

Iteratively Reweighted ℓ_1 Approaches to Sparse Composite Regularization

Phil Schniter



THE OHIO STATE UNIVERSITY

Joint work with Rizwan Ahmad (OSU)

Supported in part by NSF grant CCF-1018368.

SAHD @ Duke — July 27, 2015

Outline

- 1 Introduction and Motivation for Composite Penalties
- 2 Co-L1 and its Interpretations
- 3 Co-IRW-L1 and its Interpretations
- 4 Numerical Experiments

Introduction

- Goal: Recover signal $\mathbf{x} \in \mathbb{C}^N$ from noisy linear measurements

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{w} \in \mathbb{C}^M$$

where possibly $M \ll N$.

- Approach: Solve optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} R(\mathbf{x}) \text{ s.t. } \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$$

with δ selected based on statistics of $\|\mathbf{w}\|_2$.

- Question: How to choose penalty/regularization $R(\mathbf{x})$?

Typical Choices of Penalty

Suppose $\Psi \mathbf{x}$ is (approximately) sparse for analysis operator $\Psi \in \mathbb{C}^{L \times N}$:

ℓ_0 penalty: $R(\mathbf{x}) = \|\Psi \mathbf{x}\|_0$

- Impractical: optimization problem is NP hard

ℓ_1 penalty (generalized LASSO): $R(\mathbf{x}) = \|\Psi \mathbf{x}\|_1$

- Tightest convex relaxation of ℓ_0 penalty
- Fast algorithms: Douglas-Rachford, NESTA-UP, MFISTA, GAMP ...

Many other penalties, such as $R(\mathbf{x}) = \|\Psi \mathbf{x}\|_p$ for $p \in (0, 1)$.

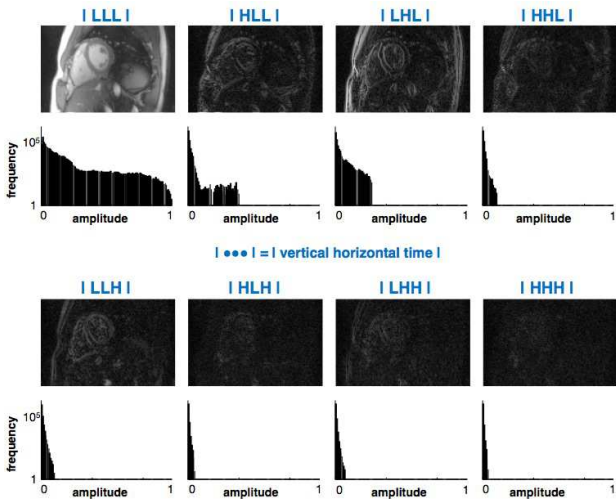
Choice of Analysis Operator

How to choose Ψ in practice?

- Maybe a wavelet dictionary? Which one?
- Maybe a concatenation of several dictionaries $\begin{bmatrix} \Psi_1 \\ \vdots \\ \Psi_D \end{bmatrix} ?$
- What if signal is more sparse in one dictionary than another?
Can we use this to our advantage?

Example: Undecimated Wavelet Transform of MRI Cine

Note different sparsity rate in each subband of 1-level UWT:



Composite ℓ_1 Penalties

We propose to use **composite ℓ_1 penalties** of the form

$$R(\mathbf{x}; \boldsymbol{\lambda}) \triangleq \sum_{d=1}^D \lambda_d \|\Psi_d \mathbf{x}\|_1, \quad \lambda_d \geq 0$$

where

- operators Ψ_d have unit-norm rows (but otherwise arbitrary),
- weights λ_d are **learned from the data**.

We propose two algorithms to jointly estimate \mathbf{x} and $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_D]^T$:

- 1 Composite- ℓ_1 minimization (**Co-L1**)
- 2 Iteratively reweighted composite- ℓ_1 minimization (**Co-IRW-L1**)

The Co-L1 Algorithm

- 1: input: $\{\Psi_d\}_{d=1}^D, \Phi, \mathbf{y}, \delta \geq 0, \epsilon \geq 0$
- 2: initialization: $\lambda_d^{(1)} = 1 \forall d$
- 3: for $t = 1, 2, 3, \dots$
- 4: $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x}} \sum_{d=1}^D \lambda_d^{(t)} \|\Psi_d \mathbf{x}\|_1$ s.t. $\|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$
- 5: $\lambda_d^{(t+1)} \leftarrow \frac{L_d}{\epsilon + \|\Psi_d \mathbf{x}^{(t)}\|_1}, d = 1, \dots, D$
- 6: end
- 7: output: $\mathbf{x}^{(t)}$

- leverages existing ℓ_1 solvers,
- applies to both real- and complex-valued cases,
- reduces to IRW-L1 algorithm [Candes,Wakin,Boyd'08] when $L_d = 1 \forall d$ (single-atom dictionaries).

The Co-IRW-L1 Algorithm

- 1: input: $\{\Psi_d\}_{d=1}^D, \Phi, \mathbf{y}, \delta \geq 0$,
- 2: if $\mathbf{x} \in \mathbb{R}^N$, use $\Lambda = (1, \infty)$ and the real version of $\log p(\mathbf{x}; \boldsymbol{\lambda}, \boldsymbol{\epsilon})$;
if $\mathbf{x} \in \mathbb{C}^N$, use $\Lambda = (2, \infty)$ and the complex version of $\log p(\mathbf{x}; \boldsymbol{\lambda}, \boldsymbol{\epsilon})$.
- 3: initialization: $\lambda_d^{(1)} = 1 \forall d, \mathbf{W}_d^{(1)} = \mathbf{I} \forall d$
- 4: for $t = 1, 2, 3, \dots$
- 5: $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x}} \sum_{d=1}^D \lambda_d^{(t)} \|\mathbf{W}_d^{(t)} \Psi_d \mathbf{x}\|_1$ s.t. $\|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$
- 6: $(\lambda_d^{(t+1)}, \epsilon_d^{(t+1)}) \leftarrow \arg \max_{\lambda_d \in \Lambda, \epsilon_d > 0} \log p(\mathbf{x}^{(t)}; \boldsymbol{\lambda}, \boldsymbol{\epsilon}), d = 1, \dots, D$
- 7: $\mathbf{W}_d^{(t+1)} \leftarrow \text{diag} \left\{ \frac{1}{\epsilon_d^{(t+1)} + |\boldsymbol{\psi}_{d,1}^\top \mathbf{x}^{(t)}|}, \dots, \frac{1}{\epsilon_d^{(t+1)} + |\boldsymbol{\psi}_{d,L_d}^\top \mathbf{x}^{(t)}|} \right\}, d = 1, \dots, D$
- 8: end
- 9: output: $\mathbf{x}^{(t)}$

- IRW version of Co-L1: tunes both λ_d and \mathbf{W}_d for all d .
- also tunes regularization parameters ϵ_d for all d .

Understanding Co-L1 and Co-IRW-L1

In the sequel, we provide four interpretations of each algorithm:

- 1 MM optimization of a particular **non-convex** penalty,
- 2 a particular approximation of ℓ_0 **minimization**,
- 3 **Bayesian** estimation according to a particular hierarchical prior,
- 4 **variational EM** algorithm under a particular prior.

Outline

- 1 Introduction and Motivation for Composite Penalties
- 2 Co-L1 and its Interpretations**
- 3 Co-IRW-L1 and its Interpretations
- 4 Numerical Experiments

Optimization Interpretations of Co-L1

Co-L1 is an MM approach to the **weighted log-sum** optimization problem

$$\arg \min_{\mathbf{x}} \sum_{d=1}^D L_d \log(\epsilon + \|\Psi_d \mathbf{x}\|_1) \text{ s.t. } \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta.$$

and

As $\epsilon \rightarrow 0$, Co-L1 aims to solve the **weighted $\ell_{1,0}$** problem

$$\arg \min_{\mathbf{x}} \sum_{d=1}^D L_d 1_{\|\Psi_d \mathbf{x}\|_1 > 0} \text{ s.t. } \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta.$$

Note: L_d is the size of dictionary Ψ_d , and 1_{\square} is the indicator function.

Bayesian Interpretations of Co-L1

As $\epsilon \rightarrow 0$, Co-L1 is an MM approach to **Bayesian MAP estimation** under an AWGN likelihood and the hierarchical prior

$$p(\mathbf{x}|\boldsymbol{\lambda}) = \prod_{d=1}^D \left(\frac{\lambda_d}{2}\right)^{L_d} \exp(-\lambda_d \|\Psi_d \mathbf{x}\|_1) \quad \text{i.i.d. Laplacian}$$

$$p(\boldsymbol{\lambda}) = \prod_{d=1}^D p(\lambda_d), \quad p(\lambda_d) \propto \begin{cases} \frac{1}{\lambda_d} & \lambda_d > 0 \\ 0 & \text{else} \end{cases}, \quad \text{Jeffrey's non-informative}$$

and

As $\epsilon \rightarrow 0$, Co-L1 is a **variational EM** approach to estimating (deterministic) $\boldsymbol{\lambda}$ under an AWGN likelihood and the prior

$$p(\mathbf{x}; \boldsymbol{\lambda}) = \prod_{d=1}^D \left(\frac{\lambda_d}{2}\right)^{L_d} \exp(-\lambda_d \|\Psi_d \mathbf{x}\|_1) \quad \text{i.i.d. Laplacian}$$

Outline

- 1 Introduction and Motivation for Composite Penalties
- 2 Co-L1 and its Interpretations
- 3 Co-IRW-L1 and its Interpretations**
- 4 Numerical Experiments

A Stepping Stone

The IRW version of real-valued Co-L1: tunes both inter-dictionary weights λ_d and intra-dictionary weights \mathbf{W}_d for given parameters ϵ_d .

1: input: $\{\Psi_d\}_{d=1}^D, \Phi, \mathbf{y}, \delta \geq 0, \epsilon_d > 0 \forall d,$

2: initialization: $\lambda_d^{(1)} = 1 \forall d, \mathbf{W}_d^{(1)} = \mathbf{I} \forall d$

3: for $t = 1, 2, 3, \dots$

4: $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x}} \sum_{d=1}^D \lambda_d^{(t)} \|\mathbf{W}_d^{(t)} \Psi_d \mathbf{x}\|_1 \text{ s.t. } \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$

5: $\lambda_d^{(t+1)} \leftarrow \left[\frac{1}{L_d} \sum_{l=1}^{L_d} \log \left(1 + \frac{|\psi_{d,l}^\top \mathbf{x}^{(t)}|}{\epsilon_d} \right) \right]^{-1} + 1, \quad d = 1, \dots, D$

6: $\mathbf{W}_d^{(t+1)} \leftarrow \text{diag} \left\{ \frac{1}{\epsilon_d + |\psi_{d,1}^\top \mathbf{x}^{(t)}|}, \dots, \frac{1}{\epsilon_d + |\psi_{d,L_d}^\top \mathbf{x}^{(t)}|} \right\}, \quad d = 1, \dots, D$

7: end

8: output: $\mathbf{x}^{(t)}$

Optimization Interpretations of real-Co-IRW-L1- ϵ

Real-Co-IRW-L1- ϵ is an MM approach to the **non-convex optimization**

$$\arg \min_{\mathbf{x}} \sum_{d=1}^D \sum_{l=1}^{L_d} \log \left[(\epsilon_d + |\boldsymbol{\psi}_{d,l}^T \mathbf{x}|) \sum_{i=1}^{L_d} \log \left(1 + \frac{|\boldsymbol{\psi}_{d,i}^T \mathbf{x}|}{\epsilon_d} \right) \right] \text{ s.t. } \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{x}\|_2 \leq \delta$$

and

As $\epsilon \rightarrow 0$, real-Co-IRW-L1- ϵ aims to solve the **$\ell_0 + \text{weighted } \ell_{0,0}$** problem

$$\arg \min_{\mathbf{x}} \left[\|\boldsymbol{\Psi} \mathbf{x}\|_0 + \sum_{d=1}^D L_d \mathbf{1}_{\|\boldsymbol{\Psi}_d \mathbf{x}\|_0 > 0} \right] \text{ s.t. } \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{x}\|_2 \leq \delta.$$

Note: L_d is the size of dictionary $\boldsymbol{\Psi}_d$, and $\mathbf{1}_{\square}$ is the indicator function.

Bayesian Interpretations of real-Co-IRW-L1- ϵ

Real-Co-IRW-L1 is an MM approach to **Bayesian MAP** estimation under an AWGN likelihood and the hierarchical prior

$$p(\mathbf{x}|\boldsymbol{\lambda}) = \prod_{d=1}^D \prod_{l=1}^{L_d} \frac{\lambda_d}{2\epsilon_d} \left(1 + \frac{|\boldsymbol{\psi}_{d,l}^\top \mathbf{x}|}{\epsilon_d} \right)^{-(\lambda_d+1)} \quad \text{i.i.d. generalized-Pareto}$$

$$p(\boldsymbol{\lambda}) = \prod_{d=1}^D p(\lambda_d), \quad p(\lambda_d) \propto \begin{cases} \frac{1}{\lambda_d} & \lambda_d > 0 \\ 0 & \text{else} \end{cases}, \quad \text{Jeffrey's non-informative}$$

and

Real-Co-IRW-L1 is a **variational EM** approach to estimating (deterministic) $\boldsymbol{\lambda}$ under an AWGN likelihood and the prior

$$p(\mathbf{x}; \boldsymbol{\lambda}) = \prod_{d=1}^D \prod_{l=1}^{L_d} \frac{\lambda_d - 1}{2\epsilon_d} \left(1 + \frac{|\boldsymbol{\psi}_{d,l}^\top \mathbf{x}|}{\epsilon_d} \right)^{-\lambda_d} \quad \text{i.i.d. generalized-Pareto}$$

The Co-IRW-L1 Algorithm

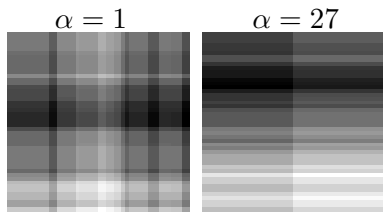
Finally, we self-tune ϵ_d and allow for real or complex quantities:

- 1: input: $\{\Psi_d\}_{d=1}^D, \Phi, \mathbf{y}, \delta \geq 0$,
- 2: if $\mathbf{x} \in \mathbb{R}^N$, use $\Lambda = (1, \infty)$ and the real version of $\log p(\mathbf{x}; \boldsymbol{\lambda}, \boldsymbol{\epsilon})$;
if $\mathbf{x} \in \mathbb{C}^N$, use $\Lambda = (2, \infty)$ and the complex version of $\log p(\mathbf{x}; \boldsymbol{\lambda}, \boldsymbol{\epsilon})$.
- 3: initialization: $\lambda_d^{(1)} = 1 \forall d, \mathbf{W}_d^{(1)} = \mathbf{I} \forall d$
- 4: for $t = 1, 2, 3, \dots$
- 5: $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x}} \sum_{d=1}^D \lambda_d^{(t)} \|\mathbf{W}_d^{(t)} \Psi_d \mathbf{x}\|_1$ s.t. $\|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$
- 6: $(\lambda_d^{(t+1)}, \epsilon_d^{(t+1)}) \leftarrow \arg \max_{\lambda_d \in \Lambda, \epsilon_d > 0} \log p(\mathbf{x}^{(t)}; \boldsymbol{\lambda}, \boldsymbol{\epsilon}), d = 1, \dots, D$
- 7: $\mathbf{W}_d^{(t+1)} \leftarrow \text{diag} \left\{ \frac{1}{\epsilon_d^{(t+1)} + |\boldsymbol{\psi}_{d,1}^\top \mathbf{x}^{(t)}|}, \dots, \frac{1}{\epsilon_d^{(t+1)} + |\boldsymbol{\psi}_{d,L_d}^\top \mathbf{x}^{(t)}|} \right\}, d = 1, \dots, D$
- 8: end
- 9: output: $\mathbf{x}^{(t)}$

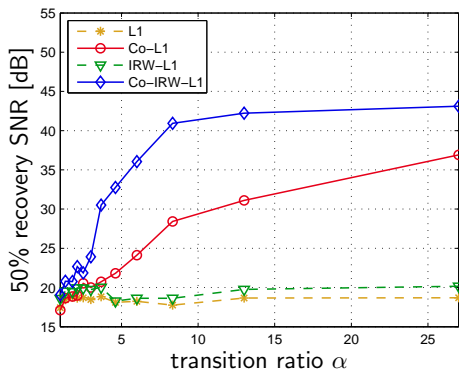
Outline

- 1 Introduction and Motivation for Composite Penalties
- 2 Co-L1 and its Interpretations
- 3 Co-IRW-L1 and its Interpretations
- 4 Numerical Experiments**

Experiment: Synthetic finite difference image



- 48×48 image with a total of 28 horiz & vert transitions.
- $\alpha \triangleq \frac{\# \text{ vertical transitions}}{\# \text{ horizontal transitions}}$
- “spread-spectrum” Φ
- sampling ratio $\frac{M}{N} = 0.3$
- AWGN @ 30 dB SNR
- $\Psi_1 =$ vertical finite difference,
 $\Psi_2 =$ horizon. finite difference

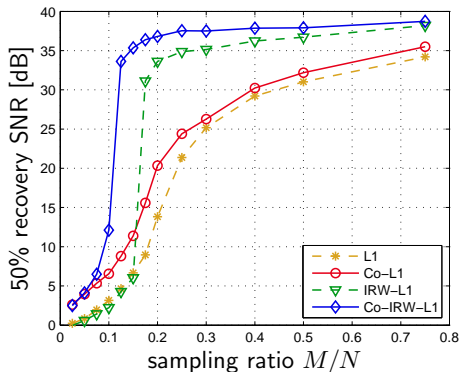


- ⇒ The composite algorithms significantly outperform the non-composite ones
- ⇒ Performance improves as sparsities become more disparate!

Experiment: Shepp-Logan Phantom



- 96×96 image
- “spread-spectrum” Φ
- AWGN @ 30 dB SNR
- $\Psi \in \mathbb{R}^{7N \times N} = 2D$ UWT-db1,
 $\Psi_d \in \mathbb{R}^{N \times N} \forall d$



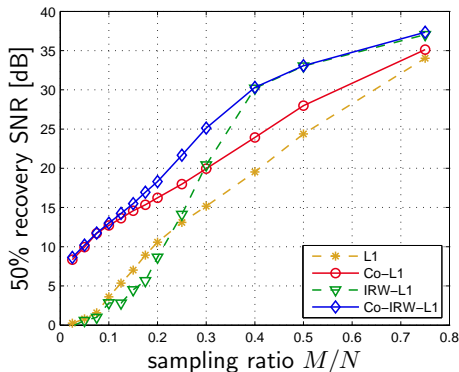
⇒ The composite algorithms significantly outperform the non-composite ones

⇒ Performance gap is larger for small M/N

Experiment: Cameraman



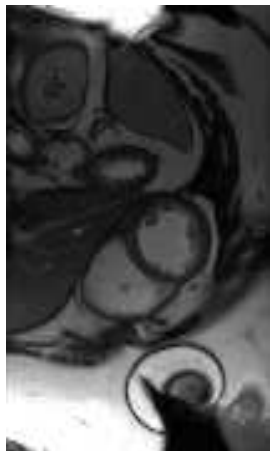
- 96×104 image
- “spread-spectrum” Φ
- AWGN @ 40 dB SNR
- $\Psi \in \mathbb{R}^{7N \times N} = 2D$ UWT-db1,
 $\Psi_d \in \mathbb{R}^{N \times N} \forall d$



⇒ The composite algorithms significantly outperform the non-composite ones

⇒ Performance gap is larger for small M/N

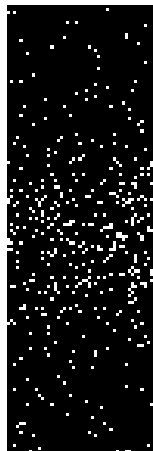
Experiment: 1D Dynamic MRI



x-y profile



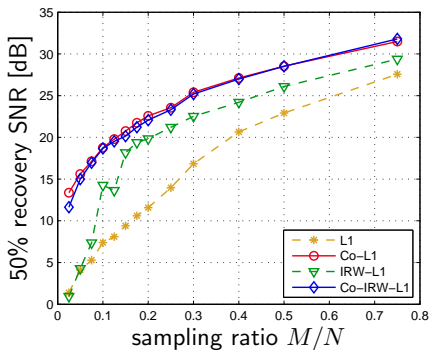
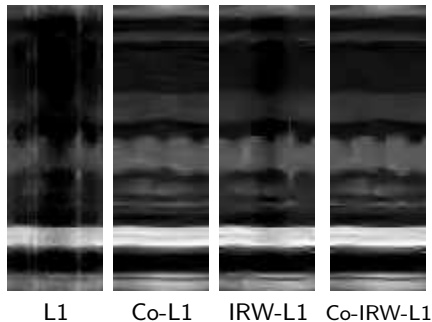
x-t profile



k-t sampling

- 144×48 spatiotemporal profile extracted from MRI cine
- Φ : variable density random Fourier
- AWGN @ 30 dB SNR
- $\Psi \in \mathbb{R}^{3N \times N}$: 2D [db1;db2;db3] DWT

Experiment: 1D Dynamic MRI (cont.)

sampling ratio $M/N = 0.3$ 

- The composite algs significantly outperform the non-composite ones
- Performance gap is larger for small M/N
- No advantage to Co-IRW-L1 over Co-L1 in this experiment

Runtimes for Previous Experiments

	Shepp-Logan	Cameraman	dMRI
L1	20.8s	23.1s	29.3s
Co-L1	32.7s	34.2s	86.4s
IRW-L1	45.9s	48.4s	54.1s
Co-IRW-L1	72.1s	96.4s	131s

The composite algs run 1.5–3× slower than the non-composite ones.

Conclusions

- We proposed a new “**composite-L1**” approach to L2-constrained signal reconstruction that **learns and exploits differences in sparsity across sub-dictionaries**.
- Relative to standard L1 methods, our composite L1 methods give **significant improvements in reconstruction SNR** at low sampling rates, at the cost of **1.5–3× slower runtimes**.
- Our algorithms can be interpreted as **MM approaches to non-convex optimization**, **approximate ℓ_0** methods, **Bayesian** methods, and **variational Bayesian** methods.

Conclusions

Thanks!

Iteratively Reweighted ℓ_1 (IRW-L1)

From [Candes, Wakin, Boyd, JFA'08] ...

-
- 1: input: $\Psi = [\psi_1, \dots, \psi_L]^T$, Φ , \mathbf{y} , $\delta \geq 0$, $\epsilon \geq 0$
 - 2: initialization: $\mathbf{W}^{(1)} = \mathbf{I}$
 - 3: for $t = 1, 2, 3, \dots$
 - 4: $\mathbf{x}^{(t)} \leftarrow \arg \min_{\mathbf{x}} \|\mathbf{W}^{(t)} \Psi \mathbf{x}\|_1$ s.t. $\|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta$
 - 5: $\mathbf{W}^{(t+1)} \leftarrow \text{diag} \left\{ \frac{1}{\epsilon + |\psi_1^T \mathbf{x}^{(t)}|}, \dots, \frac{1}{\epsilon + |\psi_L^T \mathbf{x}^{(t)}|} \right\}$
 - 6: end
 - 7: output: $\mathbf{x}^{(t)}$
-

- behaves more like ℓ_0 minimization than ℓ_1 minimization alone,
- leverages existing ℓ_1 solvers.

Majorize-Minimization (MM) Interpretation of IRW-L1

IRW-L1 is an MM approach to the log-sum optimization problem

$$\arg \min_{\mathbf{x}} \sum_{l=1}^L \log(\epsilon + |\boldsymbol{\psi}_l^\top \mathbf{x}|) \text{ s.t. } \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta.$$

How to see this? Reformulate as

$$\begin{aligned} & \arg \min_{\mathbf{x}, \mathbf{u}} \sum_l \log(\epsilon + u_l) \text{ s.t. } \begin{cases} \|\mathbf{y} - \Phi \mathbf{x}\|_2 \leq \delta \\ |\boldsymbol{\psi}_l^\top \mathbf{x}| \leq u_l \quad \forall l, \end{cases} \\ \Leftrightarrow & \arg \min_{\mathbf{v}} g(\mathbf{v}) \text{ s.t. } \mathbf{v} \in \mathcal{C} \end{aligned}$$

for $\mathbf{v} = \begin{bmatrix} \mathbf{u} \\ \mathbf{x} \end{bmatrix}$, convex \mathcal{C} , and concave g .

MM procedure: Iterate for $t = 1, 2, 3, \dots$

- 1 create surrogate $g(\mathbf{v}; \mathbf{v}^{(t)})$ that majorizes $g(\mathbf{v})$ at $\mathbf{v}^{(t)}$,
- 2 minimize the surrogate over $\mathbf{v} \in \mathcal{C}$, producing $\mathbf{v}^{(t+1)}$.

MM Interpretation of IRW-L1 (cont.)

Our concave $g(\mathbf{v})$ is **majorized by the tangent** at $\mathbf{v}^{(t)}$. So MM becomes

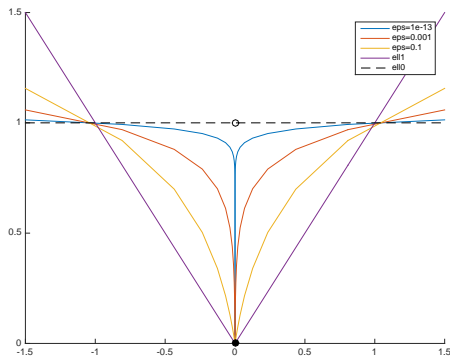
$$\begin{aligned}
 \mathbf{v}^{(t+1)} &= \arg \min_{\mathbf{v} \in \mathcal{C}} g(\mathbf{v}^{(t)}) + \nabla g(\mathbf{v}^{(t)})^\top [\mathbf{v} - \mathbf{v}^{(t)}] \\
 &= \arg \min_{\mathbf{v} \in \mathcal{C}} \nabla g(\mathbf{v}^{(t)})^\top \mathbf{v} \\
 \Leftrightarrow \mathbf{x}^{(t+1)} &= \arg \min_{\mathbf{x}} \underbrace{\sum_l \frac{1}{\epsilon + |\boldsymbol{\psi}_l^\top \mathbf{x}^{(t)}|} |\boldsymbol{\psi}_l^\top \mathbf{x}|}_{\|\mathbf{W}^{(t)} \boldsymbol{\Psi} \mathbf{x}\|_1} \text{ s.t. } \|\mathbf{y} - \boldsymbol{\Phi} \mathbf{x}\|_2 \leq \delta
 \end{aligned}$$

Implications of MM:

- IRW-L1 convergence is guaranteed
- but possibly to a suboptimal local minimum (since non-convex).

Approximate- ℓ_0 Interpretation of IRW-L1

$$\begin{aligned}
 & \sum_l \log(\epsilon + |u_l|) \\
 &= \sum_l \log(1 + |u_l|/\epsilon) + \text{const} \\
 &\propto \sum_l \frac{\log(1 + |u_l|/\epsilon)}{\log(1 + 1/\epsilon)} + \text{const} \rightarrow \\
 &= \sum_l \frac{\lim_{p \rightarrow 0} \frac{1}{p} \left[\left(1 + \frac{|u_l|}{\epsilon}\right)^p - 1 \right]}{\lim_{p \rightarrow 0} \frac{1}{p} \left[\left(1 + \frac{1}{\epsilon}\right)^p - 1 \right]} + \text{const} \\
 &= \lim_{p \rightarrow 0} \sum_l \frac{\left[\left(1 + \frac{|u_l|}{\epsilon}\right)^p - 1 \right]}{\left[\left(1 + \frac{1}{\epsilon}\right)^p - 1 \right]} + \text{const} \\
 &\approx \lim_{p \rightarrow 0} \sum_l |u_l|^p + \text{const} \quad (\text{for } \epsilon \ll 1) \\
 &= \|\mathbf{u}\|_0 + \text{const}
 \end{aligned}$$



\Rightarrow As $\epsilon \rightarrow 0$, the log-sum penalty becomes a scaled and shifted version of the ℓ_0 penalty.