

EE 5561: Image Processing and Applications

Lecture 11

Mehmet Akçakaya

Recap of Last Lecture

- General Inverse Problems
 - Least squares solution/maximum likelihood estimation in Gaussian noise
 - Regularization
 - Regularized least squares
 - Tikhonov regularization & variants
 - Energy-based regularizers
- Today:
 - More regularization based on sparsity

Sparsity

$$\arg \min_{\mathbf{x}} ||\mathbf{y} - \mathbf{H}\mathbf{x}||_2^2 + \psi(\mathbf{x})$$

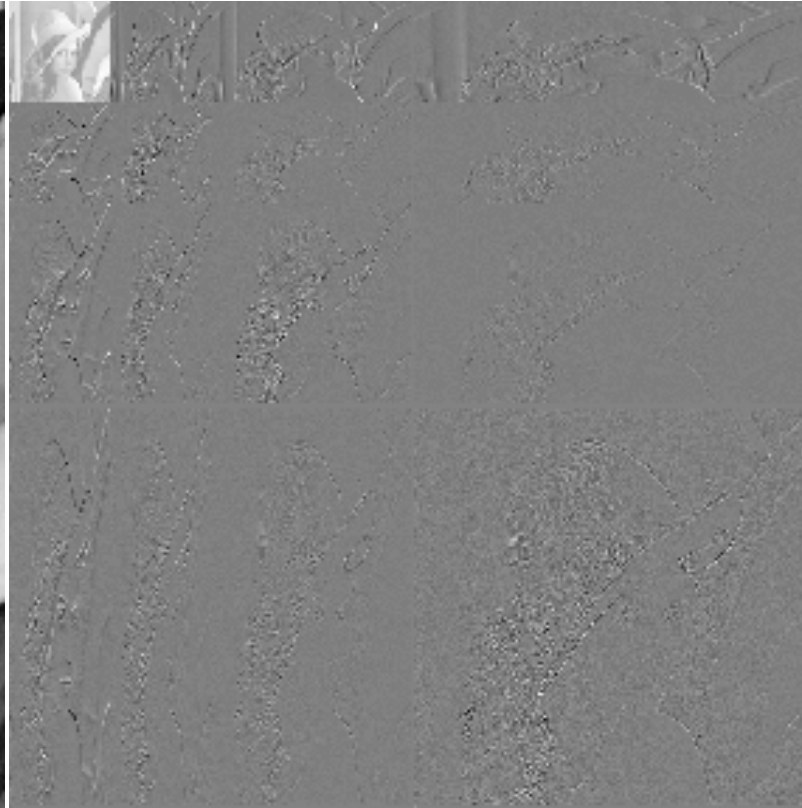
- Recall energy compaction for DCT/wavelets
- What happens when we throw away low-intensity transform coefficients?
 - Main idea behind compression
- This suggests a new regularization idea
 - Find \mathbf{x} that “best fits” the data and has the least number of transform coefficients
- Let \mathbf{W} be the transform, and $\mathbf{x} = \mathbf{W}\boldsymbol{\theta}$
- We define the l_0 “norm” as
$$||\boldsymbol{\theta}||_0 = \text{number of non-zero coefficients of } \boldsymbol{\theta}$$

Haar Wavelet

X



3-levels: $W_N X W_N^T$

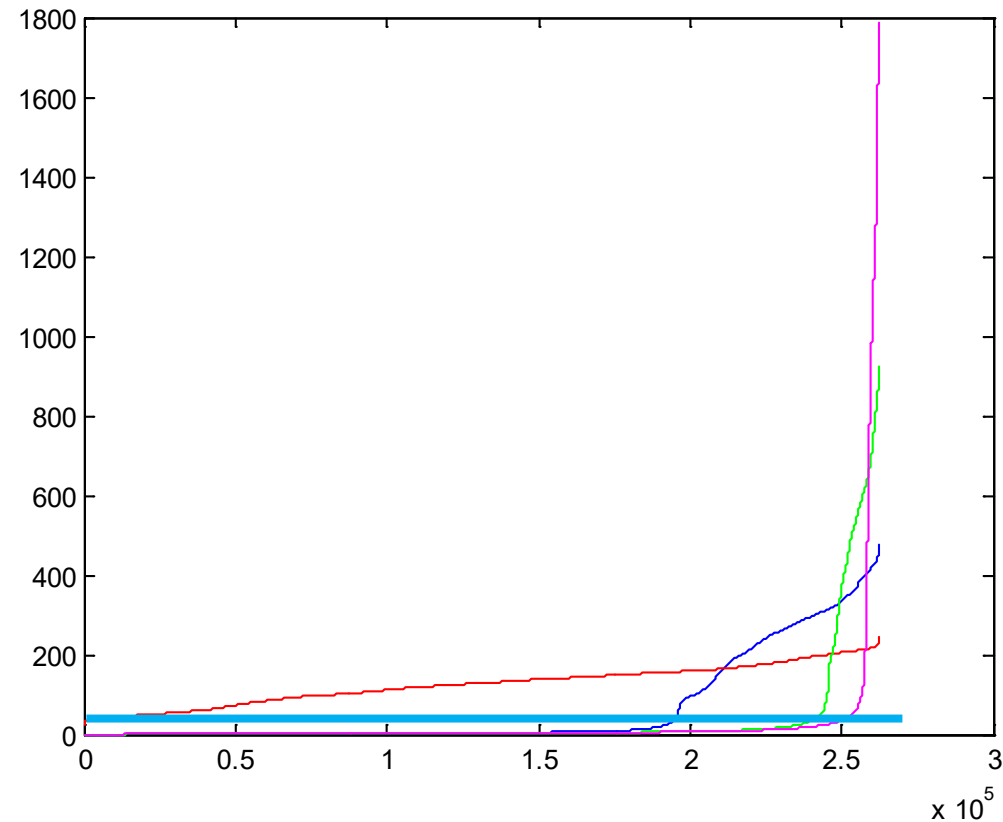


Haar Wavelet

Energy “compaction”:

sorted (abs value) image (red) vs. wavelet

1-level (blue), 2-level (green), 3-level (magenta) coefficients across the image



Wavelets

Original



Compression with Wavelet
(10-fold)



Compression with Wavelet
(20-fold)



Sparsity

- Goal: Find $\hat{\boldsymbol{\theta}}$ that minimizes

$$\min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{HW}\boldsymbol{\theta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\theta}\|_0 \leq K$$

for some $K \in \mathbb{Z}^+$



$$\min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|_0 \quad \text{s. t.} \quad \|\mathbf{y} - \mathbf{HW}\boldsymbol{\theta}\|_2^2 \leq D$$

for some $D \in \mathbb{R}^+$

- $\|\boldsymbol{\theta}\|_0$ is hard to work with

Sparsity

- The best convex approximation is $\|\boldsymbol{\theta}\|_1 = \sum_{i=1}^N |\theta_i|$
 - The reason why this is sparsity promoting has a lengthy history
 - We will try to give a high-level intuition
 - First consider the following two definitions

ℓ_p norm: $\|\boldsymbol{\theta}\|_p = (\sum_{i=1}^N |\theta_i|^p)^{\frac{1}{p}}$ This is actually a norm for $p \geq 1$

ℓ_p ball: $\{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_p = 1\}$

- Now let's see how these ℓ_p balls look for different values of p



$p = \infty$



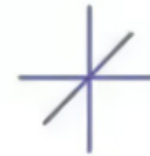
$p = 2$



$p = 1$



$0 < p < 1$

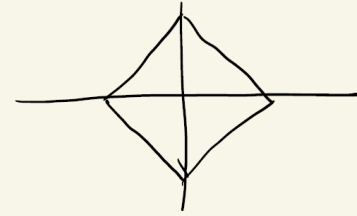


$p = 0$

Sparsity

- In particular

$$\ell_1\text{-ball} \quad \{\theta: \|\theta\|_1 = 1\}$$
$$|\theta_1| + |\theta_2| = 1$$



"diamond"

- Now let's go back to our problem, but in the noiseless scenario, i.e.

$$\mathbf{y} = \mathbf{H}\mathbf{W}\boldsymbol{\theta} \triangleq \mathbf{A}\boldsymbol{\theta}$$

- As we discussed earlier, this is an interesting problem when \mathbf{A} is either ill-conditioned or under-determined
- We will consider the latter case, i.e.

$$\mathbf{A} \in \mathbb{R}^{M \times N} \quad M < N$$

Sparsity

- Suppose we solve

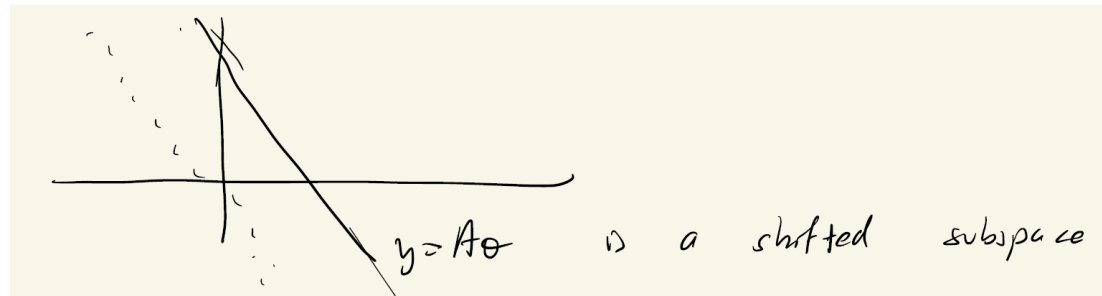
$$\min_{\boldsymbol{\theta}} \|\boldsymbol{\theta}\|_p \quad \text{s. t.} \quad \mathbf{y} = \mathbf{A}\boldsymbol{\theta} \quad \text{for } p \in \{0, 1, 2\}$$

- First look at all solutions to $\mathbf{y} = \mathbf{A}\boldsymbol{\theta}$

- How does this work?

- Find any solution $\boldsymbol{\theta}'$ such that $\mathbf{y} = \mathbf{A}\boldsymbol{\theta}'$
- Then all the solutions lie in a shifted subspace given by $\boldsymbol{\theta}' + N(\mathbf{A})$
- Why? Take any $\mathbf{w} \in N(\mathbf{A})$, then $\mathbf{A}\mathbf{w} = \mathbf{0}$
- $(\mathbf{w} + \boldsymbol{\theta}')$ is another solution since $\mathbf{A}(\mathbf{w} + \boldsymbol{\theta}') = \mathbf{A}\mathbf{w} + \mathbf{A}\boldsymbol{\theta}' = \mathbf{0} + \mathbf{y} = \mathbf{y}$

- Pictorially

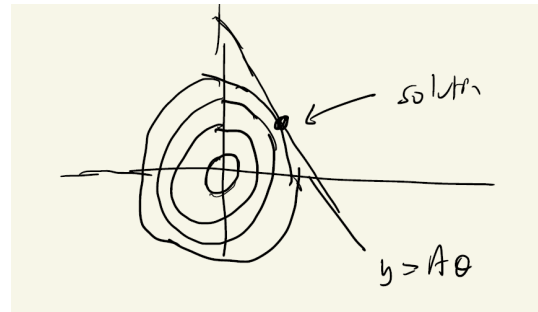


Sparsity

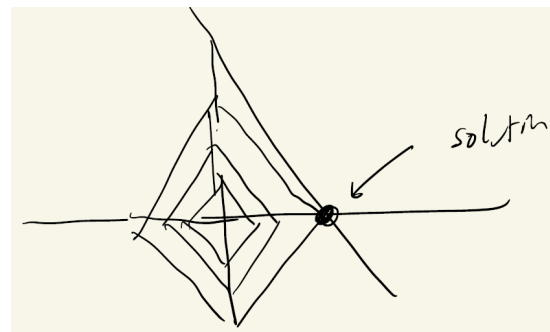
- For $p = 2$, the solution to

$$\min_{\theta} \|\theta\|_2 \quad \text{s. t.} \quad y = A\theta$$

is the intersection of a scaled l_2 ball (i.e. sphere) and this shifted subspace



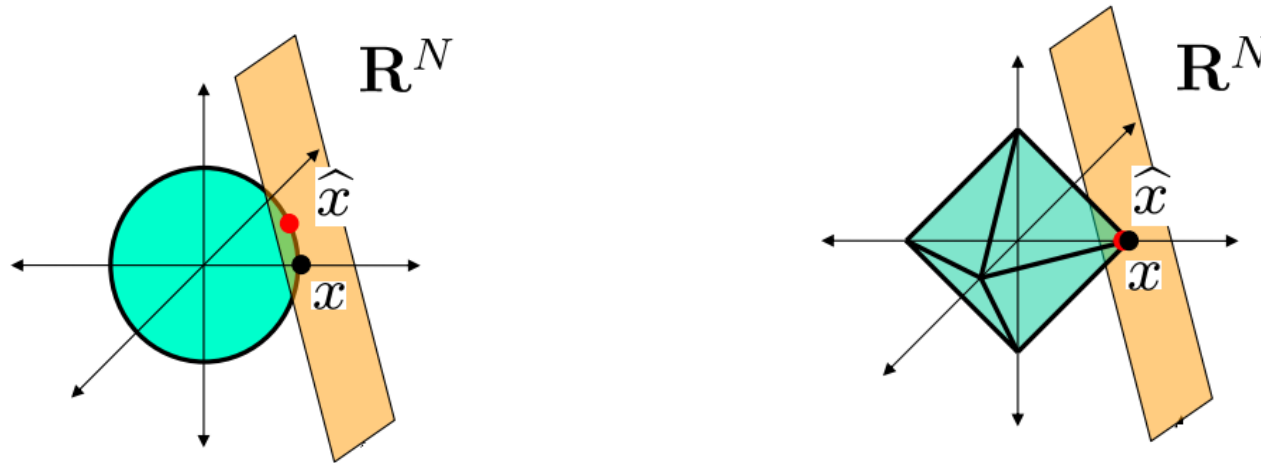
- For $p = 1$, it's the intersection of a “diamond” and this shifted subspace



Note this ends up on the axes
→ same as the l_0 solution

Sparsity

- In high dimensions (M, N large), with appropriate choice of \mathbf{A} , the null-space is oriented randomly
- Pictorially



- Same idea here: l_1 solution overlaps with the l_0 solution
- More dramatically at higher dimensions

Sparsity

- Easier to perform sparsity regularization in an unconstrained manner (i.e. in Lagrangian form)

$$\arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{H}\mathbf{W}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1$$

- Questions:
 - How to choose \mathbf{W} ?
 - How to solve this objective function?
 - Better ways than the l_1 norm of transform-domain coefficients for using sparsity?

Solving Sparsity-Regularized LS

- First let's consider denoising again, i.e. $\mathbf{H} = \mathbf{I}$, or

$$\mathbf{y} = \mathbf{x} + \mathbf{n}, \quad \mathbf{x} = \mathbf{W}\boldsymbol{\theta}, \quad \boldsymbol{\theta} \text{ “sparse”}$$

i.e. $\mathbf{y} = \mathbf{W}\boldsymbol{\theta} + \mathbf{n}$

↖
image domain

↗
transform domain

- For now assume, \mathbf{W} is orthogonal (e.g. wavelet or DCT), then

$$\mathbf{W}^T \mathbf{y} = \boldsymbol{\theta} + \mathbf{W}^T \mathbf{n}$$

- Recall if $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}) \Rightarrow \mathbf{W}^T \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{W}^T \mathbf{W}) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$
- Thus without loss of generality, let

$$\mathbf{y}' = \mathbf{W}^T \mathbf{y}, \quad \mathbf{n}' = \mathbf{W}^T \mathbf{n}$$

- We will solve $\mathbf{y}' = \boldsymbol{\theta} + \mathbf{n}'$

Denoising with Sparsity

- Here everything is in transform domain
- For ease of notation, we will just use $\mathbf{y} = \boldsymbol{\theta} + \mathbf{n}$
- Hence we are solving the following objective function (for denoising)

$$\begin{aligned} & \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1 \\ &= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \sum_{k=1}^N |y_k - \theta_k|^2 + \lambda \sum_{k=1}^N |\theta_k| \\ &= \arg \min_{\boldsymbol{\theta}} \sum_{k=1}^N \left(\frac{1}{2} |y_k - \theta_k|^2 + \lambda |\theta_k| \right) \end{aligned}$$

- Hence we can minimize this separately for each k

Denoising by Thresholding

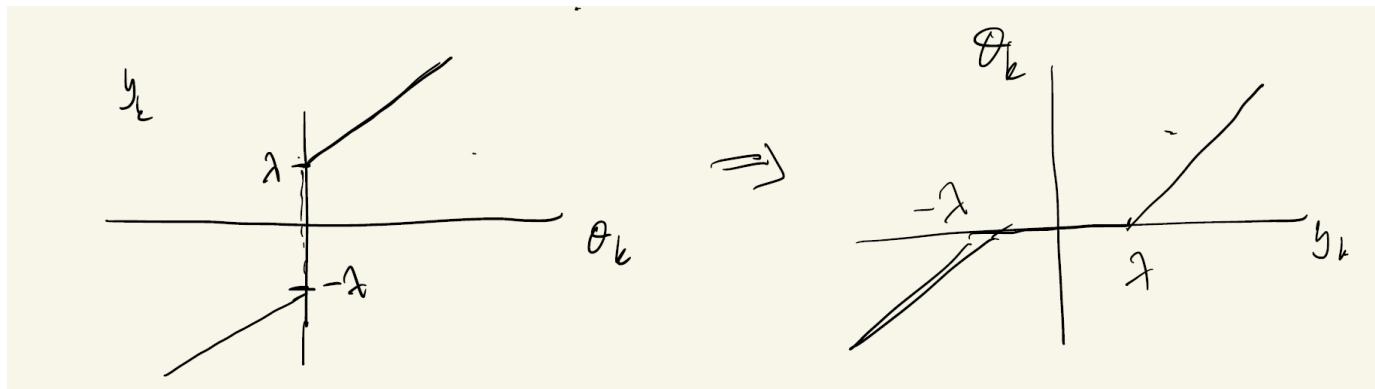
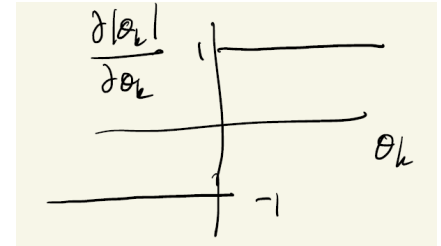
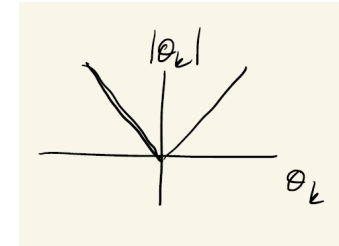
- For each k , we are solving

$$\min_{\theta_k} \frac{1}{2} |y_k - \theta_k|^2 + \lambda |\theta_k|$$

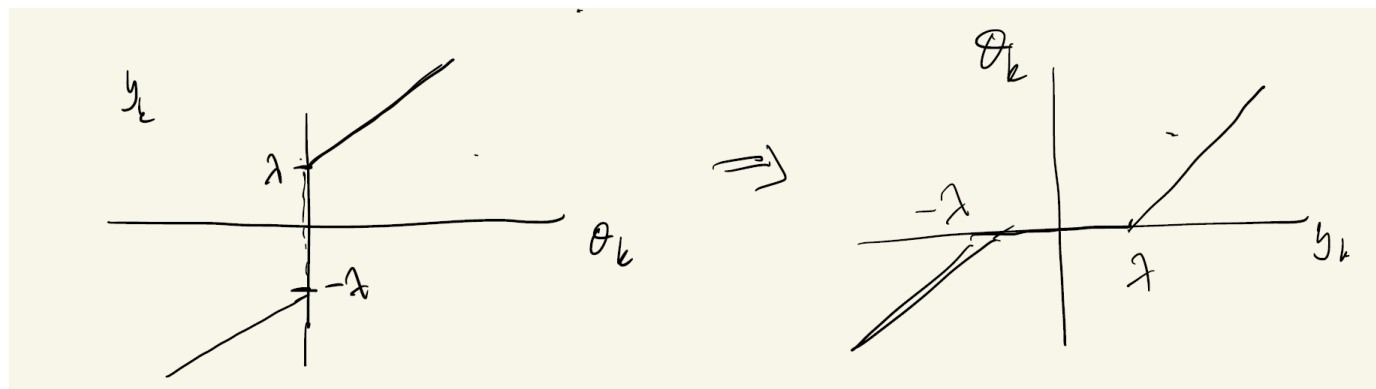
- Take derivative with respect to θ_k and set to 0

$$(\theta_k - y_k) + \lambda \operatorname{sign}(\theta_k) = 0$$

$$y_k = \theta_k + \lambda \operatorname{sign}(\theta_k)$$



Denoising by Thresholding



$$\hat{\theta}_k = \begin{cases} y_k - \lambda & y_k > \lambda \\ 0 & |y_k| < \lambda \\ y_k + \lambda & y_k < -\lambda \end{cases} = \begin{cases} |y_k| - \lambda & y_k > \lambda \\ 0 & |y_k| < \lambda \\ -(|y_k| - \lambda) & y_k < -\lambda \end{cases}$$

$$\hat{\theta}_k = \max(|y_k| - \lambda, 0) \cdot \text{sign}(y_k) \triangleq S_\lambda(y_k)$$

soft-thresholding function

Denoising by Thresholding

Original



Noisy



Soft-thresholding



Hard-thresholding



Solving Sparsity-Regularized LS

- How to choose λ ?
 - For denoising there are certain techniques, e.g. Stein's unbiased risk estimate (SURE)
 - With more complicated $\mathbf{H} \rightarrow$ usually empirical

– **Recap:**
$$\arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\mathbf{y} - \boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1$$

- Solution is
$$\hat{\theta}_k = \max(|y_k| - \lambda, 0) \cdot \text{sign}(y_k) \triangleq S_\lambda(y_k)$$

Solving Sparsity-Regularized LS

- Now let's go back to

$$\arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1 \quad (\mathbf{A} = \mathbf{H}\mathbf{W})$$

- This is of the form (more generally)

$$\arg \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) \quad \text{with } f(\cdot) \text{ convex}$$

- No closed form solution necessarily... Simplest solution is gradient descent

- Start at $\boldsymbol{\theta}^{(0)}$, move in the negative direction of the gradient

$$\boldsymbol{\theta}^{(k)} = \boldsymbol{\theta}^{(k-1)} - \eta_k \nabla f(\boldsymbol{\theta}^{(k-1)}) \quad k: \text{iteration number}$$

- Then

$$f(\boldsymbol{\theta}^{(0)}) > f(\boldsymbol{\theta}^{(1)}) > \dots$$

For appropriate choices of η_k
(+ some details)

Gradient Descent

- Not easy in our case
 - Derivative of the l_1 norm is not continuous
- Let's reinterpret gradient descent
 - Consider the quadratic approximation

$$f(\mathbf{z}) = f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{z} - \boldsymbol{\theta} \rangle + \frac{1}{2\eta} \|\mathbf{z} - \boldsymbol{\theta}\|_2^2$$

← replaces Taylor series

- Then

$$\begin{aligned} \min_{\mathbf{z}} f(\mathbf{z}) &= \min_{\mathbf{z}} \left\{ f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{z} - \boldsymbol{\theta} \rangle + \frac{1}{2\eta} \|\mathbf{z} - \boldsymbol{\theta}\|_2^2 \right\} \\ &= \min_{\mathbf{z}} \left\{ 2\eta \langle \nabla f(\boldsymbol{\theta}), \mathbf{z} \rangle + \|\mathbf{z}\|_2^2 - 2\langle \boldsymbol{\theta}, \mathbf{z} \rangle \right\} \end{aligned}$$

Gradient Descent

$$\begin{aligned}\min_{\mathbf{z}} f(\mathbf{z}) &= \min_{\mathbf{z}} \left\{ f(\boldsymbol{\theta}) + \langle \nabla f(\boldsymbol{\theta}), \mathbf{z} - \boldsymbol{\theta} \rangle + \frac{1}{2\eta} \|\mathbf{z} - \boldsymbol{\theta}\|_2^2 \right\} \\ &= \min_{\mathbf{z}} \left\{ \|\mathbf{z}\|_2^2 - 2\langle \boldsymbol{\theta} - \eta \nabla f(\boldsymbol{\theta}), \mathbf{z} \rangle \right\} \\ &= \min_{\mathbf{z}} \left\{ \|\mathbf{z}\|_2^2 - 2\langle \boldsymbol{\theta} - \eta \nabla f(\boldsymbol{\theta}), \mathbf{z} \rangle + \underbrace{\|\boldsymbol{\theta} - \eta \nabla f(\boldsymbol{\theta})\|_2^2}_{\text{does not depend on } \mathbf{z}} \right\} \\ &= \min_{\mathbf{z}} \left\| \mathbf{z} - (\boldsymbol{\theta} - \eta \nabla f(\boldsymbol{\theta})) \right\|_2^2 \\ &\quad \swarrow \text{same as our gradient step}\end{aligned}$$

Proximal Gradient Descent

- Now we will use this to get at a method called the proximal gradient descent

$$\arg \min_{\boldsymbol{\theta}} \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 + \lambda \|\boldsymbol{\theta}\|_1$$

is (also) of the form

$$\arg \min_{\boldsymbol{\theta}} g(\boldsymbol{\theta}) + h(\boldsymbol{\theta}) \quad (\text{both convex})$$

- We do the quadratic approximation on g , but not on h

$$\min_{\boldsymbol{\theta}} \left\{ g(\mathbf{z}) + \langle \nabla g(\mathbf{z}), \boldsymbol{\theta} - \mathbf{z} \rangle + \frac{1}{2\eta} \|\boldsymbol{\theta} - \mathbf{z}\|_2^2 + h(\boldsymbol{\theta}) \right\}$$

$$= \min_{\boldsymbol{\theta}} \frac{1}{2\eta} \left\| \boldsymbol{\theta} - (\mathbf{z} - \eta \nabla g(\mathbf{z})) \right\|_2^2 + h(\boldsymbol{\theta})$$

Proximal Gradient Descent

- This objective function looks like denoising with $h(\cdot)$ as regularizer

$$\min_{\boldsymbol{\theta}} \frac{1}{2\eta} \left\| \boldsymbol{\theta} - (\mathbf{z} - \eta \nabla g(\mathbf{z})) \right\|_2^2 + h(\boldsymbol{\theta})$$

- Proximal operator is exactly this!

- It's the solution to

$$\text{prox}_{h,\eta}(\mathbf{u}) = \arg \min_{\boldsymbol{\theta}} \frac{1}{2\eta} \left\| \boldsymbol{\theta} - \mathbf{u} \right\|_2^2 + h(\boldsymbol{\theta})$$

- Thus, we set iterations as

$$\boldsymbol{\theta}^{(t)} = \text{prox}_{h,\eta}(\boldsymbol{\theta}^{(t-1)} - \eta \nabla g(\boldsymbol{\theta}^{(t-1)}))$$

- Also note we know how to solve the proximal operator for $h(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1$ (soft-thresholding)

Proximal Gradient Descent

- Overall strategy can be described in two steps

$$\mathbf{z}^{(t)} = \boldsymbol{\theta}^{(t-1)} - \eta \nabla g(\boldsymbol{\theta}^{(t-1)}) \quad \longleftarrow \text{data fidelity}$$

$$\boldsymbol{\theta}^{(t)} = \text{prox}_{h,\eta}(\mathbf{z}^{(t)}) = \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta} - \mathbf{z}^{(t)}\|_2^2 + \eta h(\boldsymbol{\theta}) \quad \longleftarrow \text{regularization/proximal}$$

- For our problem
$$\begin{array}{l} g(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{A}\boldsymbol{\theta}\|_2^2 \\ h(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_1 \end{array} \quad \longrightarrow \quad \begin{array}{l} \nabla g(\boldsymbol{\theta}) = -\mathbf{A}^T(\mathbf{y} - \mathbf{A}\boldsymbol{\theta}) \\ \text{prox}_{h,\eta}(\mathbf{z}) = S_\eta(\mathbf{z}) \end{array}$$

- Thus, overall

$$\mathbf{z}^{(t)} = \boldsymbol{\theta}^{(t-1)} + \eta \mathbf{A}^T(\mathbf{y} - \mathbf{A}\boldsymbol{\theta}^{(t-1)})$$

$$\boldsymbol{\theta}^{(t)} = S_\eta(\mathbf{z}^{(t)})$$