

Last time:

- Properties of local minima of separable cost fn. - with reweighting based opt.

Today

- Convergence of reweighting based methods
- Some interpretations of the cost fn.
- Time permitting: Sparse Bayesian learning.

Recap: Separable cost fn.

$$\min_{x \in \mathbb{R}^N} \sum_{i=1}^N g(x_i) \quad \text{s.t. } y = Ax$$

$$\min_{x \in \mathbb{R}^N} \|y - Ax\|_1 + \lambda \sum_{i=1}^N g(x_i)$$

Start with $g(x) = (x^2 + \beta)^{\frac{1}{2}}$, $0 < \beta < 2$, $\rho > 0$.

Concave in x , found using 1st order Taylor expansion.

$$\Rightarrow \min_{x \in \mathbb{R}^N} \|Ax - y\|_1 + \lambda \sum_{i=1}^N \frac{(1+x_i^2)^{\frac{1}{2}}}{\beta^{\frac{1}{2}}} x_i^2$$

$$\Rightarrow x^{(k+1)} = \frac{(1+x^{(k)2})^{-\frac{1}{2}}}{W_k^{-1}} \text{diag}\{ (1+x^{(k)2})^{-\frac{1}{2}} \} y$$

When $\beta \rightarrow \infty$, we get $W_k = (1+x^{(k)2})^{-\frac{1}{2}}$ [Observed-Fn update]

When $\beta \rightarrow 0$, we get $W_k = 1$ [Focus].

$$\text{As } \beta \rightarrow 0, g(x) \rightarrow (x^2 + \beta)^{\frac{1}{2}}, \text{ so we are minimizing } \|Ax - y\|_1 + \lambda \sum_{i=1}^N |x_i|$$

\Rightarrow An intuitive explanation of why $g(x) = (x^2 + \beta)^{\frac{1}{2}}$ is a good idea.

Alternative viewpoint (connection to $g(x) = \log|x|$)

$$\|x\|_1 = \lim_{\beta \rightarrow 0} \sum_{i=1}^N (x_i^2 + \beta)^{\frac{1}{2}}$$

$$\lim_{\beta \rightarrow 0} \frac{1}{\beta} \sum_{i=1}^N (1+x_i^2) = \sum_{i=1}^N \log|x_i| \quad [\text{Show this!}]$$

Hence, \exists a 1-1 correspondence between (P) and

$$\min_{x \in \mathbb{R}^N} \sum_{i=1}^N \log|x_i| \quad \text{s.t. } y = Ax.$$

$$\text{Now, } \log|x_i| \leq \frac{\varepsilon}{2} \log|x_i^2 + \varepsilon| \quad \text{where } \varepsilon, \gamma_i \geq 0 \text{ and } \log \varepsilon \leq \gamma_i - 1 \text{ for } \varepsilon > 0 = \text{when } \gamma_i = x_i^2 + \varepsilon.$$

$$\text{Consider solving } \min_{x \in \mathbb{R}^N} \sum_{i=1}^N \frac{\varepsilon + x_i^2}{\gamma_i} + \frac{1}{2} \log \gamma_i \quad \text{s.t. } y = Ax.$$

Given γ, ε , quadratic in x w/ a linear constraint

$$x = \Gamma^{-\frac{1}{2}} (A \Gamma^{\frac{1}{2}})^{\dagger} y, \quad \text{where } \Gamma = \text{diag}(\gamma)$$

$$(A \Gamma^{\frac{1}{2}})^{\dagger} = \Gamma^{-\frac{1}{2}} A^{\dagger} (A \Gamma A^{\dagger})^{-1}$$

$$x = \Gamma^{-\frac{1}{2}} A^{\dagger} (A \Gamma A^{\dagger})^{-1} y$$

Given x, ε , the optimal γ is

$$\gamma_i = x_i^2 + \varepsilon = |x_i^2 + \varepsilon|$$

The exact schedule for updating ε is not crucial.

Weka: the update is the same as before!

$$x^{(k+1)} = W_k^{-1} A^{\dagger} (A \Gamma^k + A W_k A^{\dagger})^{-1} y$$

$$W_k^{-1} = \text{diag}\{ (1+x^{(k)2})^{-\frac{1}{2}} \} \Rightarrow \text{diag}\{ (1+x^{(k)2})^{-\frac{1}{2}} \}$$

$$(W_k^{-1})^2 = (1+x^{(k)2})^{-1} \Rightarrow W_k^{-1} = \text{diag}\{ (1+x^{(k)2})^{-\frac{1}{2}} \}$$

$\Rightarrow \sum_{i=1}^N \log|x_i|$ is a "good" cost fn. for solving (P) and

the above is an iterative approach for solving the non-separable cost minimization problem.

Zangwill's global convergence theorem:

States that, under certain conditions, which are satisfied here, the above updates are guaranteed to converge to a local min. or saddle pt. of

$$\sum_{i=1}^N \left\{ \frac{\varepsilon + x_i^2}{\gamma_i} + \log \gamma_i \right\}$$

from any starting point $x^{(0)}$.

Analysis of λ reweighting:

Focus on Candès et al.'s method: $W_k = (1+x^{(k)2})^{-\frac{1}{2}}$

$$\text{and solve } x^{(k+1)} = \arg \min_{x \in \mathbb{R}^N} \|y - Ax\|_1 + \lambda \sum_{i=1}^N W_k^{-1} |x_i|$$

$$g(x) = (1+x^2)^{\frac{1}{2}}, \quad \varepsilon > 0, \quad 0 < \rho < 1$$

Concave in $|x|$

$$g(x) \leq g'(x) x - (\text{terms that only dep. on } x_0)$$

At the $(k+1)$ th update, we solve

$$x^{(k+1)} = \arg \min_{x \in \mathbb{R}^N} \|y - Ax\|_1 + \lambda \|W_k^{-1} x\|_1$$

$$\text{where } W_k^{-1} = \text{diag}\{ g'(x^{(k)}) \} = \text{diag}\{ (1+x^{(k)2})^{-\frac{1}{2}} \}$$

So as $\beta \rightarrow 0$, we get the Candès et al. (2010) update.

\Rightarrow for small enough ε , we are solving

$$\min_{x \in \mathbb{R}^N} \|y - Ax\|_1 + \lambda \|x\|_1, \text{ as desired}$$

In the noisier case, the weight update is

$$\text{equivalent to solving } \min_{x \in \mathbb{R}^N} \sum_{i=1}^N \log(1+x_i^2 + \varepsilon) \quad \text{s.t. } y = Ax \quad \textcircled{2}$$

As before, by the Zangwill theory of global convergence,

the iterates converge to a local min or saddle pt. of $\textcircled{2}$.

[HW: Show that

$$\log(1+x^2 + \varepsilon) \leq \frac{x^2}{\gamma_i} + \log \left[\frac{(\varepsilon^2 + 2x^2) + \varepsilon}{2} \right] - \left[\frac{(\varepsilon^2 + 2x^2) + \varepsilon}{2} - \varepsilon \right]^2$$

for all $\varepsilon, \gamma_i \geq 0$, with equality iff $\gamma_i = x^2 + \varepsilon|x|$.

⇒ leads to a different iteratively reweighted z update

Sparse Bayesian Learning (SBL)

$$y = Ax + w, \quad w \sim \mathcal{N}(0, \sigma^2 I)$$

$$p(y | x; \sigma^2) = \frac{1}{\sqrt{(2\pi\sigma^2)^n}} \exp\left(-\frac{(y - Ax)^T (y - Ax)}{2\sigma^2}\right)$$

Directly finding the ML est. of x from $p(y | x; \sigma^2)$ amounts to $\min \|y - Ax\|_2^2$, which does not yield sparse solns. ⇒ Incorporate a sparsity-promoting prior on x .

FoCuss, BP etc can be cast into this framework, and can view them as finding the MAP est. of x .
For example, $p(x) \propto \exp\left(-\sum_{i=1}^n |x_i|^p\right)$, $p \in [0, 1]$.

Then, the MAP est. of x is

$$\begin{aligned} x_{\text{MAP}} &= \arg \max_{x \in \mathbb{R}^n} p(x | y) \propto p(y | x) p(x) \\ &= \arg \min_{x \in \mathbb{R}^n} \underbrace{-\log p(y | x)}_{\text{data fit}} - \underbrace{\log p(x)}_{\text{sparsity prior}} \\ &= \arg \min_{x \in \mathbb{R}^n} \|y - Ax\|_2^2 + \lambda \sum_{i=1}^n |x_i|^p \end{aligned}$$

The above is the cost fn. for BP when $p=1$, and for FoCuss when $p < 1$.

Our previous algos can be viewed as a MAP est. problem under an appropriately chosen sparsity promoting prior on x .