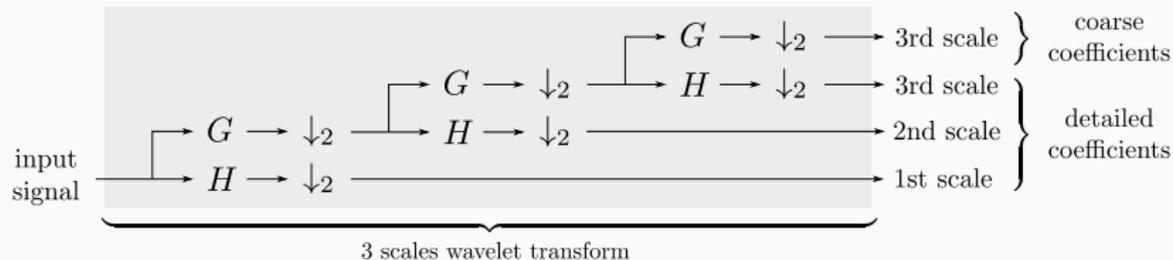
Discrete wavelet transform (DWT) (3/3)

(1d and n even)

$$\text{Let } W = \begin{pmatrix} G \downarrow_2 \\ H \downarrow_2 \end{pmatrix} \in \mathbb{R}^{n \times n}$$

Then $\left\{ \begin{array}{ll} \bullet x \mapsto Wx: & \text{orthonormal discrete wavelet transform,} \\ \bullet \text{Columns of } W: & \text{orthonormal discrete wavelet basis,} \\ \bullet z = Wx: & \text{wavelet coefficients of } x. \end{array} \right.$

Multi-scale discrete wavelets



Multi-scale DWT (1d and n multiple of 2^J)

[Mallat, 1989]

Defined recursively as
$$\mathbf{W}^{J\text{-th}} = \begin{pmatrix} \mathbf{W}^{(J-1)\text{-th}} & \mathbf{O} \\ \mathbf{0} & \text{Id} \end{pmatrix} \mathbf{W}$$

Implementation of 2D DWT

(n_1 and n_2 multiple of 2^J)

```
def dwt(x, J, h, g):                                     # 2d and multi-scale
    if J == 0:
        return x
    n1, n2 = x.shape[:2]
    m1, m2 = (int(n1 / 2), int(n2 / 2))
    z = dwt1d(x, h, g)
    z = flip(dwt1d(flip(z), h, g))
    z[:m1, :m2] = dwt(z[:m1, :m2], J - 1, h, g)
    return z
```

```
def dwt1d(x, h, g):                                     # 1d and 1scale
    coarse = convolve(x, g)
    detail = convolve(x, h)
    z = np.concatenate((coarse[:,2:], detail[:,2:]), axis=0)
    return z
```

Use separability

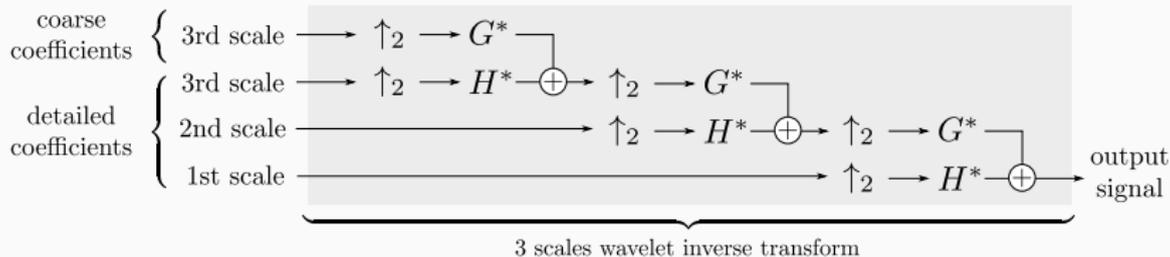
Multi-scale Inverse DWT (1d and n multiple of 2^J)

Defined recursively as $(W^{J\text{-th}})^{-1} = W^{-1} \begin{pmatrix} (W^{(J-1)\text{-th}})^{-1} & O \\ 0 & \text{Id} \end{pmatrix}$

where $W^{-1} = W^* = \begin{pmatrix} G^* \uparrow_2 & H^* \uparrow_2 \end{pmatrix} \in \mathbb{R}^{n \times n}$

and $M \uparrow_2$: remove every two columns.

$M \uparrow_2 x$: insert 0 every two entries in x and next apply M .



Implementation of 2D IDWT

(n_1 and n_2 multiple of 2^J)

```
def idwt(z, J, h, g):                                # 2d and multi-scale
    if J == 0:
        return z
    n1, n2 = z.shape[:2]
    m1, m2 = (int(n1 / 2), int(n2 / 2))
    x = z.copy()
    x[:m1, :m2] = idwt(x[:m1, :m2], J - 1, h, g)
    x = flip(idwt1d(flip(x), h, g))
    x = idwt1d(x, h, g)
    return x
```

```
def idwt1d(z, h, g):                                # 1d and 1scale
    n1 = z.shape[0]
    m1 = int(n1 / 2)
    coarse, detail = np.zeros(z.shape), np.zeros(z.shape)
    coarse[::2, :], detail[::2, :] = z[:m1, :], z[m1:, :]
    x = convolve(coarse, g[:::-1]) + convolve(detail, h[:::-1])
    return x
```

Discrete wavelet with limited support

- Consider a high pass filter with **finite support** of size $m = 2p$ (even). For instance for $m = 4$

$$H = \begin{pmatrix} h_2 & h_3 & 0 & & & \dots & & & 0 & h_0 & h_1 & & \\ h_0 & h_1 & h_2 & h_3 & 0 & & \dots & & & & & & 0 \\ & 0 & h_0 & h_1 & h_2 & h_3 & 0 & & \dots & & & & \\ & & & & & & \ddots & & & & & & \\ & & & & & & & & & & & & \\ 0 & & & \dots & & & & & 0 & h_0 & h_1 & h_2 & h_3 \\ h_2 & h_3 & 0 & & & \dots & & & & 0 & h_0 & h_1 & \end{pmatrix}$$

- Then h defines a wavelet transform if it satisfies the three conditions

$$\sum h_i = 0 \quad \text{and} \quad \sum h_i^2 = 1 \quad \text{and} \quad \sum h_i h_{i+2k} = 0 \quad \text{for } k = 1 \text{ to } p - 1$$

- This system has $2p$ unknowns and $1 + p$ independent equations.
- If $p = 1$, $2p = 1 + p$, this implies that the solution is unique (Haar).
- Otherwise, one has **$p - 1$ degrees of freedom**.

Daubechies' wavelets (1988)

- Daubechies suggests adding the $p - 1$ constraints

$$\sum_{i=0}^{2p-1} i^q h_i = 0 \quad \text{for } q = 1 \text{ to } p - 1 \quad (\text{vanishing } q\text{-order moments})$$

- For $p = 2$, the (orthonormal) Daubechies' wavelets are defined as

$$\begin{cases} h_0^2 + h_1^2 + h_2^2 + h_3^2 = 1 \\ h_0 + h_1 + h_2 + h_3 = 0 \\ h_0 h_2 + h_1 h_3 = 0 \\ h_1 + 2h_2 + 3h_3 = 0 \end{cases} \Leftrightarrow h = \pm \frac{1}{\sqrt{2}} \begin{pmatrix} \frac{1+\sqrt{3}}{4} \\ \frac{3+\sqrt{3}}{4} \\ \frac{3-\sqrt{3}}{4} \\ \frac{1-\sqrt{3}}{4} \end{pmatrix}$$

- The corresponding DWT is referred to as Daubechies-2 (or Db2).

As for the Fourier transform, there also exists a continuous version.

Continuous wavelet transform (CWT)

(1d)

- Continuum of locations $t \in \mathbb{R}$ and scales $a > 0$,
- **Continuous wavelet transform** of $x : \mathbb{R} \rightarrow \mathbb{R}$

$$\underbrace{c(a, t)}_{\text{wavelet coefficient}} = \int_{-\infty}^{+\infty} \psi_{a,t}^*(t') x(t') dt' = \langle \underbrace{x}_{\text{signal}}, \underbrace{\psi_{a,t}}_{\text{wavelet}} \rangle$$

where $*$ is the complex conjugate.

- $\psi_{a,t}$: **daughter wavelets**, translated and scaled versions of Ψ

$$\psi_{a,t}(t') = \frac{1}{\sqrt{a}} \Psi \left(\frac{t' - t}{a} \right)$$

- Ψ : **the mother wavelets** satisfying

$$\int_{-\infty}^{+\infty} \Psi(t) dt = 0 \quad \text{and} \quad \int_{-\infty}^{+\infty} |\Psi(t)|^2 dt = 1 < \infty$$

(zero-mean) (unit-norm / square-integrable)

Inverse CWT

(1d)

- The **inverse continuous wavelet transform** is given by

$$x(t) = \frac{1}{C_\Psi} \int_{-\infty}^{+\infty} \int_0^{+\infty} \frac{1}{|a|^2} c(a, t') \psi_{a,t}(t') \, da \, dt'$$

with $C_\Psi = \int_0^{+\infty} \frac{|\hat{\Psi}(u)|^2}{u} \, du$ where $\hat{\Psi}$ is the Fourier transform of Ψ .

Relation between CWT/DWT

(1d)

- The DWT can be seen as the discretization of the CWT
 - **Diadic discretization** in scale: $a = 1, 2, 4, \dots, 2^J$
 - Uniform discretization in time at scale j with step 2^j : $t = 1:2^j:n$

Twin-scale relation

(1d)

- The CWT is orthogonal (inverse = adjoint), if and only if Ψ satisfies

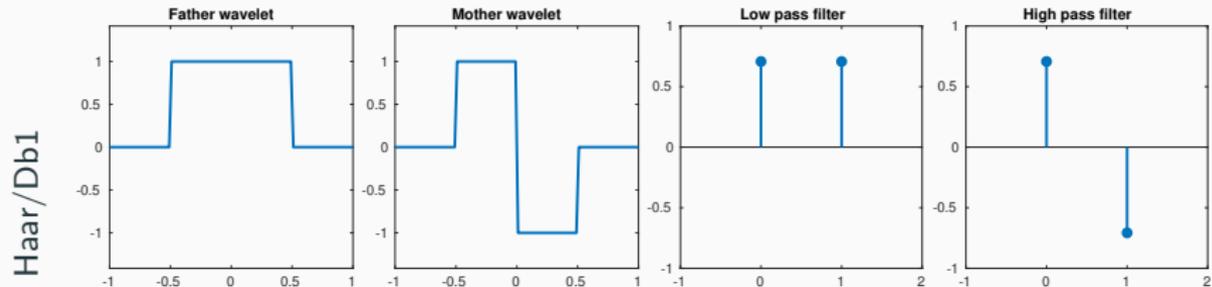
$$\Psi(t) = \sqrt{2} \sum_{i=0}^{m-1} h_i \Phi(2t - i) \quad \text{and} \quad \Phi(t) = \sqrt{2} \sum_{i=0}^{m-1} g_i \Phi(2t - i)$$

where h and g are high- and low-pass filters defining a DWT.

- Φ is called **father wavelet or scaling function**.
- Note: potentially $m = \infty$.

Twin-scale relation: allows to define a CWT from DWT and vice-versa.
The CWT may not have a closed form (approximated by the **cascade algorithm**)

Continuous and discrete wavelets



Popular wavelets are:

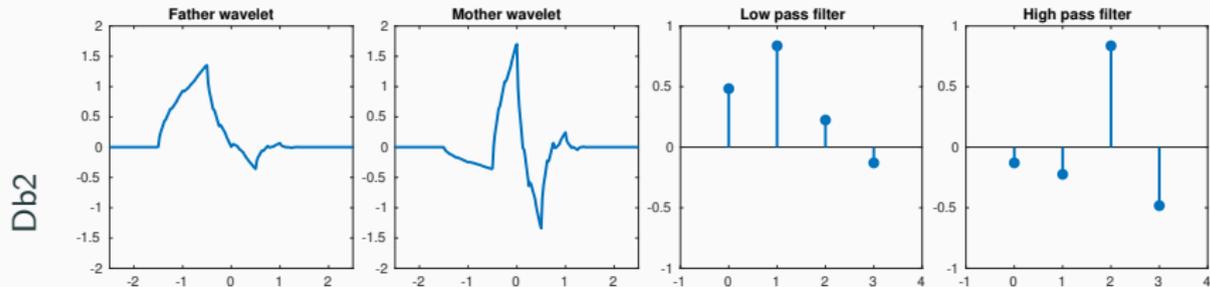
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

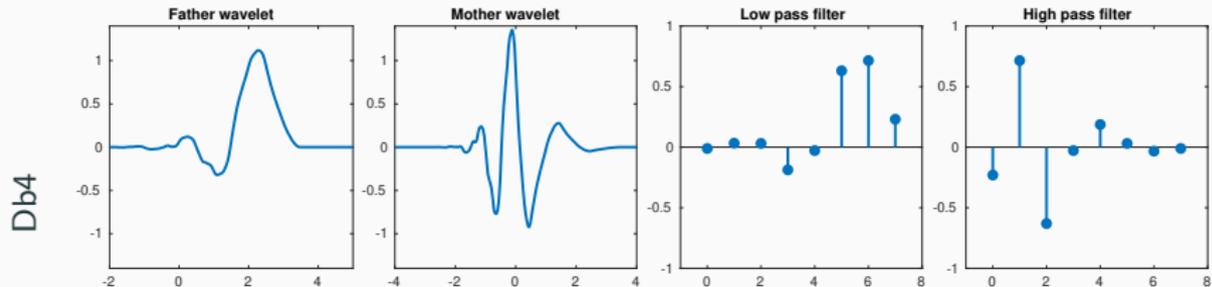
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

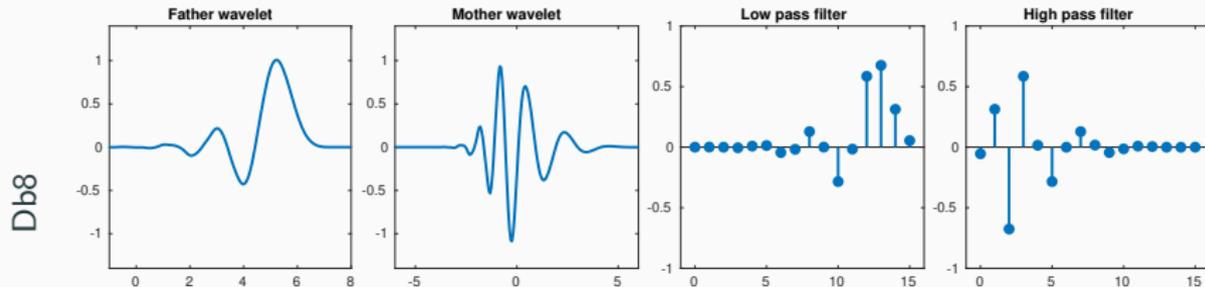
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

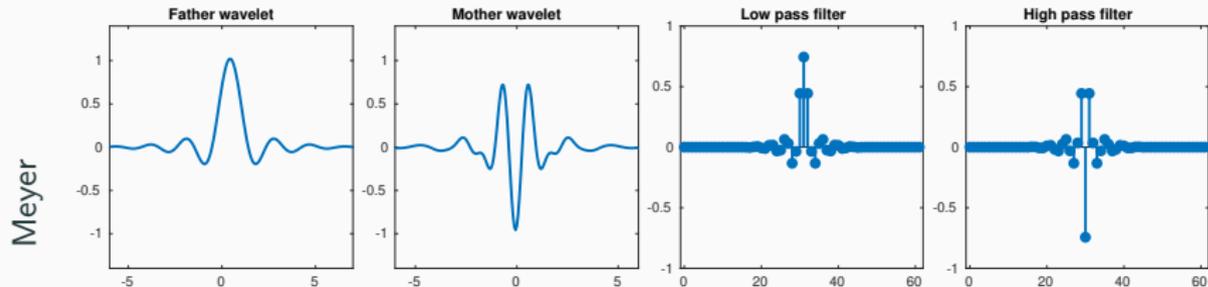
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

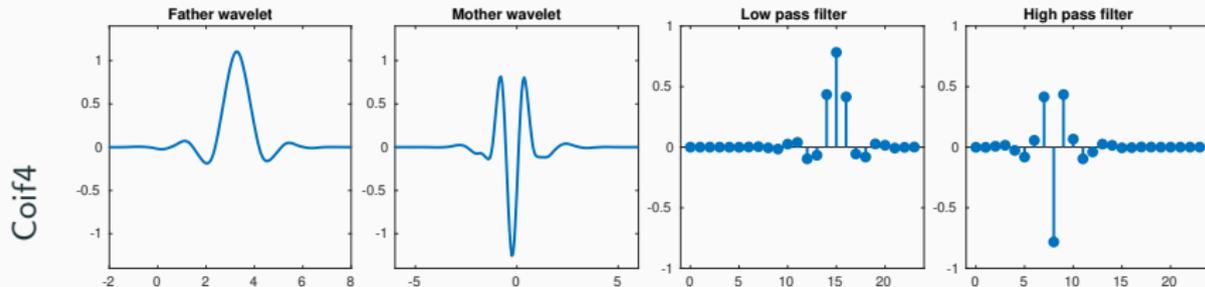
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

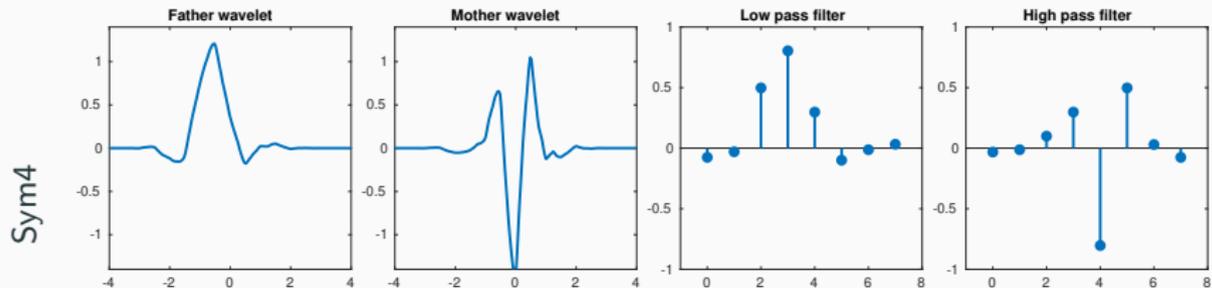
Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Continuous and discrete wavelets



Popular wavelets are:

Haar (1909); Gabor wavelet (1946); Mexican hat/Marr wavelet (1980); Morlet wavelet (1984); Daubechies (1988); Meyer wavelet (1990); Binomial quadrature mirror filter (1990); Coiflets (1991); Symlets (1992).

Some classical wavelet transforms are not orthogonal.

Bi-orthogonal: Non orthogonal but invertible (ex: for symmetric wavelets).

Two filters for the direct, and two others for the inverse.

Wavelets perform image compression

- Haar encodes constant signals with one coefficient,
- Db- p encodes $(p-1)$ -order polynomials with p coefficients.

Consequences:

- Polynomial/Smooth signals are encoded with very few coefficients,
- Coarse coefficients encode the smooth underlying signal,
- Detailed coefficients encode non-smooth content of the signal,
- Typical signals are concentrated on few coefficients,
- The remaining coefficients capture only noise components.

Wavelets perform image compression

- Haar encodes constant signals with one coefficient,
- Db- p encodes $(p-1)$ -order polynomials with p coefficients.

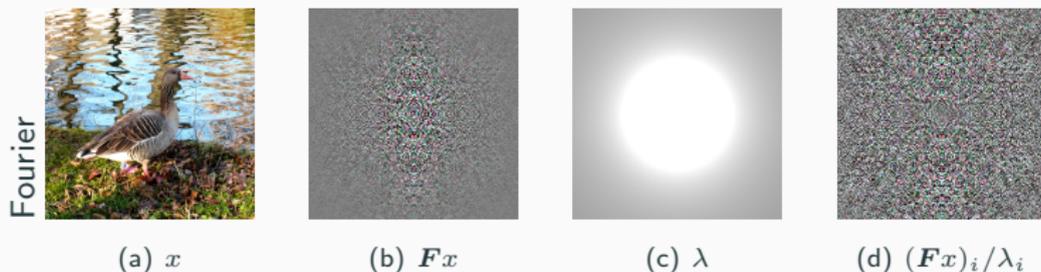
Consequences:

- Polynomial/Smooth signals are encoded with very few coefficients,
- Coarse coefficients encode the smooth underlying signal,
- Detailed coefficients encode non-smooth content of the signal,
- Typical signals are concentrated on few coefficients,
- The remaining coefficients capture only noise components.

⇒ **Heavy tailed distribution with a peak at zero,**

i.e., **wavelets favor sparsity.**

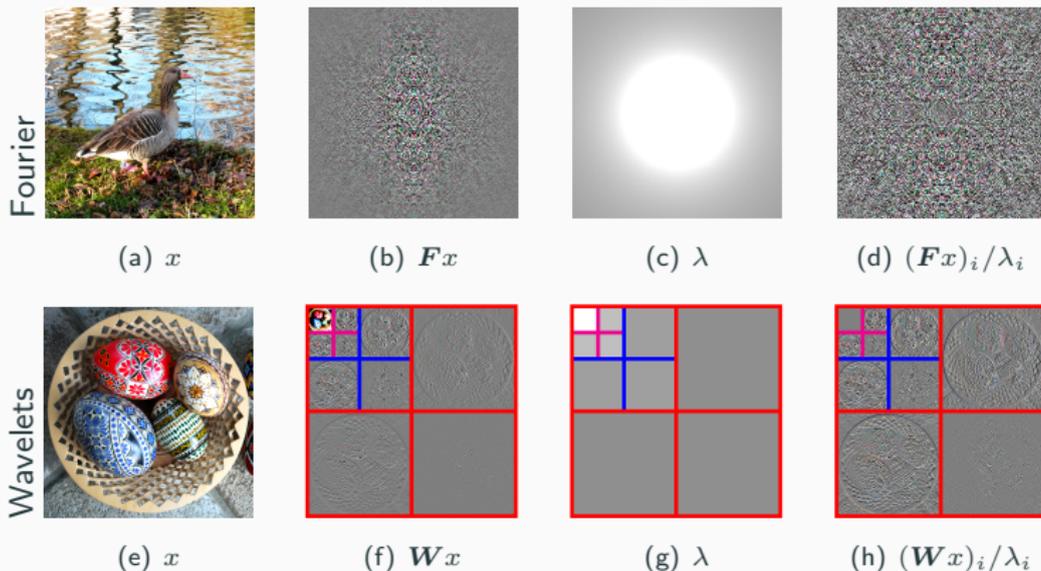
Wavelets as a sparsifying transform



Fourier (u_i, v_i freq. of component i)

- $\mathbf{E}^* = \mathbf{F}/\sqrt{n}$
- $\lambda_i^2 = n^{-1}\text{MPSD}$ and ∞ if $i = 0$
- Arbitrary DC component

Wavelets as a sparsifying transform



Fourier (u_i, v_i freq. of component i)

- $E^* = F/\sqrt{n}$
- $\lambda_i^2 = n^{-1}$ MPSD and ∞ if $i = 0$
- Arbitrary DC component

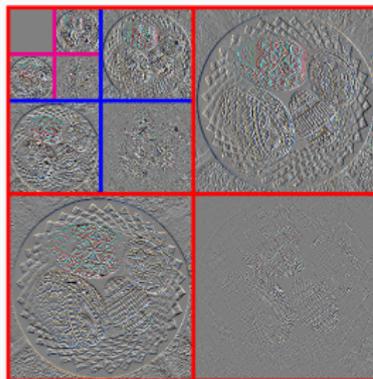
Wavelets (j_i scale of component i)

- $E^* = W$
- $\lambda_i = \alpha 2^{j_i-1}$ and ∞ if $j_i = J$
- Arbitrary coarse component

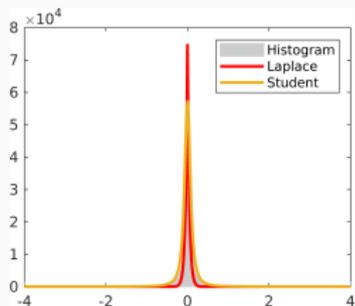
Distribution of wavelet coefficients



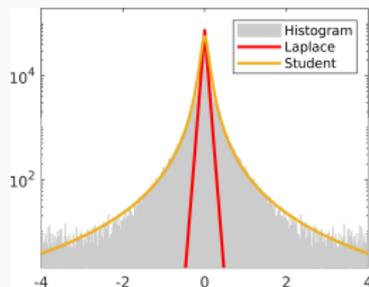
(a) x



(b) $\eta_i = (Wx)_i / \lambda_i$



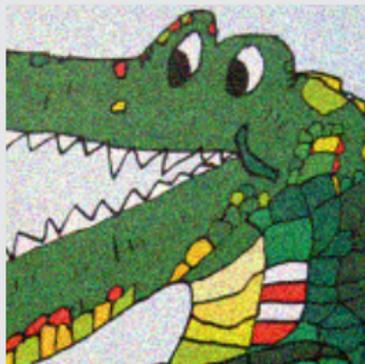
(c) Histogram of η



(d) Histogram of η

Shrinkage in the wavelet domain

Shrinkage in the discrete wavelet domain



(a) y

```
sig = 20
▶ y = x + sig * nr.randn(*x.shape)

z = im.dwt(y, 3, h, g)
zhat = shrink(z, lbd, sig)
xhat = im.idwt(zhat, 3, h, g)
```

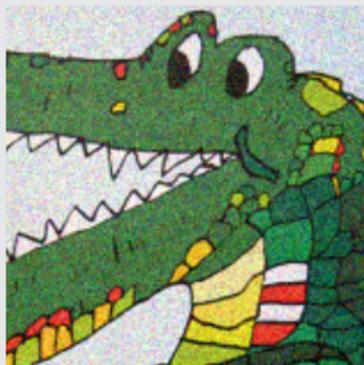
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

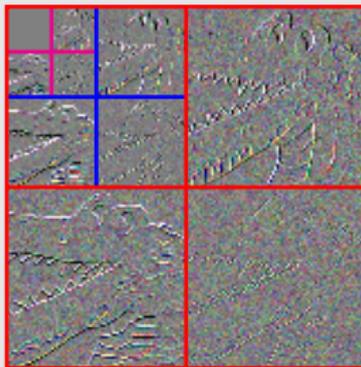
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

Shrinkage in the discrete wavelet domain



(a) y



(b) z (Haar)

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

```
► z = im.dwt(y, 3, h, g)  
zhat = shrink(z, lbd, sig)  
xhat = im.idwt(zhat, 3, h, g)
```

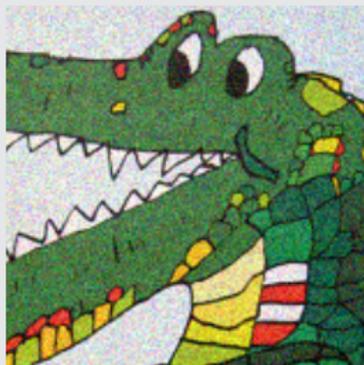
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

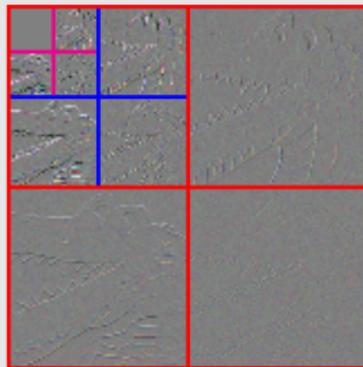
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

Shrinkage in the discrete wavelet domain



(a) y



(b) \hat{z} (Haar+LMMSE)

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

```
z = im.dwt(y, 3, h, g)  
▶ zhat = shrink(z, lbd, sig)  
xhat = im.idwt(zhat, 3, h, g)
```

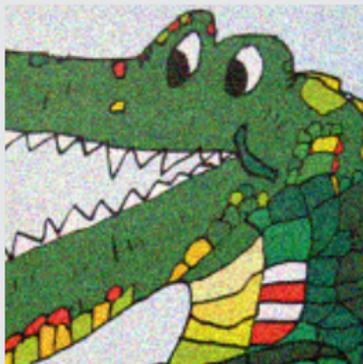
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

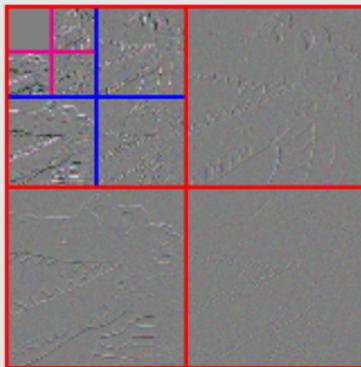
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

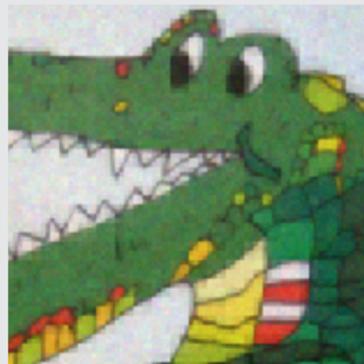
Shrinkage in the discrete wavelet domain



(a) y



(b) \hat{z} (Haar+LMMSE)



(c) \hat{x}

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

```
z = im.dwt(y, 3, h, g)  
zhat = shrink(z, lbd, sig)  
▶ xhat = im.idwt(zhat, 3, h, g)
```

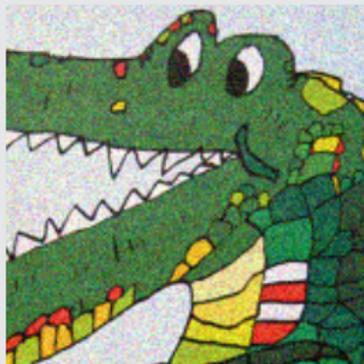
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

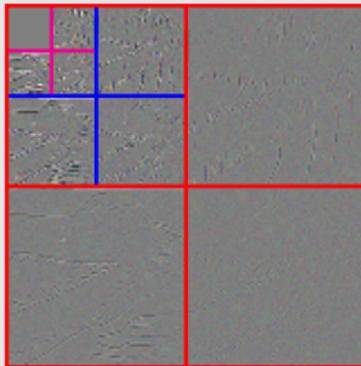
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

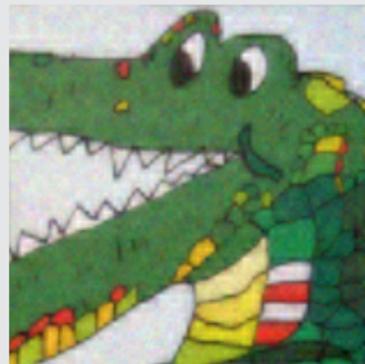
Shrinkage in the discrete wavelet domain



(a) y



(b) \hat{z} (Daubechies+LMMSE)



(c) \hat{x}

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

```
z = im.dwt(y, 3, h, g)  
zhat = shrink(z, lbd, sig)  
▶ xhat = im.idwt(zhat, 3, h, g)
```

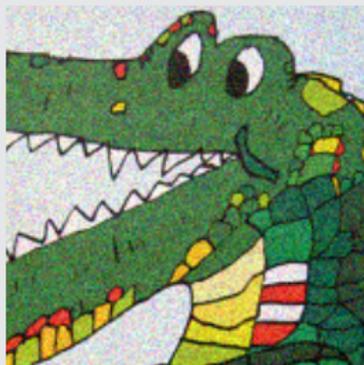
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

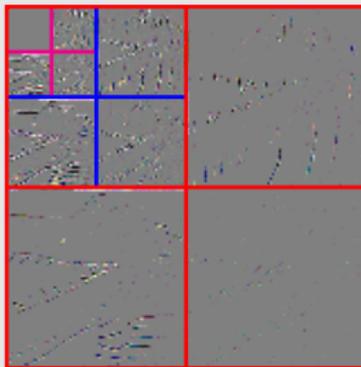
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

Shrinkage in the discrete wavelet domain



(a) y



(b) \hat{z} (Daubechies+Soft-T)



(c) \hat{x}

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

```
z = im.dwt(y, 3, h, g)  
zhat = shrink(z, lbd, sig)  
▶ xhat = im.idwt(zhat, 3, h, g)
```

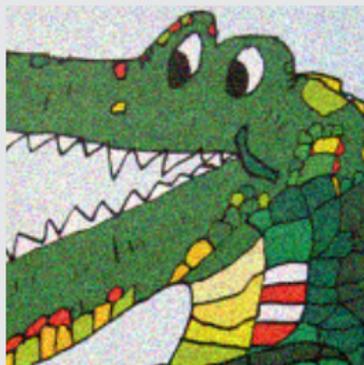
$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

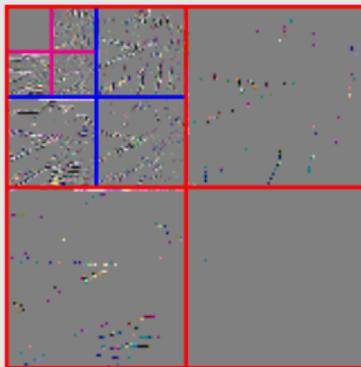
$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

Shrinkage in the wavelet domain

Shrinkage in the discrete wavelet domain



(a) y



(b) \hat{z} (Daubechies+Hard-T)



(c) \hat{x}

```
sig = 20  
y = x + sig * nr.randn(*x.shape)
```

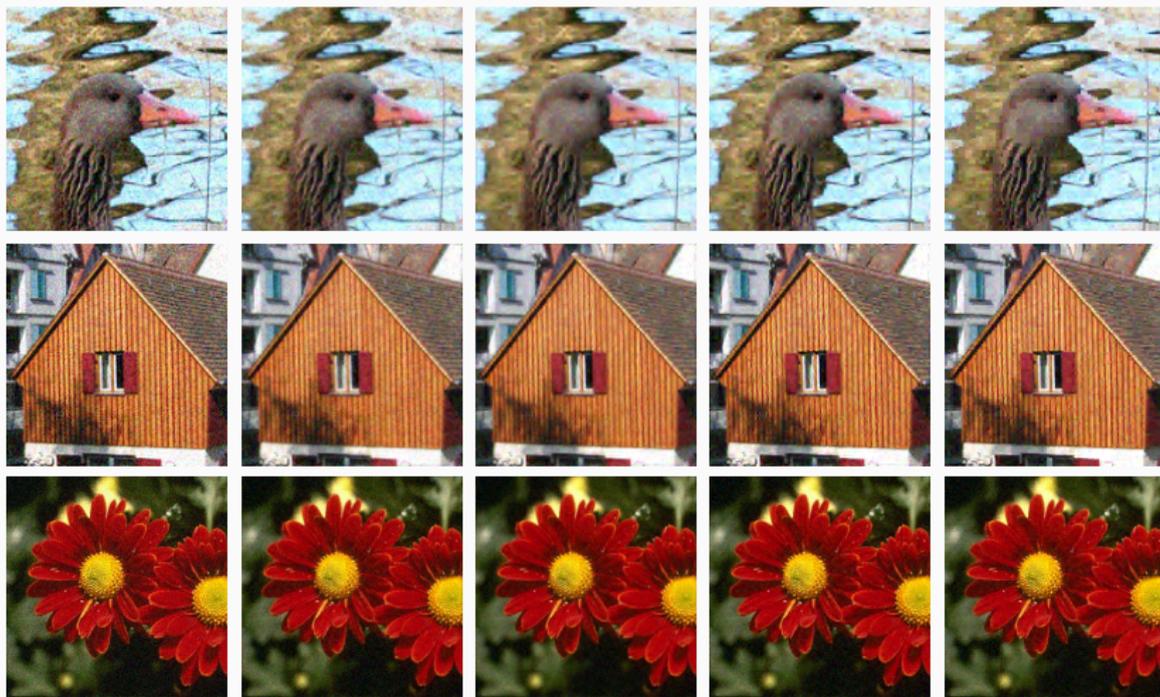
```
z = im.dwt(y, 3, h, g)  
zhat = shrink(z, lbd, sig)  
▶ xhat = im.idwt(zhat, 3, h, g)
```

$$z = \mathbf{W}y$$

$$\hat{z}_i = s(z_i; \lambda_i, \sigma)$$

$$\hat{x} = \mathbf{W}^{-1}\hat{z}$$

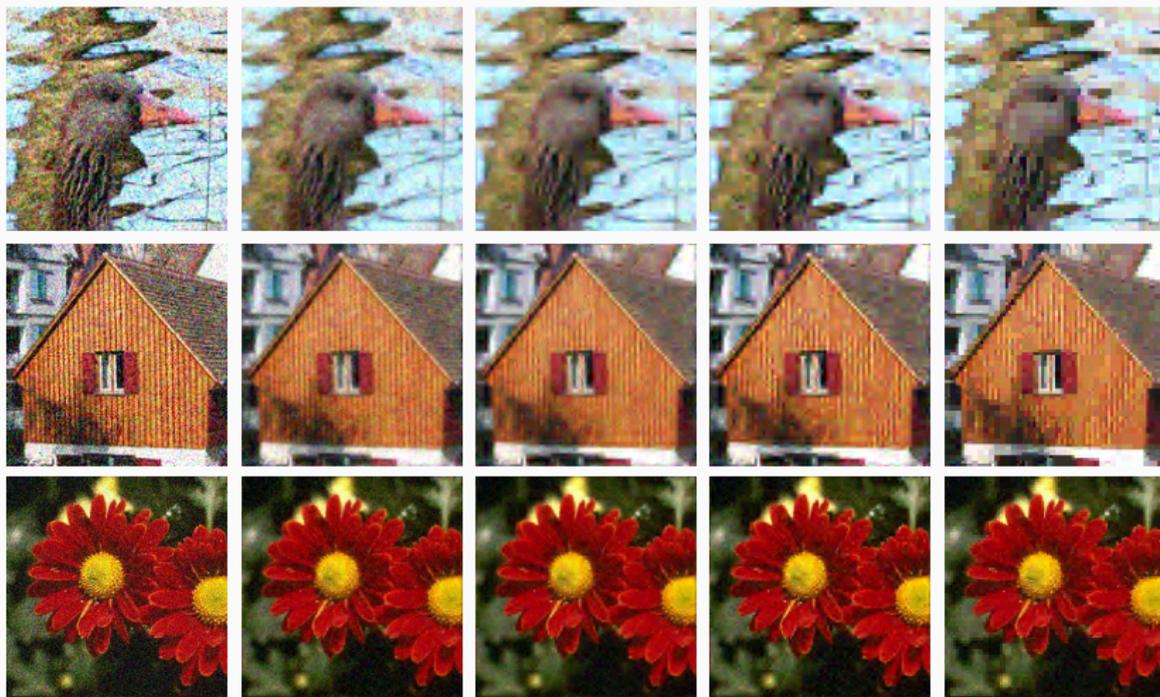
Shrinkage in the wavelet domain



(a) y ($\sigma = 20$) (b) Db2+LMMSE (c) Db2+Soft-T (d) Db2+Hard-T (e) Haar+Hard-T

For large noise: **Blocky effects** and **Ringing artifacts**

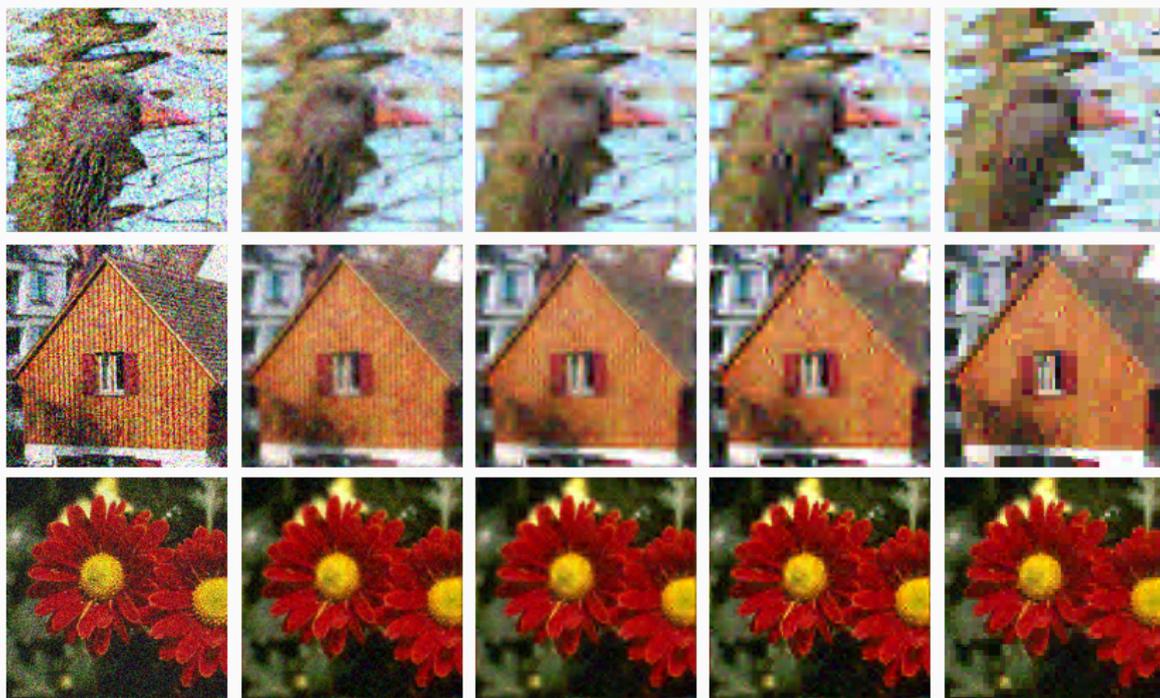
Shrinkage in the wavelet domain



(a) y ($\sigma = 40$) (b) Db2+LMMSE (c) Db2+Soft-T (d) Db2+Hard-T (e) Haar+Hard-T

For large noise: **Blocky effects** and **Ringing artifacts**

Shrinkage in the wavelet domain



(a) y ($\sigma = 60$) (b) Db2+LMMSE (c) Db2+Soft-T (d) Db2+Hard-T (e) Haar+Hard-T

For large noise: **Blocky effects** and **Ringing artifacts**

Shrinkage in the wavelet domain

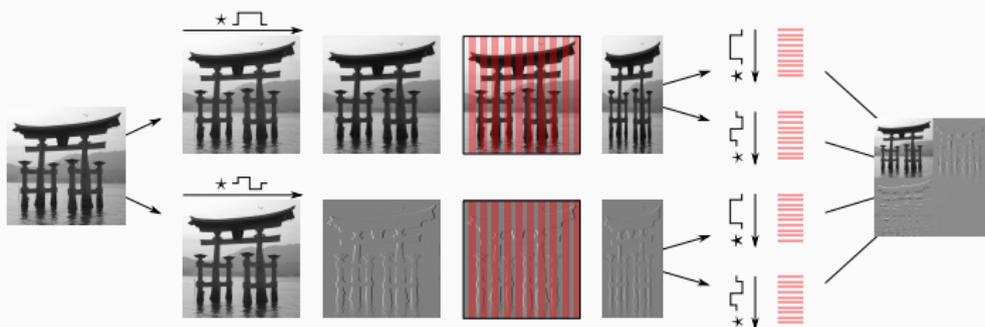


(a) y ($\sigma = 120$) (b) Db2+LMMSE (c) Db2+Soft-T (d) Db2+Hard-T (e) Haar+Hard-T

For large noise: **Blocky effects** and **Ringing artifacts**

Undecimated wavelet transforms

Limits of the discrete wavelet transform



- While Fourier shrinkage is translation invariant:

$$\psi(y^\tau) = \psi(y)^\tau \quad \text{where} \quad y^\tau(s) = y(s + \tau)$$

- Wavelet shrinkage is **not translation invariant**.
- This is due to the decimation step:

$$\mathbf{W} = \begin{pmatrix} G \downarrow_2 \\ H \downarrow_2 \end{pmatrix} \in \mathbb{R}^{n \times n} \quad \text{where} \quad M \downarrow_2 = "M[: :2, :]"$$

- This explains the **blocky artifacts** that we observe.

Undecimated discrete wavelet transform (UDWT)

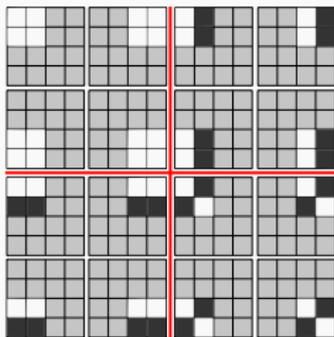


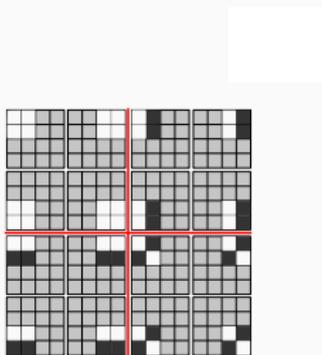
Figure 1 – Haar DWT

- Haar transform groups pixels by clusters of 4.
- Blocks are treated independently to each other.
- When similar neighbor blocks are shrunk differently, it becomes clearly visible in the image.
- This arises all the more as the noise level is large.

What if we do not decimate?

⇒ **UDWT, aka, stationary or translation-invariant wavelet transform.**

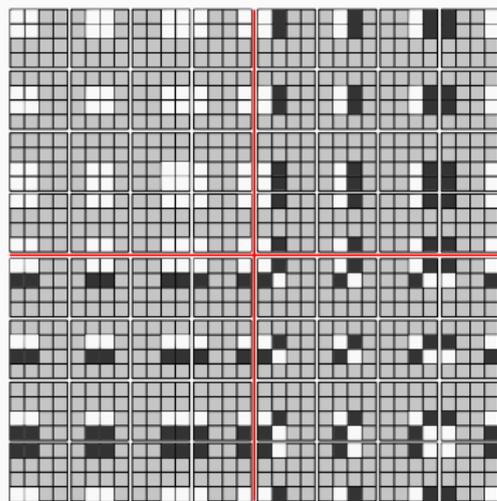
Undecimated discrete wavelet transform (UDWT)



Haar discrete wavelet transform (DWT)

1-scale DWT

- For a 4×4 image:
 4×4 coefficients.
- For n pixels: $K = n$ coefficients.



Haar undecimated discrete wavelet transform (UDWT)

1-scale UDWT

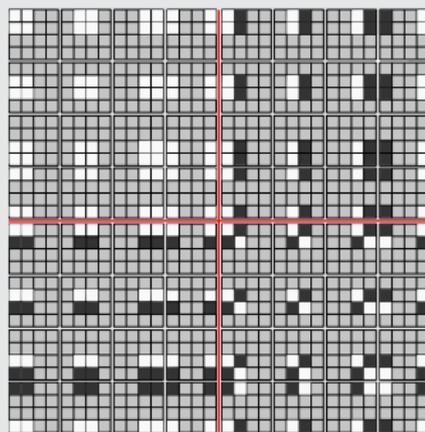
- For a 4×4 image:
 8×8 coefficients.
- For n pixels: $K = 4n$ coeffs.

What about multi-scale?

Undecimated discrete wavelet transform (UDWT)

A trous algorithm (with holes)

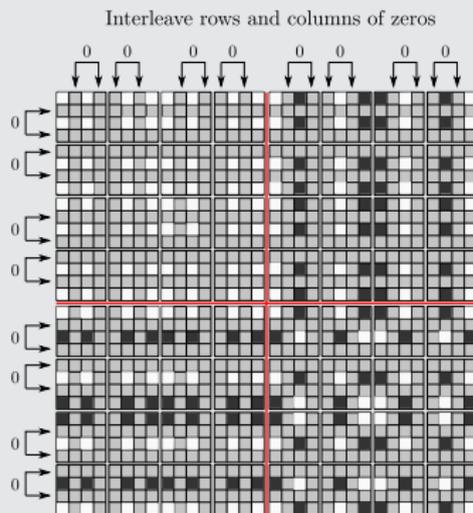
(Holschneider *et al.*, 1989)



$$g = \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

$$h = \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

Haar UDWT, first scale



$$g^{:1} = \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

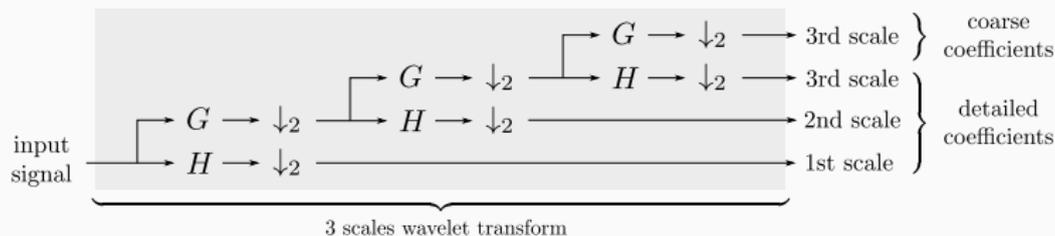
$$h^{:1} = \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

Haar UDWT, second scale

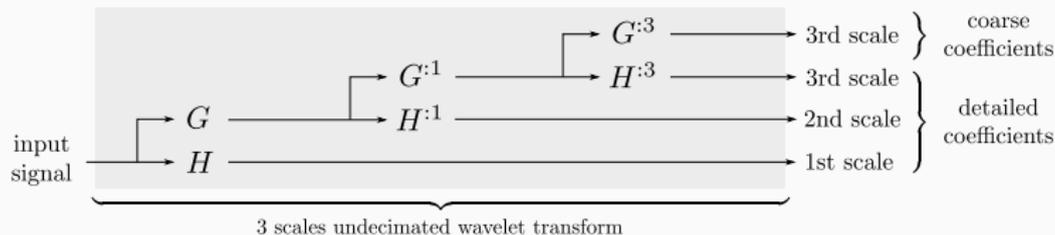
Instead of decimating the coefficients at each scale j , upsample the filters h and g by injecting $2^j - 1$ zeros between each entries.

Undecimated discrete wavelet transform (UDWT)

DWT: Mallat's dyadic pyramidal multi-resolution scheme



UDWT: A trous algorithm - G^p : inject p zeros between each filter coeffs



Multi-scales: $K = (1 + J(2^d - 1))n$ coeffs (J : #scales, $d = 2$ for images)

Undecimated discrete wavelet transform (UDWT)

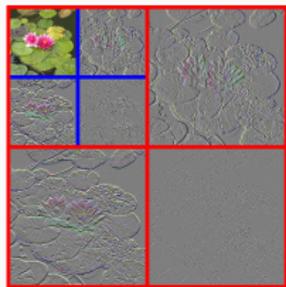
Implementation of 2D UDWT (*A trous* algorithm)

```
def udwt(x, J, h, g):  
    if J == 0:  
        return x[:, :, np.newaxis]  
    tmp_h = flip(convolve(flip(x), h)) / 2  
    tmp_g = flip(convolve(flip(x), g)) / 2  
    detail = np.stack((convolve(tmp_g, h),  
                      convolve(tmp_h, g),  
                      convolve(tmp_h, h)), axis=2)  
    coarse = convolve(tmp_g, g)  
    h2 = interleave0(h)  
    g2 = interleave0(g)  
    z = np.concatenate((udwt(coarse, J - 1, h2, g2), detail), axis=2)  
    return z
```

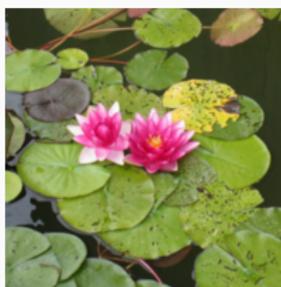
Linear complexity.

Can be easily modified to reduce memory usage.

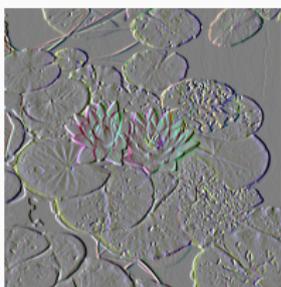
Undecimated discrete wavelet transform (UDWT)



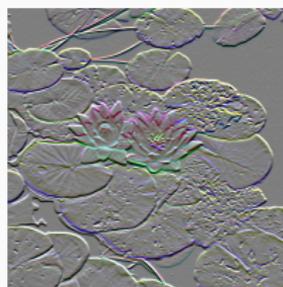
(a) DWT (J=2)



(b) UDWT Coarse Sc.



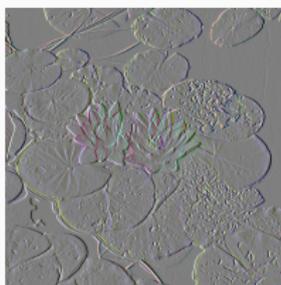
(c) Detailed Scale #1



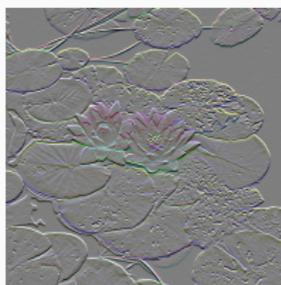
(d) Detailed Scale #2



(e) Detailed Scale #3



(f) Detailed Scale #4



(g) Detailed Scale #5



(h) Detailed Scale #6

What about its inverse transform?

Undecimated discrete wavelet transform (UDWT)

DWT – Wavelet basis – and inverse DWT

- The DWT $\mathbf{W} \in \mathbb{R}^{n \times n}$ has n columns and n rows.
- The n columns/rows of \mathbf{W} are orthonormal.
- The inverse DWT is $\mathbf{W}^{-1} = \mathbf{W}^*$.
- One-to-one relationship between an image and its wavelet coefficients.

UDWT – Redundant wavelet dictionary

- The UDWT $\bar{\mathbf{W}} \in \mathbb{R}^{K \times n}$ has $K = (1 + J(2^d - 1))n$ rows and n columns.
- The rows of $\bar{\mathbf{W}}$ cannot be linearly independent: not a basis.
- They are said to form a **redundant/overcomplete wavelet dictionary**.
- Since $\bar{\mathbf{W}}$ is non square, it is not invertible.

Note: **redundant dictionaries necessarily favor sparsity**.

Pseudo-inverse UDWT

- Nevertheless, the n columns are orthonormal, then: $\bar{W}^* = \bar{W}^+$
- It satisfies $\bar{W}^+ \bar{W} = \text{Id}_n$, but $\bar{W} \bar{W}^+ \neq \text{Id}_K$
 - image $\xrightarrow{\bar{W}}$ coefficients $\xrightarrow{\bar{W}^+}$ back to the original image,
 - coefficients $\xrightarrow{\bar{W}^+}$ image $\xrightarrow{\bar{W}}$ not necessarily the same coefficients.
- Satisfies the Parseval equality

$$\langle \bar{W}x, \bar{W}y \rangle = \langle x, \bar{W}^* \bar{W}y \rangle = \langle x, \bar{W}^+ \bar{W}y \rangle = \langle x, y \rangle$$

- In the vocabulary of linear algebra: \bar{W} is called a **tight-frame**.

Consequence: an algorithm for \bar{W}^+ can be obtained.

Undecimated discrete wavelet transform (UDWT)

Implementation of 2D Inverse UDWT

```
def iudwt(z, J, h, g):  
    if J == 0:  
        return z[:, :, 0]  
    h2 = interleave0(h)  
    g2 = interleave0(g)  
    coarse = iudwt(z[:, :, :-3], J - 1, h2, g2)  
    tmpg = convolve(coarse, g[:, :-1]) + \  
            convolve(z[:, :, -3], h[:, :-1])  
    tmpg = convolve(z[:, :, -2], g[:, :-1]) + \  
            convolve(z[:, :, -1], h[:, :-1])  
    x = (flip(convolve(flip(tmpg), g[:, :-1])) + \  
         flip(convolve(flip(tmpg), h[:, :-1]))) / 2  
    return x
```

Linear complexity again.

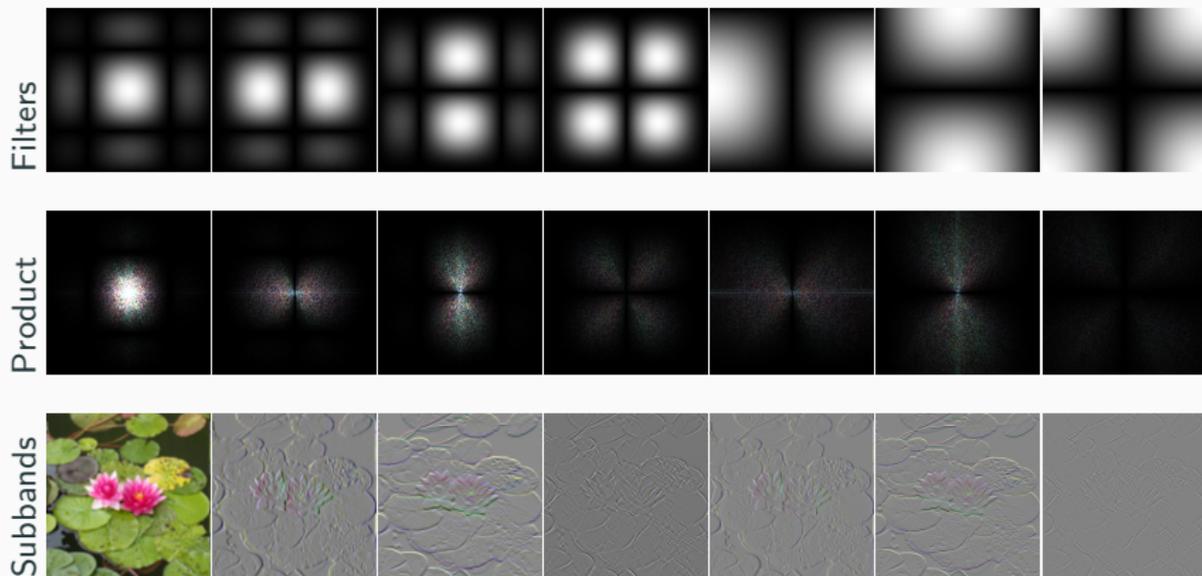
Can also be easily modified to reduce memory usage.

Can we be more efficient?

Filter bank

- The UDWT of x for subband k , $x \mapsto (\mathbf{W}x)_k$ is
linear and translation invariant (LTI)
 \Rightarrow It's a convolution.
- The UDWT is a filter bank:
a **set of band-pass filters** that separates the input image into multiple components.
- Each filter can be **represented by its frequential response**.
- Direct and inverse transform: **implementation in the Fourier domain**.

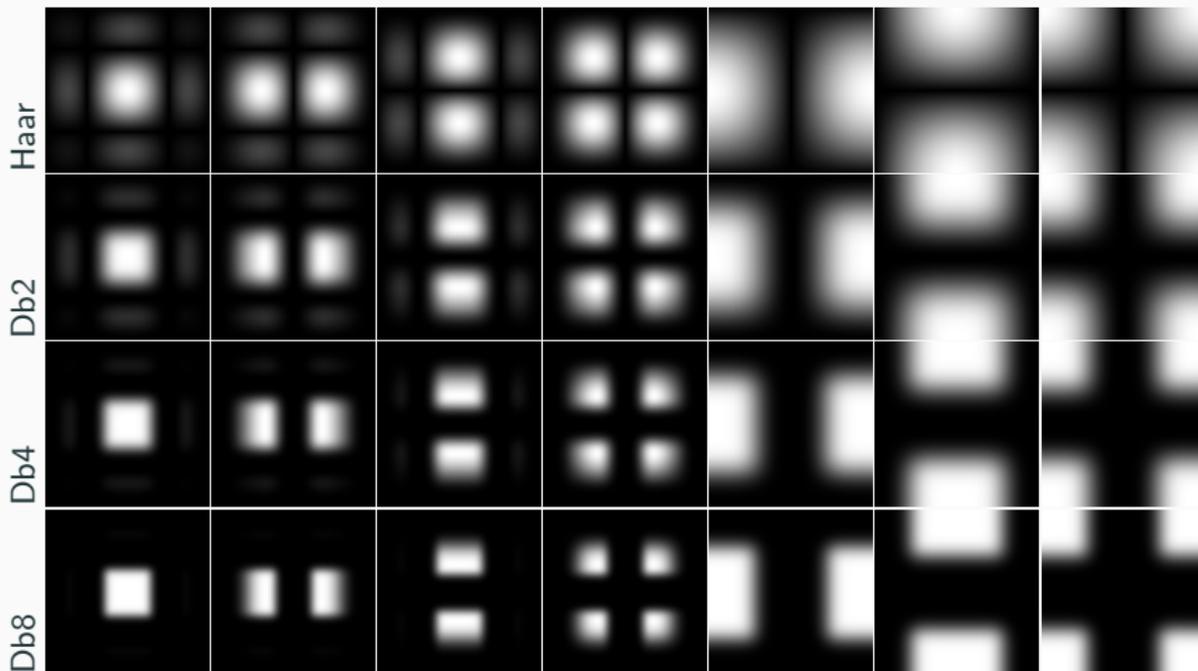
Undecimated discrete wavelet transform (UDWT)



(a) Coarse 2 (b) Details 2 (c) Details 2 (d) Details 2 (e) Details 1 (f) Details 1 (g) Details 1

Haar with $J = 2$ levels of decomposition

Undecimated discrete wavelet transform (UDWT)

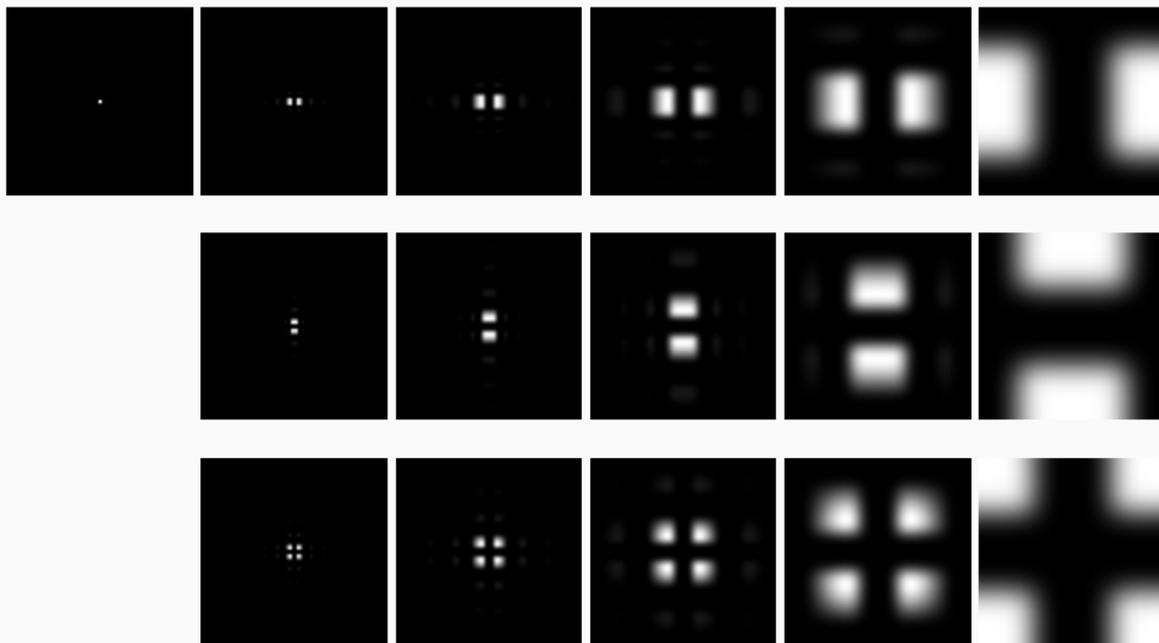


(a) Coarse 2 (b) Details 2 (c) Details 2 (d) Details 2 (e) Details 1 (f) Details 1 (g) Details 1

Haar: band pass with side lobes.

Db8: closer to ideal band pass.

Undecimated discrete wavelet transform (UDWT)



Db4 with $J = 6$ levels of decomposition

How to create such a filter bank?

Undecimated discrete wavelet transform (UDWT)

UDWT: Creation of the filter bank (offline)

```
def udwt_create_fb(n1, n2, J, h, g, ndim=3):
    if J == 0:
        return np.ones((n1, n2, 1, *[1] * (ndim - 2)))
    h2 = interleave0(h)
    g2 = interleave0(g)
    fbrec = udwt_create_fb(n1, n2, J - 1, h2, g2, ndim=ndim)
    gf1 = nf.fft(fftpad(g, n1), axis=0)
    hf1 = nf.fft(fftpad(h, n1), axis=0)
    gf2 = nf.fft(fftpad(g, n2), axis=0)
    hf2 = nf.fft(fftpad(h, n2), axis=0)
    fb = np.zeros((n1, n2, 4), dtype=np.complex128)
    fb[:, :, 0] = np.outer(gf1, gf2) / 2
    fb[:, :, 1] = np.outer(gf1, hf2) / 2
    fb[:, :, 2] = np.outer(hf1, gf2) / 2
    fb[:, :, 3] = np.outer(hf1, hf2) / 2
    fb = fb.reshape(n1, n2, 4, *[1] * (ndim - 2))
    fb = np.concatenate((fb[:, :, 0:1] * fbrec, fb[:, :, -3:]),
                        axis=2)
    return fb
```

Undecimated discrete wavelet transform (UDWT)

UDWT: Direct transform using the filter bank (online)

```
def fb_apply(x, fb):  
    x = nf.fft2(x, axes=(0, 1))  
    z = fb * x[:, :, np.newaxis]  
    z = np.real(nf.ifft2(z, axes=(0, 1)))  
    return z
```

UDWT: Inverse transform using the filter bank (online)

```
def fb_adjoint(z, fb):  
    z = nf.fft2(z, axes=(0, 1))  
    x = (np.conj(fb) * z).sum(axis=2)  
    x = np.real(nf.ifft2(x, axes=(0, 1)))  
    return x
```

Much more efficient than previous implementation when $J > 1$

Shrinkage with UDWT

- Consider a denoising problem $y = x + w$ with noise variance σ^2 .
- Shrink the $K \geq n$ coefficients independently.

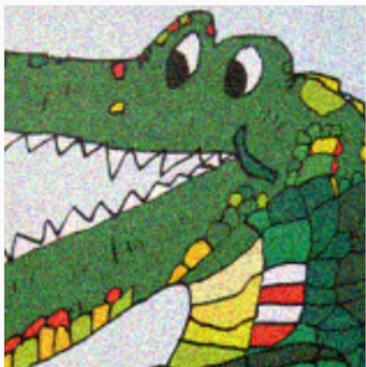
$$\hat{x}^* = \underbrace{\bar{\mathbf{W}}^+}_{\text{Pseudo-inverse}} \hat{z} \quad \text{where} \quad \underbrace{\hat{z}_i = s(z_i; \lambda_i, \sigma_i)}_{\text{shrinkage}} \quad \text{and} \quad z = \underbrace{\bar{\mathbf{W}} y}_{\text{Redundant representation}}$$

Rule of thumb for soft-thresholding:

- For the orthonormal DWT \mathbf{W} : increase λ_i as $\sqrt{2}^{d(j_i-1)}$.
- For the tight-frame UDWT $\bar{\mathbf{W}}$: increase λ_i as: $2^{d(j_i-1/2)}$.

(j_i scale for coefficient i , $d = 2$ for images).

Reconstruction with the UDWT



(a) y



(b) $DWT(3)+Haar+HT$

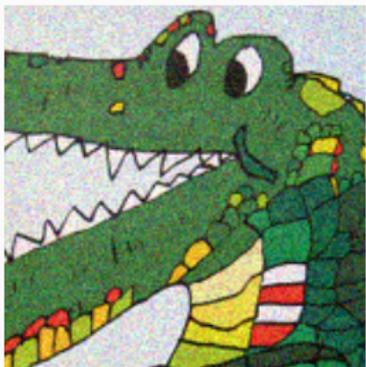


(c) $DWT(3)+Db2+HT$



(e) $UDWT(3)+Db2+HT$

Reconstruction with the UDWT



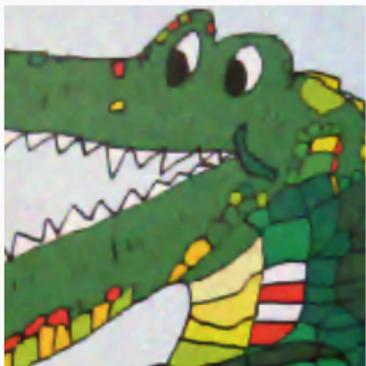
(a) y



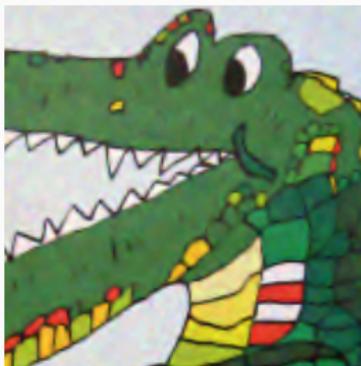
(b) $\text{DWT}(3)+\text{Haar}+\text{HT}$



(c) $\text{DWT}(3)+\text{Db2}+\text{HT}$



(d) $\text{UDWT}(3)+\text{Haar}+\text{HT}$

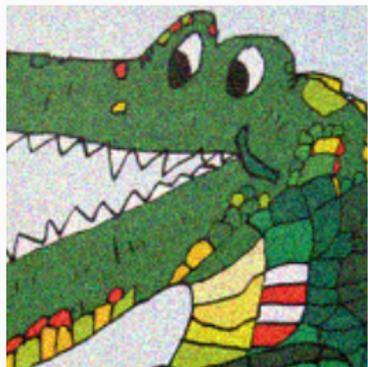


(e) $\text{UDWT}(3)+\text{Db2}+\text{HT}$



(f) $\text{UDWT}(3)+\text{Db8}+\text{HT}$

Reconstruction with the UDWT



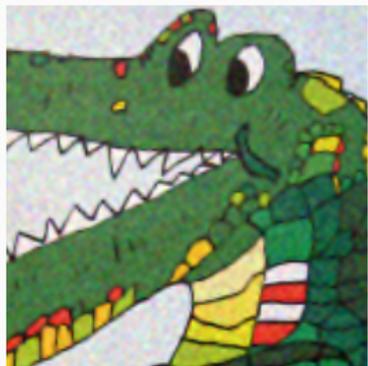
(a) y



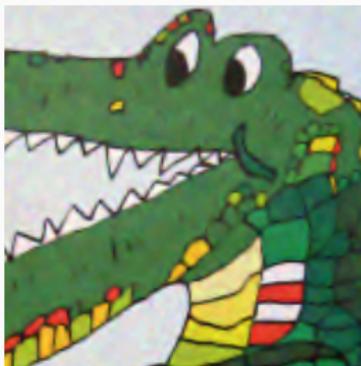
(b) $\text{DWT}(3)+\text{Haar}+\text{HT}$



(c) $\text{DWT}(3)+\text{Db}2+\text{HT}$



(d) $\text{UDWT}(1)+\text{Db}2+\text{HT}$

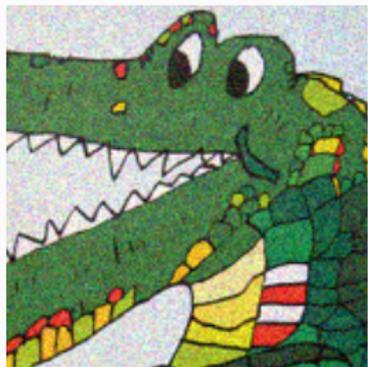


(e) $\text{UDWT}(3)+\text{Db}2+\text{HT}$



(f) $\text{UDWT}(5)+\text{Db}2+\text{HT}$

Reconstruction with the UDWT



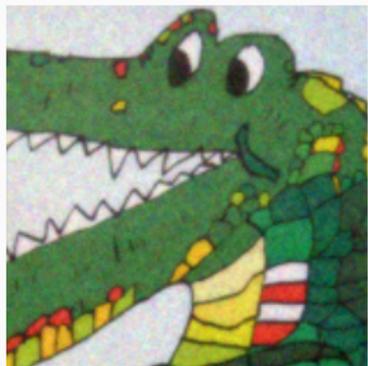
(a) y



(b) $DWT(3)+Haar+HT$



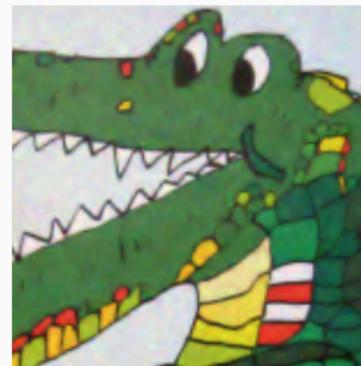
(c) $DWT(3)+Db2+HT$



(d) $UDWT(3)+Db2+Linear$



(e) $UDWT(3)+Db2+HT$



(f) $UDWT(3)+Db2+ST$

Reconstruction with the UDWT



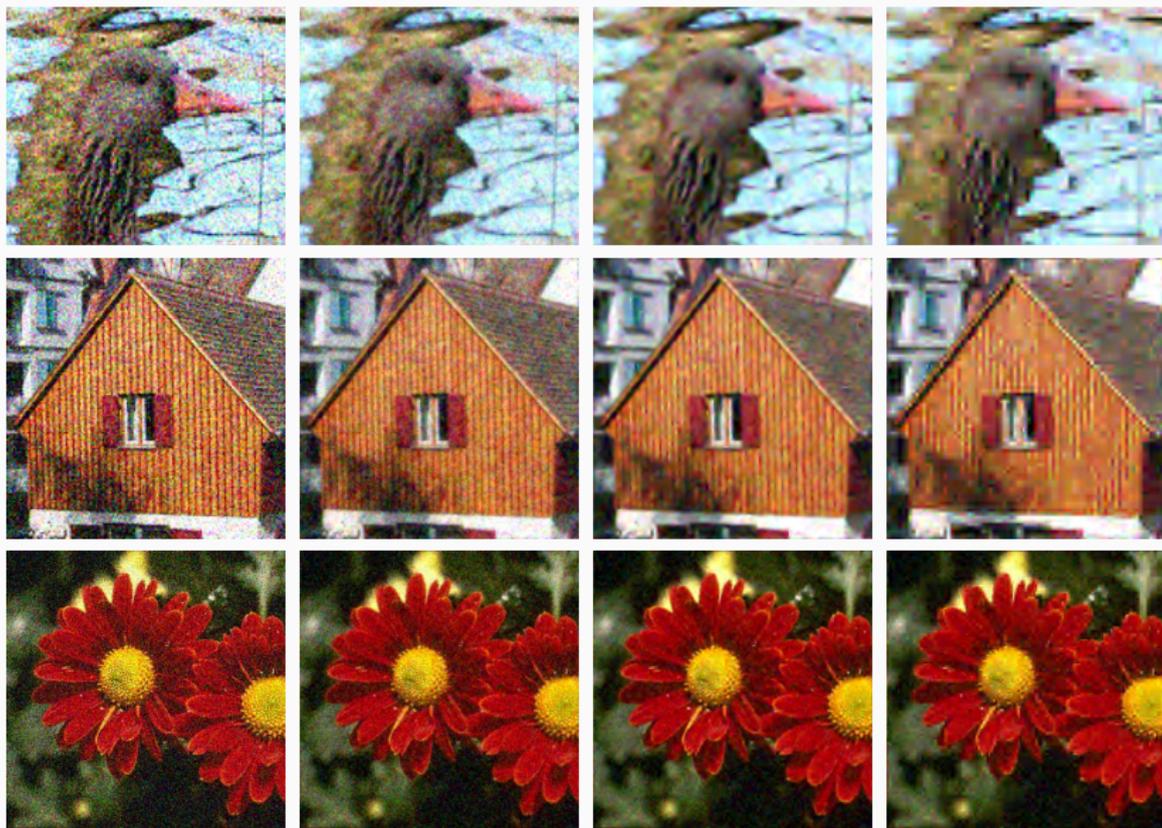
(a) y ($\sigma = 20$)

(b) UDWT+Lin.

(c) UDWT+HT

(d) DWT+HT

Reconstruction with the UDWT



(a) y ($\sigma = 40$)

(b) UDWT+Lin.

(c) UDWT+HT

(d) DWT+HT

Reconstruction with the UDWT



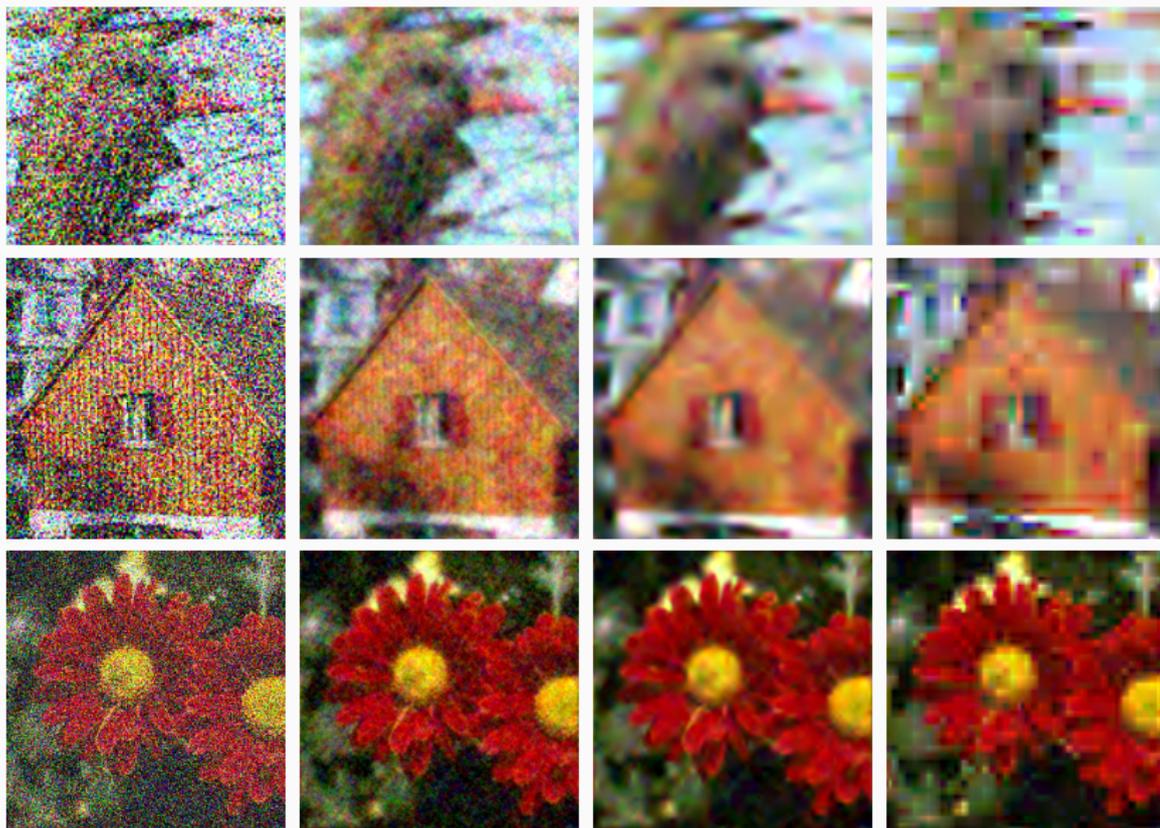
(a) y ($\sigma = 60$)

(b) UDWT+Lin.

(c) UDWT+HT

(d) DWT+HT

Reconstruction with the UDWT



(a) y ($\sigma = 120$)

(b) UDWT+Lin.

(c) UDWT+HT

(d) DWT+HT

Reconstruction with the UDWT

$$\hat{x}^* = \underbrace{\bar{W}^+}_{\text{Pseudo-inverse}} \hat{z} \quad \text{where} \quad \underbrace{\hat{z}_i = s(z_i; \lambda_i, \sigma_i)}_{\text{shrink } K \text{ coefficients}} \quad \text{and} \quad z = \underbrace{\bar{W}y}_{\text{Redundant representation}}$$

Connection with Bayesian shrinkage?

- Since the rows of \bar{W} are linearly dependent, the **coefficients z_i are necessarily correlated** (non-white).
- Shrink the $K \geq n$ coefficients independently, even though they cannot be assumed independent.
- This estimator has **no Bayesian interpretation**, it does not correspond to the MMSE or MAP.

How to use the UDWT in the Bayesian context?

Bayesian analysis model

Whitening model: Consider $\eta = \mathbf{\Lambda}^{-1/2} \mathbf{W}x$ (η coeffs)
such that $\mathbb{E}[\eta] = 0_n$ and $\text{Var}[\eta] = \text{Id}_n$

Analysis: images can be transformed to white coeffs.

⚠ Non-sense when rows of \mathbf{W} are redundant.

Bayesian synthesis model

Generative model: Consider $x = \bar{\mathbf{W}}^+ \mathbf{\Lambda}^{1/2} \eta$ (η code)
such that $\mathbb{E}[\eta] = 0_K$ and $\text{Var}[\eta] = \text{Id}_K$

Synthesis: images can be generated from a white code.

☺ Always well-founded.

Forward model: $y = x + w$

Maximum a Posteriori for the Synthesis model

- Instead of looking for x , consider the MAP for the code η

$$\begin{aligned}\hat{\eta}^* &\in \operatorname{argmax}_{\eta \in \mathbb{R}^K} p(\eta|y) \\ &= \operatorname{argmin}_{\eta \in \mathbb{R}^K} [-\log p(y|\eta) - \log p(\eta)] \\ &= \operatorname{argmin}_{\eta \in \mathbb{R}^K} \left[\frac{1}{2} \|y - \bar{\mathbf{W}}^+ \mathbf{\Lambda}^{1/2} \eta\|_2^2 - \log p(\eta) \right]\end{aligned}$$

- Once you get $\hat{\eta}^*$, generate the image \hat{x}^* as

$$\hat{x}^* = \bar{\mathbf{W}}^+ \mathbf{\Lambda}^{1/2} \hat{\eta}^*$$

What interpretation?

Penalized least square with redundant dictionary

- Consider the **redundant wavelet dictionary** $D = \bar{W}^+ \Lambda^{1/2}$

$$D = (\underbrace{d_1, d_2, \dots, d_K}_{\text{linearly dependent atoms}}), \quad \|d_i\| = \lambda_i, \quad K \geq n$$

- Goal: Look for a code $\eta \in \mathbb{R}^K$, such that \hat{x} close to y

$$\hat{x} = D\eta = \sum_{i=1}^K \eta_i d_i = \text{"linear comb. of the redundant atoms } d_i \text{ of } D\text{"}$$

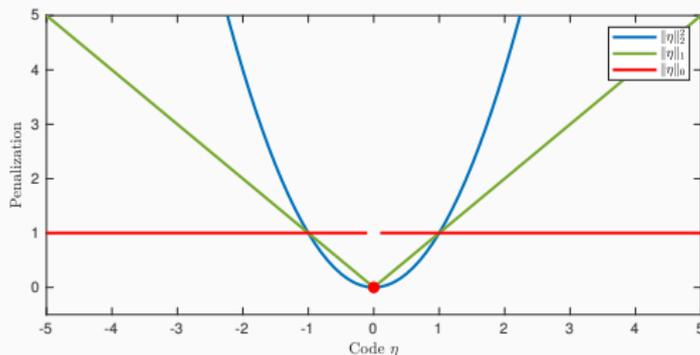
- Since D is redundant, different codes η produce the same image x .
- Penalize independently each η_i to select a relevant one

$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \left[\frac{1}{2} \|y - \bar{W}^+ \Lambda^{1/2} \eta\|_2^2 - \sum_{i=1}^K \log p(\eta_i) \right]$$

What choice for $\log p(\eta_i)$?

Penalized least square with redundant dictionary

- $\frac{1}{2} \|y - D\eta\|_2^2 + \frac{\tau^2}{2} \|\eta\|_2^2, \quad \|\eta\|_2^2 = \sum_i \eta_i^2 \quad \leftarrow$ Ridge regression
- $\frac{1}{2} \|y - D\eta\|_2^2 + \tau \|\eta\|_1, \quad \|\eta\|_1 = \sum_i |\eta_i| \quad \leftarrow$ LASSO
- $\frac{1}{2} \|y - D\eta\|_2^2 + \frac{\tau^2}{2} \|\eta\|_0, \quad \|\eta\|_0 = \sum_i \mathbf{1}_{\{\eta_i \neq 0\}} \quad \leftarrow$ Sparse regression



When D is redundant, these problems are no longer separable.
They require **large-scale optimization techniques**.

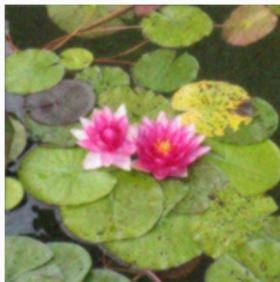
Regularizations and optimization

Ridge/Smooth regression

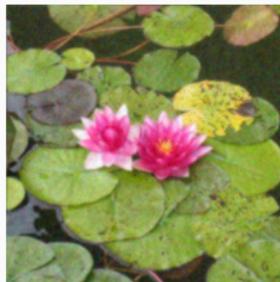
- Convex energy: $E(\eta) = \frac{1}{2} \|y - D\eta\|_2^2 + \frac{\tau^2}{2} \|\eta\|_2^2$
- Gradient: $\nabla E(\eta) = D^*(D\eta - y) + \tau^2 \eta$
- Optimality conditions: $\hat{\eta}^* = (D^*D + \tau^2 \text{Id}_K)^{-1} D^*y$
- For UDWT: this is an LTI filter \equiv convolution (non adaptive)



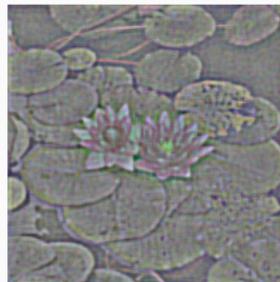
(a) y



(b) Linear shrink



(c) Ridge

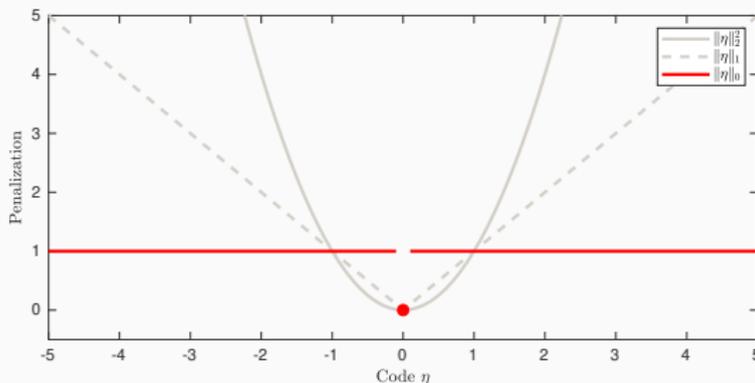


(d) Difference

Ridge \neq Linear shrinkage (except if D is orthogonal).

Sparse regression / ℓ_0 regularization (1/3)

- Energy: $E(\eta) = \frac{1}{2}\|y - D\eta\|_2^2 + \frac{\tau^2}{2}\|\eta\|_0$
- Penalty: $\|\eta\|_0 = \#\text{non zero elements in } \eta$
- Non-convex: $0.5 = \frac{1}{2}(\|0\|_0 + \|1\|_0) < \|0.5\|_0 = 1$
- Produces optimal sparse solutions adapted to the signal ☺
- But, non-differentiable and discontinuous. ☹



Sparse regression / ℓ_0 regularization (2/3)

- If D is orthogonal: solution given by the Hard-Thresholding.
- Otherwise, exact solution obtained by **brute force**:
 - For all possible support $\mathcal{I} \subseteq \{1, \dots, K\}$ (set of non-zero coefficients)
 - Solve the least square estimation problem:

$$\operatorname{argmin}_{(\eta_i)_{i \in \mathcal{I}}} \frac{1}{2} \|y - \sum_{i \in \mathcal{I}} \eta_i a_i\|_2^2$$

- Pick the solution that minimizes E .
- **NP-hard combinatorial problem:**

$$\#\text{subsets} = \sum_{k=0}^K \binom{K}{k} = 2^K$$

Sparse regression / ℓ_0 regularization (3/3)

- Sub-optimal solutions can be obtained by **greedy** algorithms.
- **Matching pursuit (MP):** (Mallat, 1993)
 - ① Initialization: $r \leftarrow y, \eta \leftarrow 0, k \leftarrow 0$
 - ② Choose i maximizing $|\mathbf{D}^* r|_i = |\langle d_i, r \rangle|$
 - ③ Compute $\alpha = \langle r, d_i \rangle / \|d_i\|_2^2$
 - ④ Update $r \leftarrow r - \alpha d_i$
 - ⑤ Update $\eta_i = \alpha$
 - ⑥ Update $k \leftarrow k + 1$
 - ⑦ Back to step 2 while $E(\eta) = \frac{1}{2}\|r\|_2^2 + \frac{\tau^2}{2}k$ decreases
- Lots of iterations: complexity $O(kn)$, with k the sparsity of the solution.
- Each iteration requires to compute an UDWT.
- **Extensions:** OMP (Tropp & Gilbert, 2007), CoSaMP (Needel & Tropp, 2009)

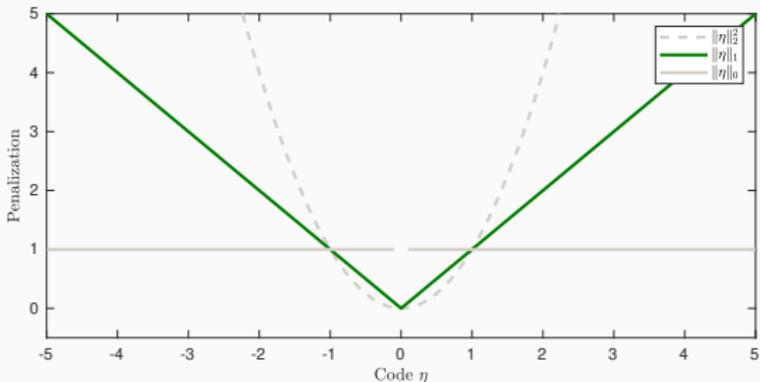
Least Absolute Shrinkage and Selection Operator (LASSO)

Convex relaxation: **Take the best of both worlds: sparsity and convexity**

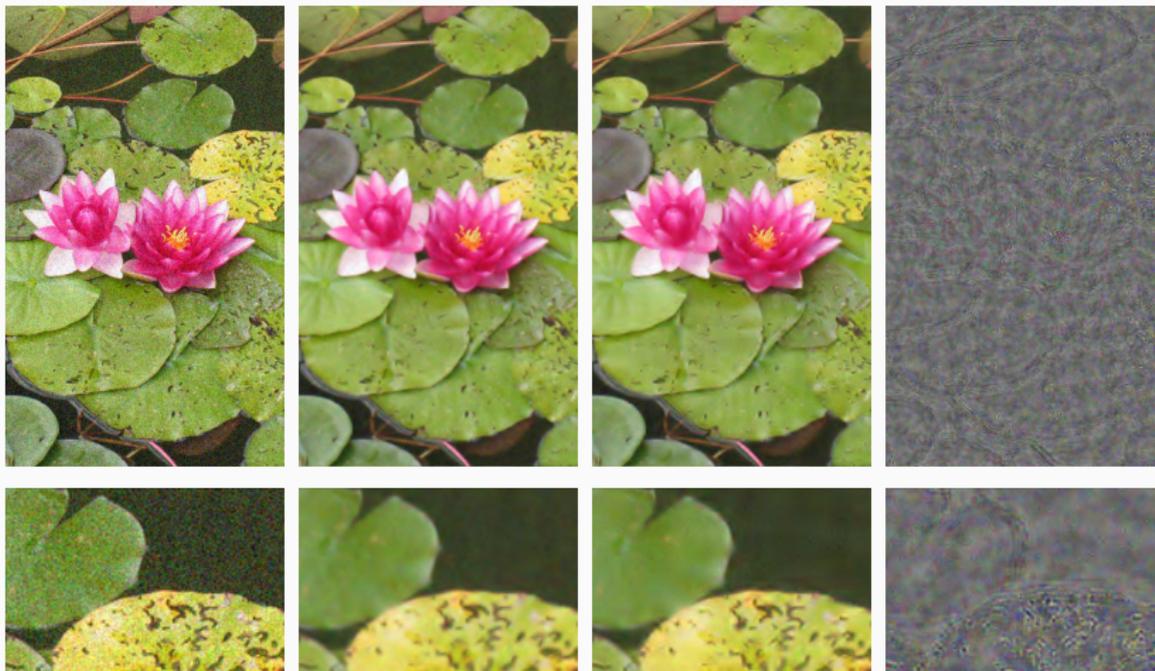
LASSO / ℓ_1 regularization

(Tibshirani 1996)

- Convex energy: $E(\eta) = \frac{1}{2} \|y - D\eta\|_2^2 + \tau \|\eta\|_1$
- Non-smooth penalty: $\|\eta\|_1 = \sum_{i=1}^K |\eta_i|$
- If D is orthogonal: solution given by the Soft-Thresholding.
- Produces also sparse solutions adapted to the signal ☺



Least Absolute Shrinkage and Selection Operator



(a) Input

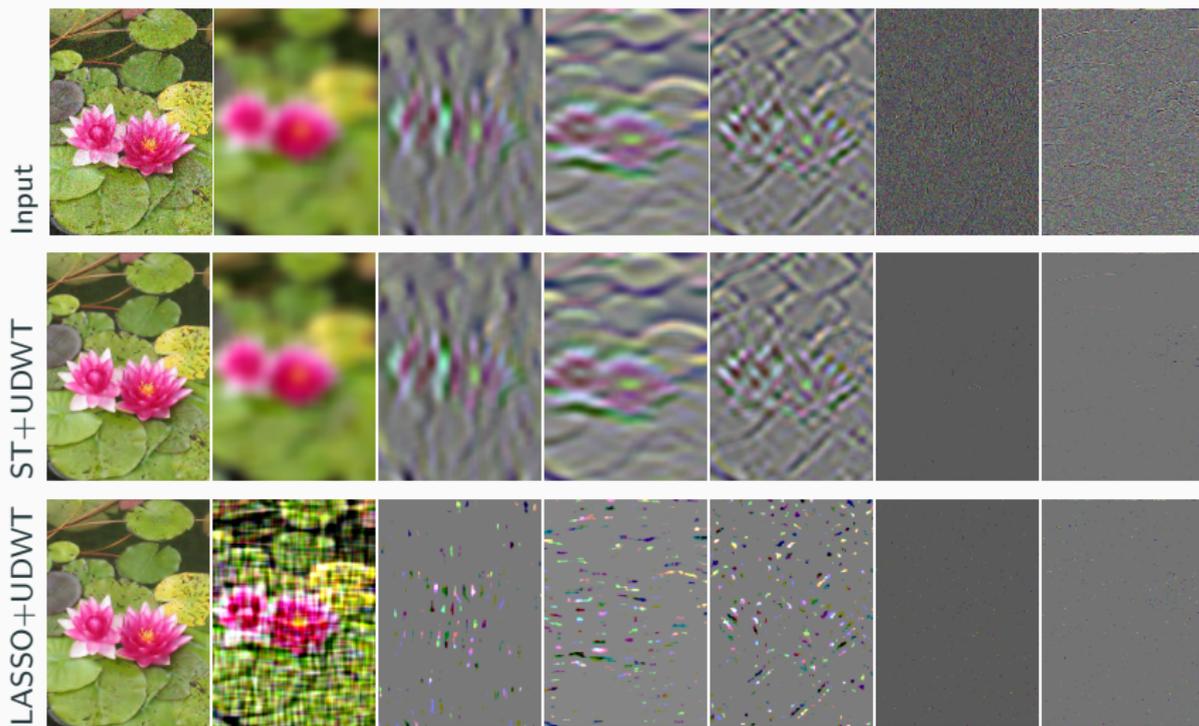
(b) ST+UDWT (1s)

(c) LASSO+UDWT (30s)

(d) Difference

Though the solutions look alike, their codes η are very different.

Least Absolute Shrinkage and Selection Operator



(a) Image (b) Coarse 5 (c) Scale 5 (d) Scale 5 (e) Scale 5 (f) Scale 1 (g) Scale 1

The LASSO creates much sparser codes than ST only.

Why use the LASSO if shrinkage in the UDWT provides similar results?

- Shrinkage in the UDWT domain can only be applied for denoising problems.

- The LASSO can be adapted to **inverse-problems**:

$$\hat{x}^* = D\hat{\eta}^* \quad \text{with} \quad \hat{\eta}^* \in \underset{\eta \in \mathbb{R}^K}{\operatorname{argmin}} \left[\frac{1}{2} \|y - HD\eta\|_2^2 + \tau \|\eta\|_1 \right]$$

But it requires solving a non-smooth convex optimization problem.

Solution: use **sub-differential** and **Fermat's rule**.

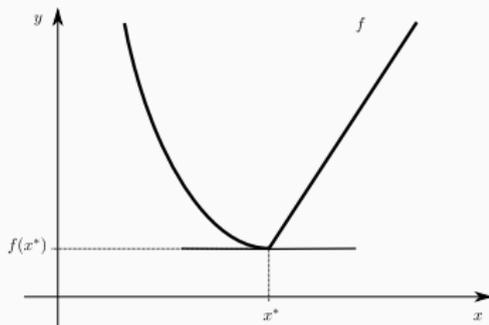
Definition (Sub-differential)

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, $u \in \mathbb{R}^n$ is a **sub-gradient** of f at x^* , if for all $x \in \mathbb{R}^n$

$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle.$$

- The **sub-differential** is the set of sub-gradients

$$\partial f(x^*) = \{u \in \mathbb{R}^n : \forall x \in \mathbb{R}^n, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$



If the sub-gradient is unique, f is differentiable and $\partial f(x) = \{\nabla f(x)\}$.

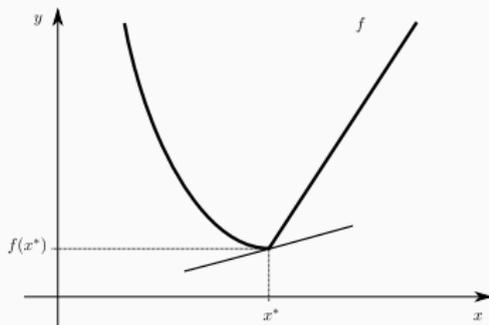
Definition (Sub-differential)

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, $u \in \mathbb{R}^n$ is a **sub-gradient** of f at x^* , if for all $x \in \mathbb{R}^n$

$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle.$$

- The **sub-differential** is the set of sub-gradients

$$\partial f(x^*) = \{u \in \mathbb{R}^n : \forall x \in \mathbb{R}^n, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$



If the sub-gradient is unique, f is differentiable and $\partial f(x) = \{\nabla f(x)\}$.

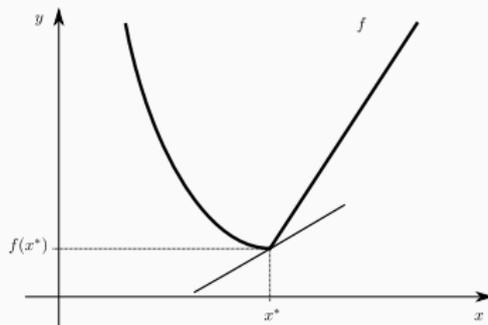
Definition (Sub-differential)

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, $u \in \mathbb{R}^n$ is a **sub-gradient** of f at x^* , if for all $x \in \mathbb{R}^n$

$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle.$$

- The **sub-differential** is the set of sub-gradients

$$\partial f(x^*) = \{u \in \mathbb{R}^n : \forall x \in \mathbb{R}^n, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$



If the sub-gradient is unique, f is differentiable and $\partial f(x) = \{\nabla f(x)\}$.

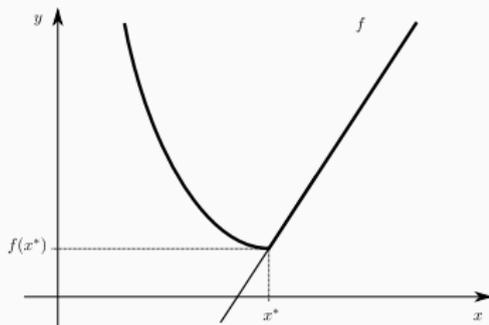
Definition (Sub-differential)

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, $u \in \mathbb{R}^n$ is a **sub-gradient** of f at x^* , if for all $x \in \mathbb{R}^n$

$$f(x) \geq f(x^*) + \langle u, x - x^* \rangle.$$

- The **sub-differential** is the set of sub-gradients

$$\partial f(x^*) = \{u \in \mathbb{R}^n : \forall x \in \mathbb{R}^n, f(x) \geq f(x^*) + \langle u, x - x^* \rangle\}.$$



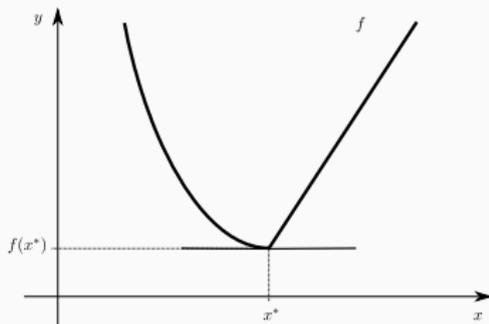
If the sub-gradient is unique, f is differentiable and $\partial f(x) = \{\nabla f(x)\}$.

Theorem (Fermat's rule)

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function, then

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} f(x) \quad \Leftrightarrow \quad 0_n \in \partial f(x^*)$$

If f is also differentiable, this corresponds to the standard rule $\nabla f(x^*) = 0_n$.

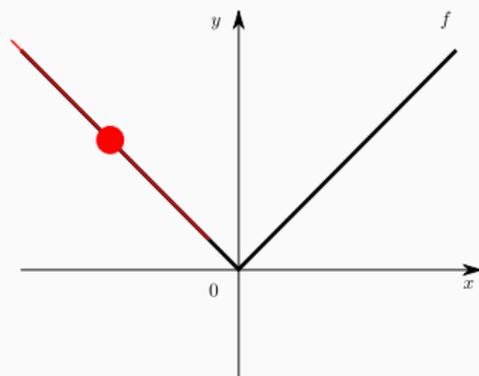


Minimizers are the only points with a horizontal tangent

Non-smooth convex optimization

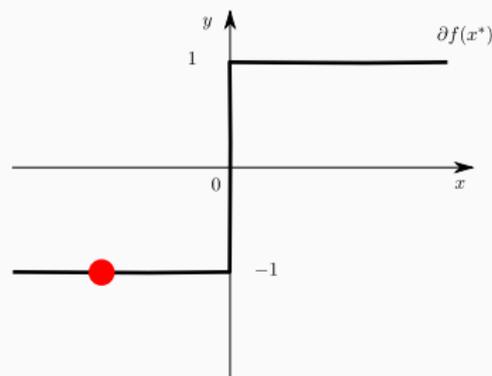
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



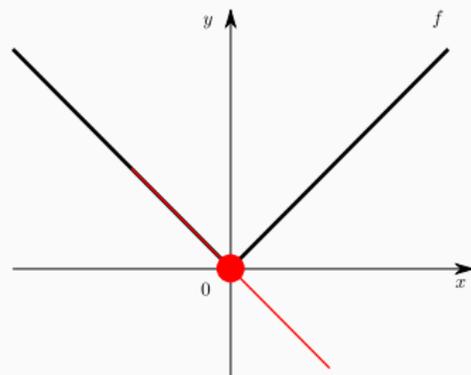
Sub-differential (sign)

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



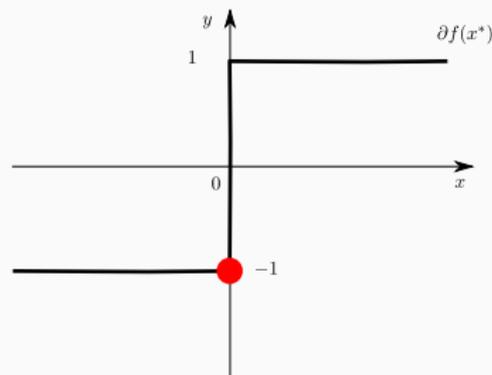
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



Sub-differential (sign)

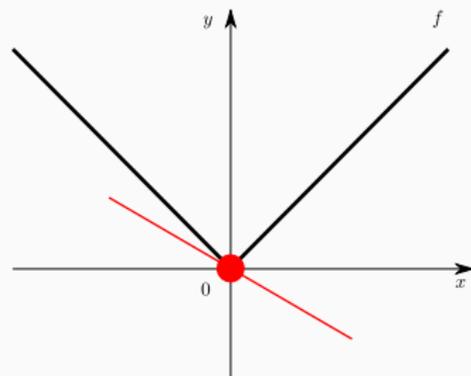
$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



Non-smooth convex optimization

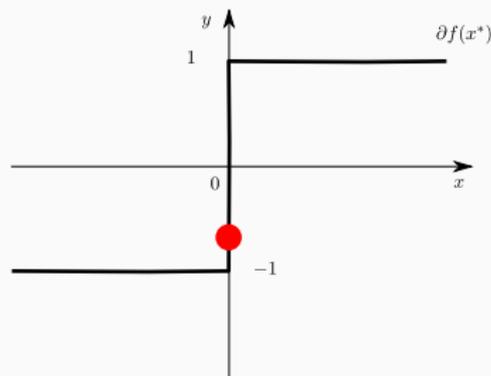
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



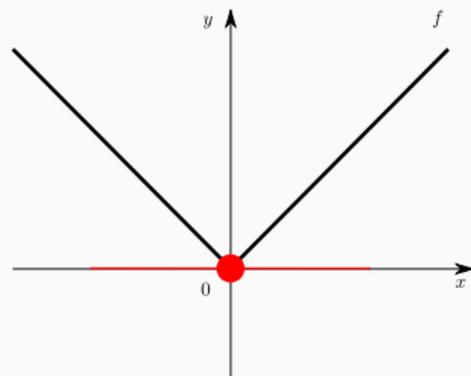
Sub-differential (sign)

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



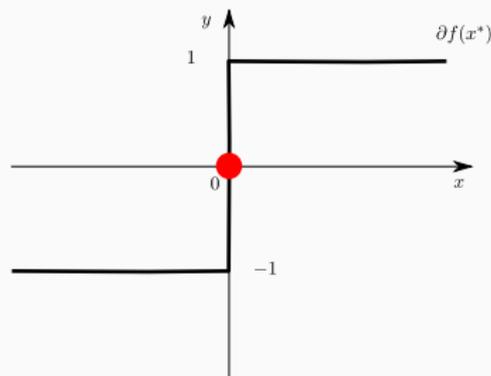
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



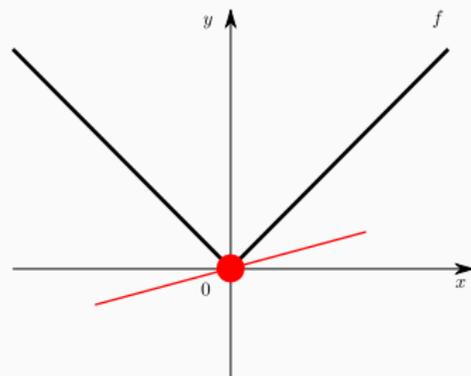
Sub-differential (sign)

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



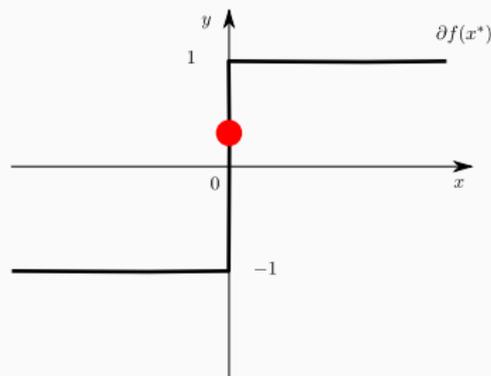
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



Sub-differential (sign)

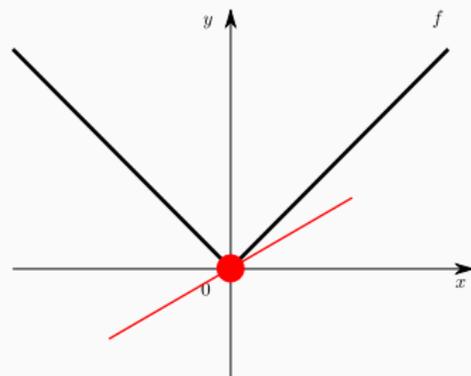
$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



Non-smooth convex optimization

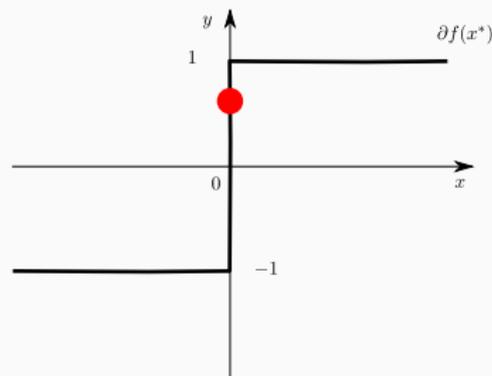
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



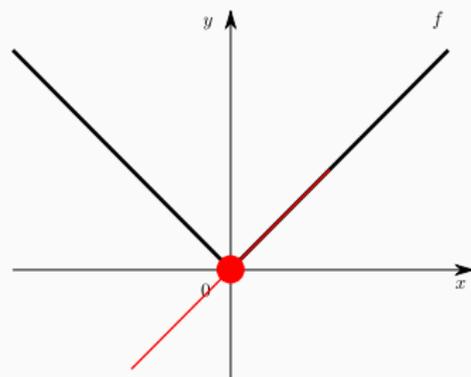
Sub-differential (sign)

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



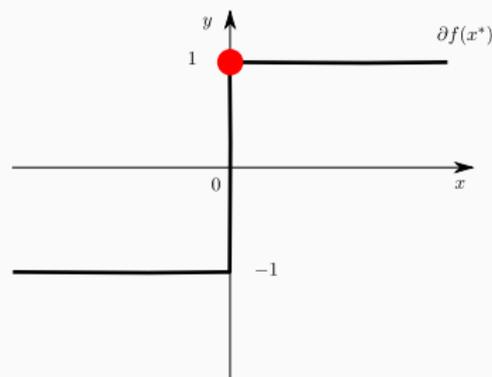
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



Sub-differential (sign)

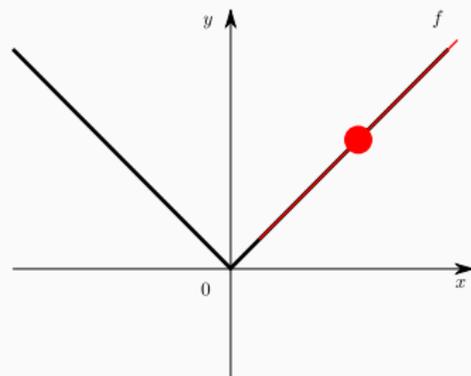
$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



Non-smooth convex optimization

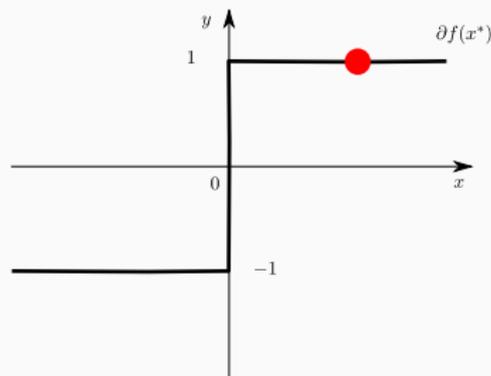
Function (abs):

$$f : \begin{cases} \mathbb{R} & \rightarrow \mathbb{R} \\ x & \mapsto |x| \end{cases}$$



Sub-differential (sign)

$$\partial f(x^*) = \begin{cases} \{-1\} & \text{if } x^* \in (-\infty, 0) \\ \{+1\} & \text{if } x^* \in (0, \infty) \\ [-1, 1] & \text{if } x^* = 0 \end{cases}$$



Proximal operator

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function (+ some technical conditions). The **proximal operator** of f is

$$\text{Prox}_f(x) = \underset{z \in \mathbb{R}^n}{\text{argmin}} \frac{1}{2} \|z - x\|_2^2 + f(z)$$

- Remark: this minimization problem always has a unique solution, so the proximal operator is without ambiguity a function $\mathbb{R}^n \rightarrow \mathbb{R}^n$.
- Always non-expansive:

$$\| \text{Prox}_f(x_1) - \text{Prox}_f(x_2) \| \leq \|x_1 - x_2\|$$

- Can be interpreted as a denoiser/shrinkage for the regularity f .

Property

$$\text{Prox}_{\gamma f}(x) = (\text{Id} + \gamma \partial f)^{-1} x$$

Proof.

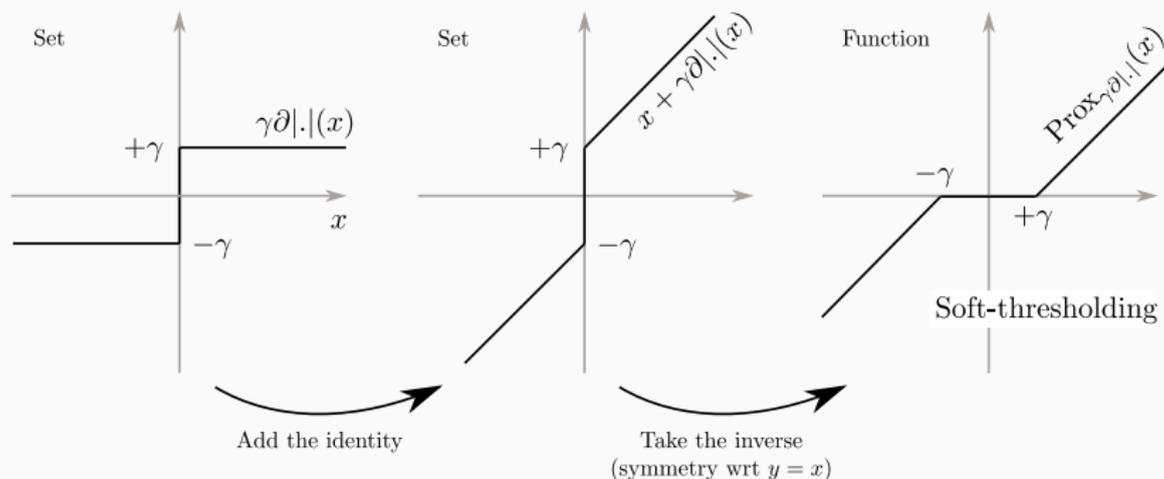
$$\begin{aligned} \operatorname{argmin}_z \frac{1}{2} \|z - x\|_2^2 + \gamma f(z) &\Leftrightarrow 0 \in \partial \left[\frac{1}{2} \|z - x\|_2^2 + \gamma f(z) \right] \\ &\Leftrightarrow 0 \in \partial \left[\frac{1}{2} \|z - x\|_2^2 \right] + \gamma \partial f(z) \\ &\Leftrightarrow 0 \in z - x + \gamma \partial f(z) \\ &\Leftrightarrow x \in z + \gamma \partial f(z) \\ &\Leftrightarrow x \in (\operatorname{Id} + \gamma \partial f)(z) \\ &\Leftrightarrow z = (\operatorname{Id} + \gamma \partial f)^{-1} x \end{aligned}$$

Even though $\partial f(x)$ is a set, the pre-image by $\operatorname{Id} + \gamma \partial f$ is unique. □

Soft-thresholding

$$\text{Prox}_{\gamma|\cdot|}(x) = \underset{z \in \mathbb{R}}{\text{argmin}} \frac{1}{2}(z - x)^2 + \gamma|z|$$

$$= (\text{Id} + \gamma\partial|\cdot|)^{-1}x = \begin{cases} x - \gamma & \text{if } x > \gamma \\ x + \gamma & \text{if } x < -\gamma \\ 0 & \text{otherwise} \end{cases}$$



Proximal operator of simple functions

Name	$f(x)$	$\text{Prox}_{\gamma f}(x)$
Indicator of convex set \mathcal{C}	$\begin{cases} 0 & \text{if } x \in \mathcal{C} \\ \infty & \text{otherwise} \end{cases}$	$\text{Proj}_{\mathcal{C}}(x)$
Square	$\frac{1}{2}\ x\ _2^2$	$\frac{x}{1+\gamma}$
Abs	$\ x\ _1$	$\text{Soft-T}(x, \gamma)$
Euclidean	$\ x\ _2$	$\left(1 - \frac{\gamma}{\max(\ x\ _2, \gamma)}\right)x$
Square+Affine	$\frac{1}{2}\ Ax + b\ _2^2$	$(\text{Id} + \gamma A^*A)^{-1}(x - \gamma A^*b)$
Separability for $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$	$g(x_1) + h(x_2)$	$\begin{pmatrix} \text{Prox}_{\gamma g}(x_1) \\ \text{Prox}_{\gamma h}(x_2) \end{pmatrix}$

More exhaustive list: <http://proximity-operator.net>

Proximal minimization

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a convex function (+ some technical conditions). Then, whatever the initialization x^0 and $\gamma > 0$, the sequence

$$x^{k+1} = \text{Prox}_{\gamma f}(x^k)$$

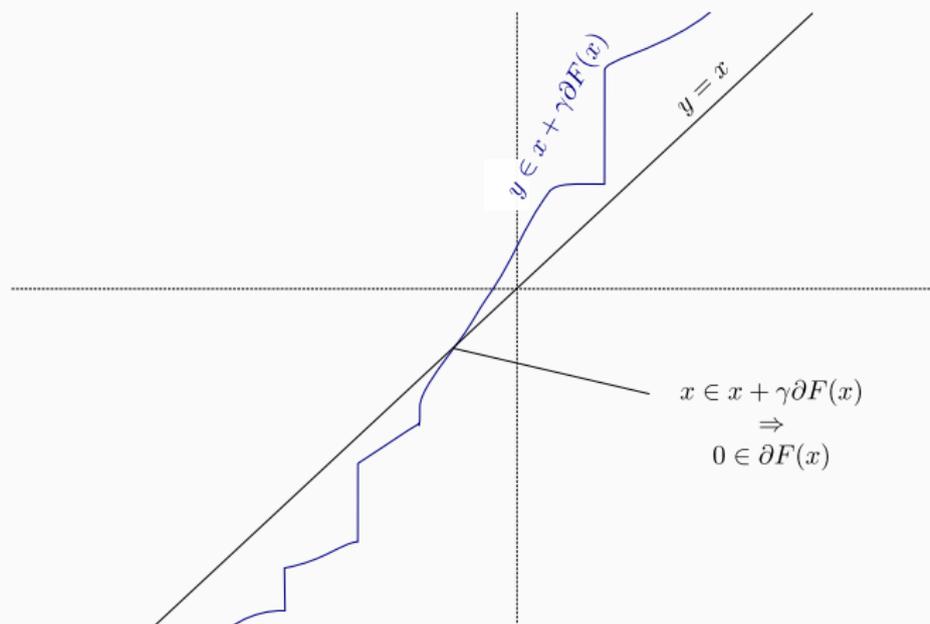
converges towards a **global minimizer** of f .

$$\text{Prox}_{\gamma f}(x^k) = (\text{Id} + \gamma \partial f)^{-1} x^k = \underset{z}{\operatorname{argmin}} \frac{1}{2} \|z - x^k\|_2^2 + \gamma f(z)$$

Compared to gradient descent

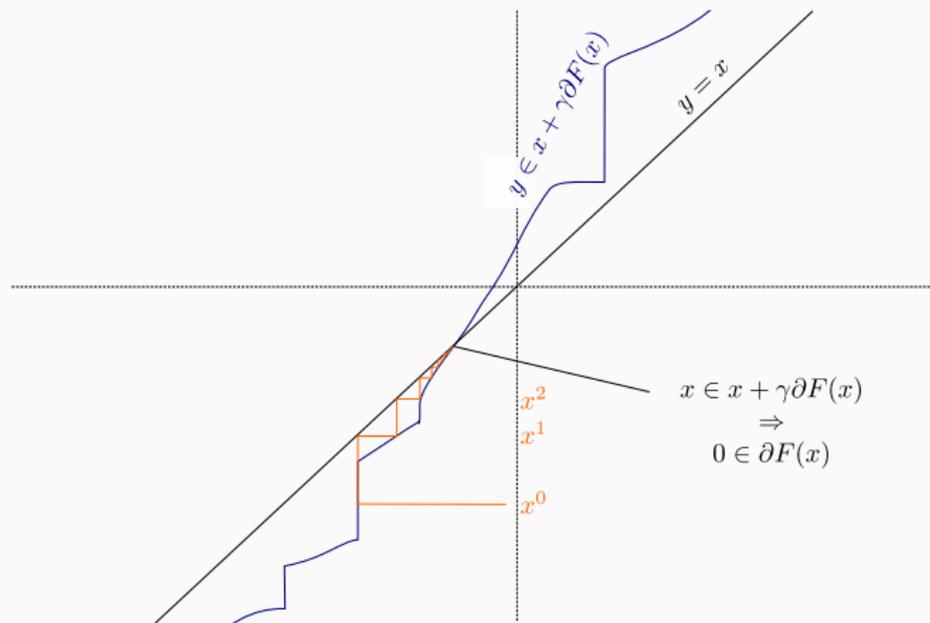
- No need to be differentiable,
- No need to have Lipschitz gradient,
- Works whatever the parameter γ ,
- Requires to solve an optimization problem at each step.

Non-smooth convex optimization



Proximal minimization:
Look at the set $x + \gamma \partial F(x)$

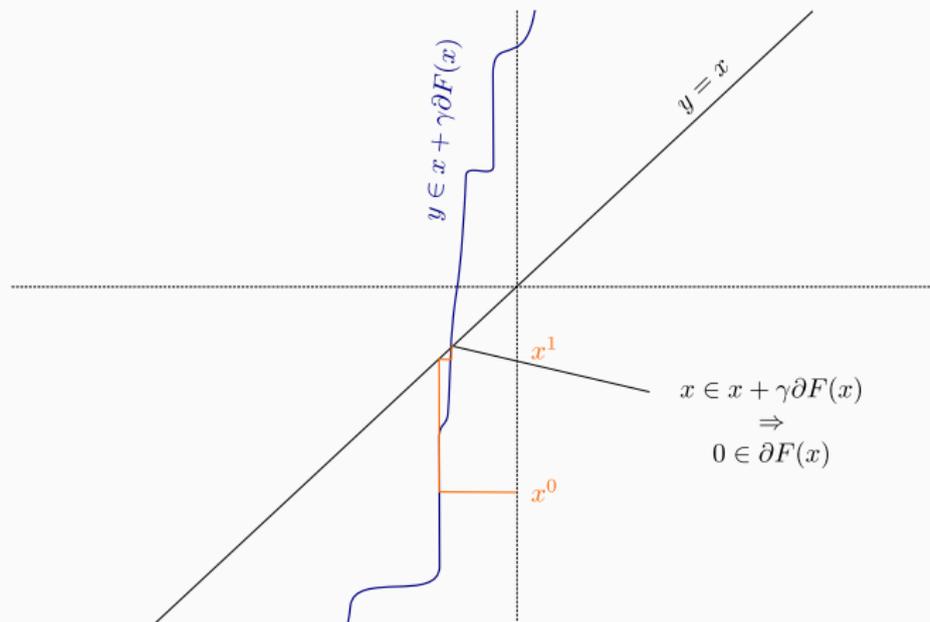
Non-smooth convex optimization



Proximal minimization:

read x^k on the y-axis and evaluate its pre-image by $x + \gamma \nabla F(x)$.

Non-smooth convex optimization



Proximal minimization:

the larger γ the faster, but the inversion becomes harder (ill-conditioned).

Toy example

- Consider the smoothing regularization problem

$$F(x) = \frac{1}{2} \|\nabla x\|_{2,2}^2$$

- Its sub-gradient is thus given by

$$\partial F(x) = \{\nabla F(x) = -\Delta x\}$$

- The proximal minimization reads as

$$\begin{aligned} x^{k+1} &= (\text{Id} + \gamma \partial F)^{-1} x^k \\ &= (\text{Id} - \gamma \Delta)^{-1} x^k \end{aligned}$$

- This is exactly the implicit Euler scheme for the Heat equation.

Can we apply proximal minimization for the LASSO?

Proximal splitting methods

- The proximal operator may not have a closed form.
- Computing it may be as difficult as solving the original problem ☹️
- Solution: use **proximal splitting methods**, a family of techniques developed for non-smooth convex problems.
- Idea: split the problem into subproblems, that involve
 - gradient descent steps for smooth terms,
 - proximal steps for simple convex terms.

$$\min_{x \in \mathbb{R}^n} \{E(x) = F(x) + G(x)\}$$

Proximal forward-backward algorithm

- Assume F is convex and differentiable with **L -Lipschitz gradient**

$$\|\nabla F(x_1) - \nabla F(x_2)\|_2 \leq L\|x_1 - x_2\|_2, \quad \text{for all } x_1, x_2 .$$

- Assume G is convex and **simple**, i.e., its prox is known in closed form

$$\text{Prox}_{\gamma G}(x) = \underset{z}{\operatorname{argmin}} \frac{1}{2}\|z - x\|_2^2 + \gamma G(z)$$

- The **proximal forward-backward algorithm** reads

$$x^{k+1} = \text{Prox}_{\gamma G}(x^k - \gamma \nabla F(x^k))$$

- For $0 < \gamma < 2/L$, it converges to a minimizer of $E = F + G$.

Aka, explicit-implicit scheme by analogy with PDE discretization schemes.

Non-smooth convex optimization

The LASSO problem:
$$E(\eta) = \underbrace{\frac{1}{2}\|y - A\eta\|_2^2}_{F(\eta)} + \underbrace{\tau\|\eta\|_1}_{G(\eta)=\sum_i |\eta_i|}, \quad A = HD$$

Iterative Soft-Thresholding Algorithm (ISTA) (Daubechies, 2004)

- F is convex and differentiable with **L -Lipschitz gradient**

$$\nabla F(\eta) = A^*(A\eta - y) \quad \text{with} \quad L = \|A\|_2^2$$

- G is convex and **simple**, in fact separable:

$$\text{Prox}_{\gamma G}(\eta)_i = \text{Soft-T}(\eta_i, \gamma\tau)$$

- The proximal forward-backward algorithm reads for $0 < \gamma < 2/L$

$$\eta^{k+1} = \text{Soft-T}(\eta^k - \gamma(A^*A\eta^k - A^*y), \gamma\tau)$$

and is known as **Iterative Soft-Thresholding Algorithm (ISTA)**.

- Finally:

$$\hat{x}^* = \bar{D}\hat{\eta}^*$$

Preconditioned ISTA (1/2)

- Remark

$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - A\eta\|_2^2 + \tau \|\eta\|_1, \quad A = H \underbrace{\bar{W}^+ \Lambda^{1/2}}_D$$

$$\in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - H\bar{W}^+ \Lambda^{1/2} \eta\|_2^2 + \tau \|\eta\|_1$$

- $\Lambda^{1/2}$ invertible: bijection between $z = \Lambda^{1/2} \eta$ and $\eta = \Lambda^{-1/2} z$
- Solving for η is equivalent to solve a weighted LASSO for z

$$\hat{z}^* \in \operatorname{argmin}_{z \in \mathbb{R}^K} \frac{1}{2} \|y - H\bar{W}^+ z\|_2^2 + \tau \|\Lambda^{-1/2} z\|_1$$

$$\in \operatorname{argmin}_{z \in \mathbb{R}^K} \frac{1}{2} \|y - Bz\|_2^2 + \sum_{i=1}^K \frac{\tau}{\lambda_i} |z_i|, \quad B = H\bar{W}^+$$

- In practice, this equivalent problem has **better conditioning**.

Equivalent to:
$$E(z) = \underbrace{\frac{1}{2} \|y - Bz\|_2^2}_{F(z)} + \underbrace{\tau \|\Lambda^{-1/2} z\|_1}_{G(z) = \sum_i \frac{\tau}{\lambda_i} |z_i|}, \quad B = H\bar{W}^+$$

Preconditioned ISTA (2/2)

$$\nabla F(z) = B^*(Bz - y) \quad \text{with} \quad L = \|B\|_2^2$$

$$\text{Prox}_{\gamma G}(z)_i = \text{Soft-T} \left(z_i, \frac{\gamma\tau}{\lambda_i} \right)$$

- ISTA becomes for $0 < \gamma < 2/L$

$$z^{k+1} = \text{Soft-T} \left(z^k - \gamma(B^*Bz^k - B^*y), \frac{\gamma\tau}{\lambda_i} \right)$$

- Finally:

$$\hat{x}^* = \bar{W}^+ \hat{z}^*$$

- Leads to **larger steps γ , better conditioning, and faster convergence.**

$$\begin{aligned}z^{k+1} &= \text{Prox}_{\gamma G}(z^k - \gamma \nabla F(z^k)) \\ &= \text{Soft-T} \left(z^k - \gamma(B^* B z^k - B^* y), \frac{\gamma \tau}{\lambda_i} \right) \quad \text{with } B = H \bar{W}^+\end{aligned}$$

Bredies & Lorenz (2007): $E(z^k) - E(z^*)$ decays with rate $O(1/k)$

Fast ISTA (FISTA)

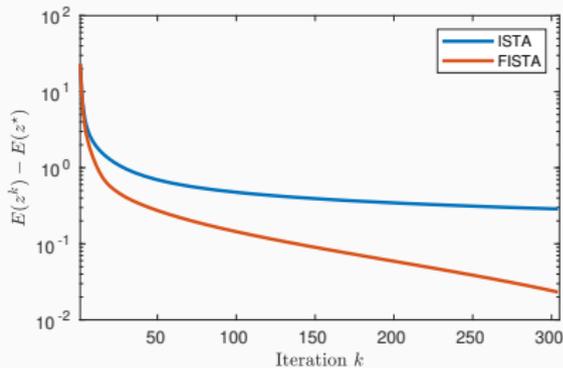
$$\begin{aligned}z^{k+1} &= \text{Prox}_{\gamma G}(\tilde{z}^k - \gamma \nabla F(\tilde{z}^k)) \\ \tilde{z}^{k+1} &= z^{k+1} + \frac{t_k - 1}{t_{k+1}}(z^{k+1} - z^k) \\ t_{k+1} &= \frac{1 + \sqrt{1 + 4t_k^2}}{2}, t_0 = 1\end{aligned}$$

Beck & Teboulle (2009): $E(z^k) - E(z^*)$ decays with rate $O(1/k^2)$

Non-smooth convex optimization



(a) Input y : motion blur + noise ($\sigma = 2$)



(b) Convergence profiles

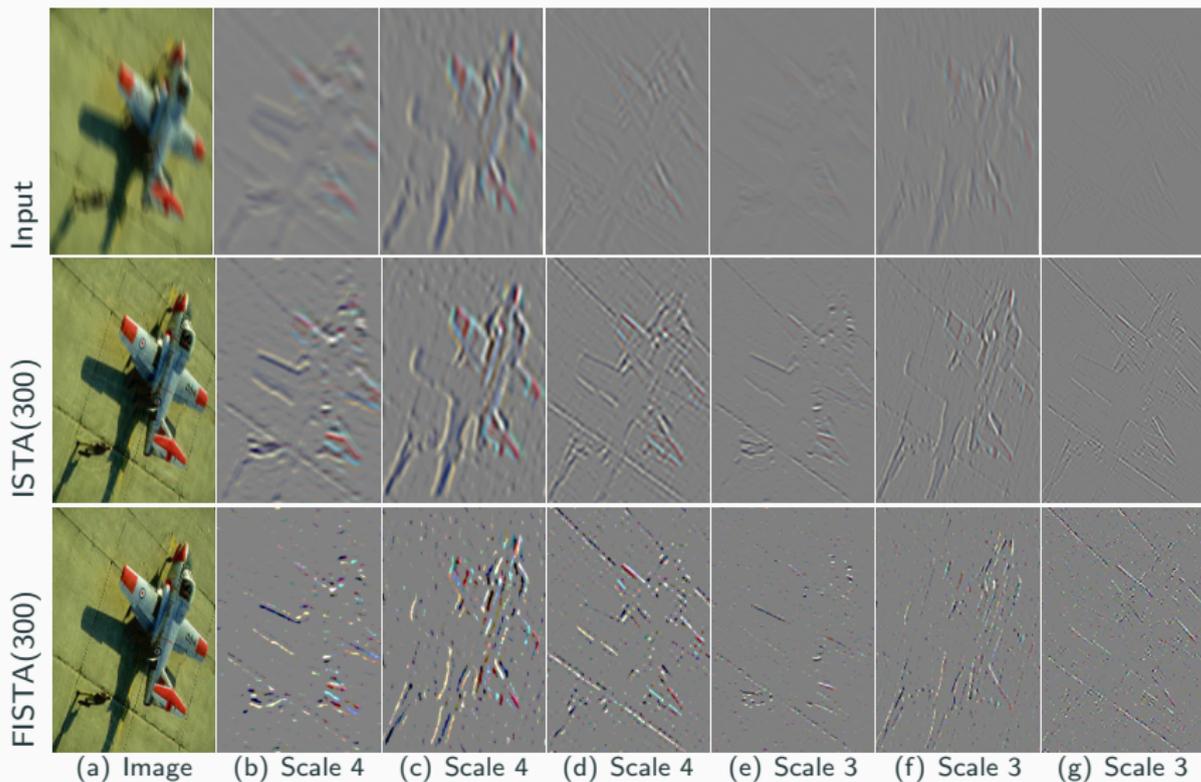


(c) Deconvolution ISTA(300)+UDWT



(d) Deconvolution FISTA(300)+UDWT

Non-smooth convex optimization



FISTA converges faster: sparser codes given a limited time budget

Sparsity: synthesis vs analysis

Sparse synthesis model with UDWT

- LASSO:
$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - H\bar{W}^+ \Lambda^{1/2} \eta\|_2^2 + \tau \|\eta\|_1$$
- Using the change of variable $\eta = \Lambda^{-1/2} z$:

$$\hat{z}^* \in \operatorname{argmin}_{z \in \mathbb{R}^K} \frac{1}{2} \|y - H\bar{W}^+ z\|_2^2 + \tau \|\Lambda^{-1/2} z\|_1$$

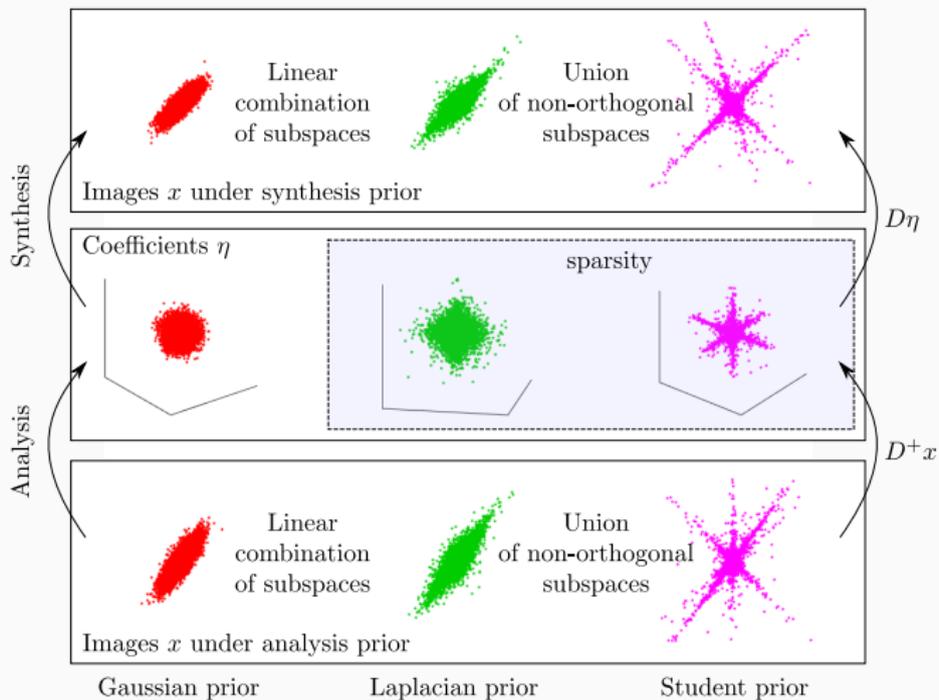
Sparse analysis model with UDWT

- What about?

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - Hx\|_2^2 + \tau \|\Lambda^{-1/2} \bar{W} x\|_1$$

- The change of variable $\eta = \Lambda^{-1/2} \bar{W} x$ is not one-to-one.
- **The two problems are not equivalent** (unless \bar{W} is invertible).

Sparse reconstruction: synthesis vs analysis



Analysis versus synthesis (Elad, Milanfar, Rubinstein, 2007)

Generative: generate good images

$$\hat{x}^* = \mathbf{D}\hat{\eta}^* \quad \text{with} \quad \hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - \mathbf{H}\mathbf{D}\eta\|_2^2 + \tau \|\eta\|_p^p, \quad p \geq 0$$

Synthesis: images are linear combinations of a few columns of \mathbf{D} .

Bayesian interpretation: MAP for the sparse code η .

Discriminative: discriminate between good and bad images

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_p^p, \quad p \geq 0$$

Analysis: images are correlated with a few rows of $\mathbf{\Gamma}$.

Bayesian interpretation: MAP for x with an improper *Gibbs prior*.

Sparse reconstruction: synthesis vs analysis

$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - \mathbf{H}\mathbf{D}\eta\|_2^2 + \tau \|\eta\|_p^p \quad (\ell_p^p\text{-synthesis})$$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_p^p \quad (\ell_p^p\text{-analysis})$$

Common properties

	Solution	Problem
$p = 0$	Optimal sparse	Non-convex & discontinuous (NP-hard)
$0 < p < 1$	Sparse	Non-convex & continuous but non-smooth
$p = 1$	Sparse	Convex & continuous but non-smooth
$p > 1$	Smooth	Convex & differentiable
$p = 2$	Linear	Quadratic

- $\mathbf{\Gamma}$ square and invertible \Rightarrow equivalent for $\mathbf{D} = \mathbf{\Gamma}^{-1}$.
- $\mathbf{\Gamma}$ full-rank and $p = 2 \Rightarrow$ equivalent for $\mathbf{D} = \mathbf{\Gamma}^+$.
- LTI dictionaries \Rightarrow redundant filter bank.

Sparse reconstruction: synthesis vs analysis

$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - \mathbf{H}\mathbf{D}\eta\|_2^2 + \tau \|\eta\|_p^p \quad (\ell_p^p\text{-synthesis})$$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\Gamma x\|_p^p \quad (\ell_p^p\text{-analysis})$$

Synthesis

- \mathbf{D} : synthesis dictionary.
- Atoms need to span images.
 - ⇒ Low- & high-pass filters
 - ⇒ $\operatorname{Im}[\mathbf{D}] \approx \mathbb{R}^n$
- Redundancy favor sparsity.
- K dimensional problem ($> n$).
- Prior separable.

Analysis

- Γ : analysis dictionary.
- Atoms need to sparsify images.
 - ⇒ High-pass filters only
 - ⇒ $\operatorname{Ker}[\Gamma] \neq \emptyset$ (\supset DC, coarse)
- Redundancy decreases sparsity.
- n dimensional problem ($< K$).
- Prior non-separable.

Quiz: What analysis dictionary is LTI and not too redundant?

Sparse reconstruction: synthesis vs analysis

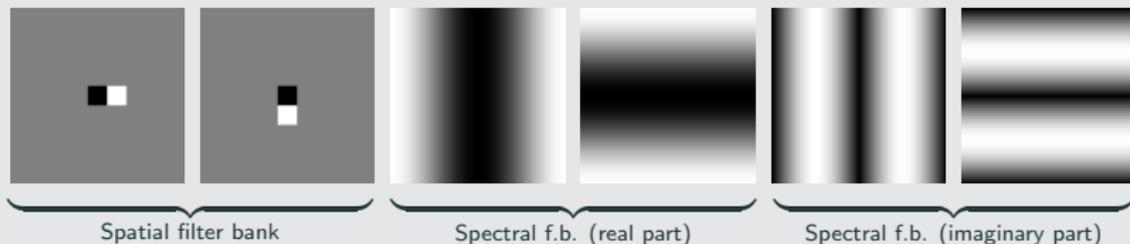
$$\hat{\eta}^* \in \operatorname{argmin}_{\eta \in \mathbb{R}^K} \frac{1}{2} \|y - \mathbf{H}\mathbf{D}\eta\|_2^2 + \tau \|\eta\|_p^p \quad (\ell_p^p\text{-synthesis})$$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_p^p \quad (\ell_p^p\text{-analysis})$$

Link between analysis models and variational methods

- $p = 2$: Analysis model = Tikhonov regularization.
- $p = 1$ & $\mathbf{\Gamma} = \nabla$: Analysis model = anisotropic Total-Variation (TV)

TV filter bank = Horizontal and vertical gradient



Can we use proximal forward-backward for ℓ_1 -analysis prior?

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \underbrace{\frac{1}{2} \|y - \mathbf{H}x\|_2^2}_{F(x)} + \underbrace{\tau \|\Gamma x\|_1}_{G(x)} \quad (\ell_1\text{-analysis})$$

Proximal forward-backward for the ℓ_1 -analysis problem?

- F convex and differentiable
- G convex but **not simple** (not separable)

→ cannot use proximal forward backward ☹

- **Exception:** for denoising $\mathbf{H} = \operatorname{Id}_n$ (see: Chambolle algorithm, 2004)

Need another proximal optimization technique.

$$\min_{x \in \mathbb{R}^n} \{E(x) = F(x) + G(x)\}$$

Alternating direction method of multipliers (ADMM) (~1970)

- Assume F and G are **convex and simple** (+ some mild conditions).
- For any initialization x^0, \tilde{x}^0 and d^0 , the ADMM algorithm reads as

$$x^{k+1} = \text{Prox}_{\gamma F}(\tilde{x}^k + d^k)$$

$$\tilde{x}^{k+1} = \text{Prox}_{\gamma G}(x^{k+1} - d^k)$$

$$d^{k+1} = d^k - x^{k+1} + \tilde{x}^{k+1}$$

- For $\gamma > 0$, x^k converges to a minimizer of $E = F + G$.

Fast version: FADMM, similar idea as for FISTA (Goldstein *et al.*, 2014).

Related concepts: Lagrange multipliers, Duality, Legendre transform.

How to use it for ℓ_1 analysis priors?

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_1 \quad (\ell_1\text{-analysis})$$

ADMM + Variable splitting (1/3)

- Define: $X = \begin{pmatrix} x \\ z \end{pmatrix} \in \mathbb{R}^{n+K}$

- Consider: $E(X) = F(X) + G(X)$

with:

$$\begin{cases} F \begin{pmatrix} x \\ z \end{pmatrix} = \|y - \mathbf{H}x\|_2^2 + \tau \|z\|_1 \\ G \begin{pmatrix} x \\ z \end{pmatrix} = \begin{cases} 0 & \text{if } \mathbf{\Gamma}x = z \\ \infty & \text{otherwise} \end{cases} \end{cases}$$

- Remark 1: Minimizing E solves the ℓ_1 -analysis problem.

- Remark 2: F and G are convex and simple \Rightarrow ADMM applies.

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_1 \quad (\ell_1\text{-analysis})$$

ADMM + Variable splitting (2/3)

Applying formula from slide 92:

$$F \begin{pmatrix} x \\ z \end{pmatrix} = \|y - \mathbf{H}x\|_2^2 + \tau \|z\|_1 \quad \longrightarrow \quad \operatorname{Prox}_{\gamma F} \begin{pmatrix} x \\ z \end{pmatrix} = \begin{pmatrix} (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (x + \gamma \mathbf{H}^* y) \\ \operatorname{Soft-T}(z, \gamma \tau) \end{pmatrix}$$

$$G \begin{pmatrix} x \\ z \end{pmatrix} = \begin{cases} 0 & \text{if } \mathbf{\Gamma}x = z \\ \infty & \text{otherwise} \end{cases} \quad \longrightarrow \quad \operatorname{Prox}_{\gamma G} \begin{pmatrix} x \\ z \end{pmatrix} = \underbrace{\begin{pmatrix} \operatorname{Id}_n \\ \mathbf{\Gamma} \end{pmatrix}}_{\text{Projection on } \mathcal{C}} (\operatorname{Id}_n + \mathbf{\Gamma}^* \mathbf{\Gamma})^{-1} (x + \mathbf{\Gamma}^* z)$$

Indicator of the convex set
 $\mathcal{C} = \{(x, z) ; \mathbf{\Gamma}x = z\}$
Projection on \mathcal{C}

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\mathbf{\Gamma}x\|_1 \quad (\ell_1\text{-analysis})$$

ADMM + Variable splitting (3/3)

$$x^{k+1} = (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (\tilde{x}^k + d_x^k + \gamma \mathbf{H}^* y)$$

$$z^{k+1} = \operatorname{Soft-T}(\tilde{z}^k + d_z^k, \gamma \tau)$$

$$\tilde{x}^{k+1} = (\operatorname{Id}_n + \mathbf{\Gamma}^* \mathbf{\Gamma})^{-1} (x^{k+1} - d_x^k + \mathbf{\Gamma}^* (z^{k+1} - d_z^k))$$

$$\tilde{z}^{k+1} = \mathbf{\Gamma} \tilde{x}^{k+1}$$

$$d_x^{k+1} = d_x^k - x^{k+1} + \tilde{x}^{k+1}$$

$$d_z^{k+1} = d_z^k - z^{k+1} + \tilde{z}^{k+1}$$

If \mathbf{H} is a blur, and $\mathbf{\Gamma}$ a filter bank, $(\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1}$ and $(\operatorname{Id}_n + \mathbf{\Gamma}^* \mathbf{\Gamma})^{-1}$ can be computed in the Fourier domain in $O(n \log n)$.

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\Gamma x\|_1 \quad (\ell_1\text{-analysis})$$

Application to Total-Variation

$$\Gamma = \nabla$$

$$x^{k+1} = (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (\tilde{x}^k + d_x^k + \gamma \mathbf{H}^* y)$$

$$z^{k+1} = \operatorname{Soft-T}(\tilde{z}^k + d_z^k, \gamma \tau)$$

$$\tilde{x}^{k+1} = (\operatorname{Id}_n + \nabla^* \nabla)^{-1} (x^{k+1} - d_x^k + \nabla^* (z^{k+1} - d_z^k))$$

$$\tilde{z}^{k+1} = \nabla \tilde{x}^{k+1}$$

$$d_x^{k+1} = d_x^k - x^{k+1} + \tilde{x}^{k+1}$$

$$d_z^{k+1} = d_z^k - z^{k+1} + \tilde{z}^{k+1}$$

$$\nabla^* = -\operatorname{div} \quad \text{and} \quad \nabla^* \nabla = -\Delta$$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\Gamma x\|_1 \quad (\ell_1\text{-analysis})$$

Application to Total-Variation

$$\Gamma = \nabla$$

$$x^{k+1} = (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (\tilde{x}^k + d_x^k + \gamma \mathbf{H}^* y)$$

$$z^{k+1} = \operatorname{Soft-T}(\tilde{z}^k + d_z^k, \gamma \tau)$$

$$\tilde{x}^{k+1} = (\operatorname{Id}_n - \Delta)^{-1} (x^{k+1} - d_x^k - \operatorname{div}(z^{k+1} - d_z^k))$$

$$\tilde{z}^{k+1} = \nabla \tilde{x}^{k+1}$$

$$d_x^{k+1} = d_x^k - x^{k+1} + \tilde{x}^{k+1}$$

$$d_z^{k+1} = d_z^k - z^{k+1} + \tilde{z}^{k+1}$$

$$\nabla^* = -\operatorname{div} \quad \text{and} \quad \nabla^* \nabla = -\Delta$$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\Gamma x\|_1 \quad (\ell_1\text{-analysis})$$

Application to sparse analysis with UDWT

$$\Gamma = \Lambda^{-1/2} \bar{\mathbf{W}}$$

$$x^{k+1} = (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (\tilde{x}^k + d_x^k + \gamma \mathbf{H}^* y)$$

$$z^{k+1} = \operatorname{Soft-T}(\tilde{z}^k + d_z^k, \frac{\gamma \tau}{\lambda_i})$$

$$\tilde{x}^{k+1} = (\operatorname{Id}_n + \bar{\mathbf{W}}^* \bar{\mathbf{W}})^{-1} (x^{k+1} - d_x^k + \bar{\mathbf{W}}^* (z^{k+1} - d_z^k))$$

$$\tilde{z}^{k+1} = \bar{\mathbf{W}} \tilde{x}^{k+1}$$

$$d_x^{k+1} = d_x^k - x^{k+1} + \tilde{x}^{k+1}$$

$$d_z^{k+1} = d_z^k - z^{k+1} + \tilde{z}^{k+1}$$

Tight-frame: $\bar{\mathbf{W}}^* \bar{\mathbf{W}} = \operatorname{Id}_n$

$$\hat{x}^* \in \operatorname{argmin}_{x \in \mathbb{R}^n} \frac{1}{2} \|y - \mathbf{H}x\|_2^2 + \tau \|\Gamma x\|_1 \quad (\ell_1\text{-analysis})$$

Application to sparse analysis with UDWT

$$\Gamma = \Lambda^{-1/2} \bar{\mathbf{W}}$$

$$x^{k+1} = (\operatorname{Id}_n + \gamma \mathbf{H}^* \mathbf{H})^{-1} (\tilde{x}^k + d_x^k + \gamma \mathbf{H}^* y)$$

$$z^{k+1} = \operatorname{Soft-T}(\tilde{z}^k + d_z^k, \frac{\gamma \tau}{\lambda_i})$$

$$\tilde{x}^{k+1} = \frac{1}{2} (x^{k+1} - d_x^k + \bar{\mathbf{W}}^* (z^{k+1} - d_z^k))$$

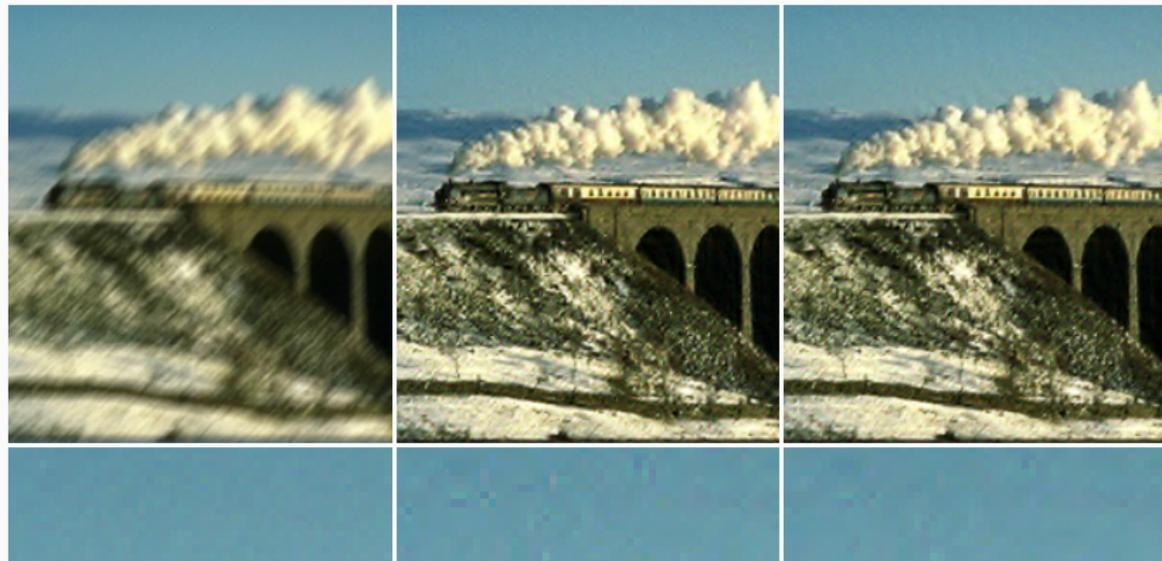
$$\tilde{z}^{k+1} = \bar{\mathbf{W}} \tilde{x}^{k+1}$$

$$d_x^{k+1} = d_x^k - x^{k+1} + \tilde{x}^{k+1}$$

$$d_z^{k+1} = d_z^k - z^{k+1} + \tilde{z}^{k+1}$$

Tight-frame: $\bar{\mathbf{W}}^* \bar{\mathbf{W}} = \operatorname{Id}_n$

Deconvolution with UDWT (5 levels, Db2)

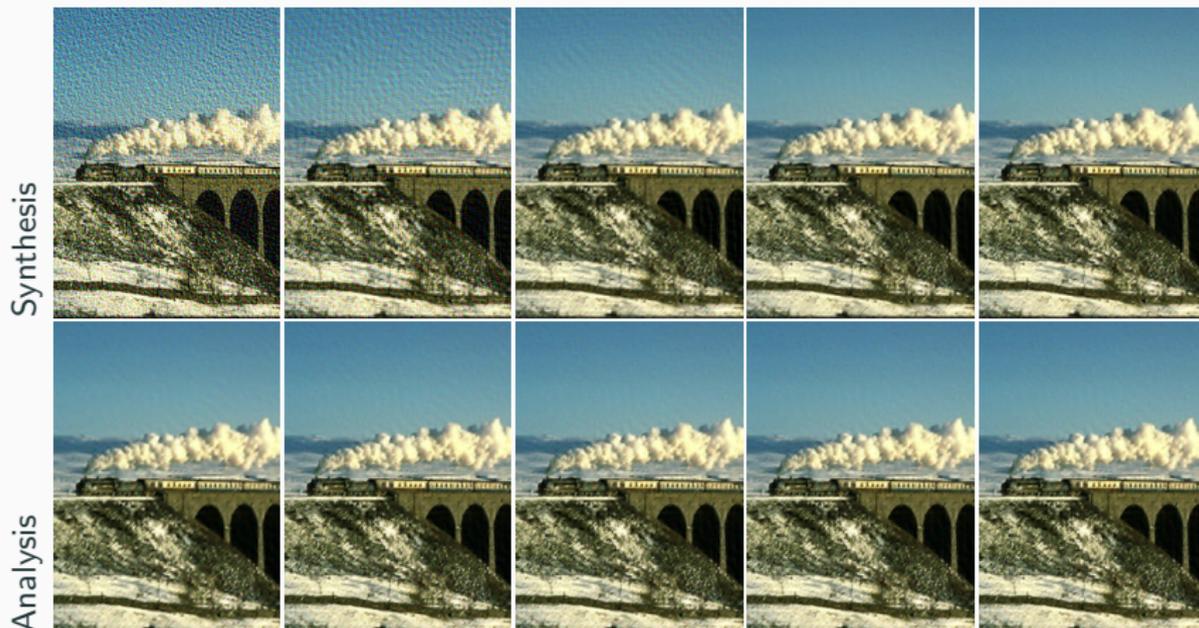


(a) Blurry image y (noise $\sigma = 2$)

(b) Synthesis (FISTA)

(c) Analysis (FADMM)

Sparse analysis – Results



Analysis allows for less decomposition levels.
⇒ leads to faster algorithms.

Sparse analysis – Results



(a) Noisy ($\sigma = 40$) (b) Analysis UDWT(4) (c) +block (orien.+col.) (d) Difference

- As for TV, **group coefficients** across orientations/color using $\ell_{2,1}$ norms:

$$\|\Gamma z\|_{2,1}$$

- The soft-thresholding becomes the **group soft-thresholding**:

$$[\text{Prox}_{\gamma\|\cdot\|_{2,1}}(z)]_i = \begin{cases} z_i - \gamma \frac{z_i}{\|z_i\|_2} & \text{if } \|z_i\|_2 > \gamma \\ 0 & \text{otherwise} \end{cases}$$

Reminder from last class:

Modeling the distribution of images is complex (large degree of freedom).
Applying LMMSE on patches \rightarrow increase performance

Next class:

**What if we use sparse priors, not for the distribution of images,
but for the distribution of patches?**

For further reading

Sparsity, shrinkage and recovery guarantee:

- Donoho & Johnstone (1994); Moulin & Liu (1999); Donoho and Elad (2003); Gribonval and Nielsen (2003); Candès and Tao (2005); Zhang (2008); Candès and Romberg (2007).
- Book: Statistical Learning with Sparsity (Hastie, Tibshirani, Wainwright, 2015).

Wavelet related transforms:

- Warblet/Chirplet (Mann, Mihovilovic et al., 1991–1992), Curvelet (Candès & Donoho, 2000), Noiselet (Coifman, 2001), Contourlet (Do & Vetterli, 2002), Ridgelet (Do & Vetterli, 2003), Shearlets (Kanghui et al., 2005), Bandelet (Le Pennec, Peyré, Mallat, 2005), Empirical wavelets (Gilles, 2013).
- Book: A wavelet tour of signal processing (Mallat, 2008)

Non-smooth convex optimization:

- Douglas-Rachford splitting (Combettes & Pesquet, 2007), Split Bregman (Goldstein & Osher, 2009), Primal-Dual (Chambolle & Pock, 2011), Generalized FB (Raguet et al., 2013), Condat algorithm (2014).
- Book: Convex Optimization (Boyd, 2004).

Questions?

Next class: Patch models and dictionary learning

Sources, images courtesy and acknowledgment

L. Condat

A. Horodniceanu

G. Peyré

J. Salmon

Wikipedia