

Super-Linear Convergence of Dual Augmented Lagrangian Algorithm for Sparse Learning

Ryota Tomioka¹, Taiji Suzuki¹, and Masashi Sugiyama²

¹University of Tokyo

²Tokyo Institute of Technology

2009-12-12 @ NIPS workshop OPT09

Objective

Develop an optimization algorithm for the optimization problem:

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{f_\ell(\mathbf{A}\mathbf{w})}_{\text{loss}} + \underbrace{\phi_\lambda(\mathbf{w})}_{\text{regularizer}} .$$

For example, lasso:

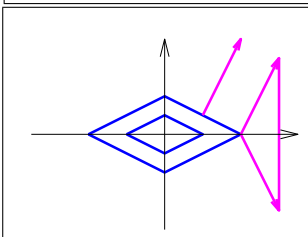
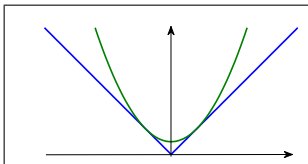
$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{A}\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 .$$

- $\mathbf{A} \in \mathbb{R}^{m \times n}$: design matrix (m : #observations, n : #unknowns) .
- f_ℓ is convex and twice differentiable.
- $\phi_\lambda(\mathbf{w})$ is convex but possibly non-differentiable. $\eta\phi_\lambda = \phi_{\eta\lambda}$.
- We are interested in algorithms for general f_ℓ and ϕ_λ (\leftrightarrow LARS).

Where does the difficulty come from?

Conventional view: the **non-differentiability** of $\phi_\lambda(\mathbf{w})$.

- Upper bound the regularizer from above with a differentiable function.
 - FOCUSS (Rao & Kreutz-Delgado, 99)
 - Majorization-Minimization (Figueiredo et al., 07)
 - Iteratively reweighted least squares (IRLS).
- Explicitly handle the non-differentiability.
 - Sub-gradient L-BFGS (Andrew & Gao, 07; Yu et al., 08)



Our view: the **coupling between variables** introduced by \mathbf{A} .

Where does the difficulty come from?

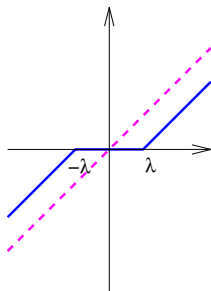
Our view: the coupling between variables introduced by \mathbf{A} .

In fact, when $\mathbf{A} = \mathbf{I}_n$

$$\min_{\mathbf{w} \in \mathbb{R}^n} \left(\frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_1 \right) = \sum_{j=1}^n \min_{w_j \in \mathbb{R}} \left(\frac{1}{2} (y_j - w_j)^2 + \lambda |w_j| \right).$$

$$\begin{aligned} \Rightarrow w_j^* &= \text{ST}_\lambda(y_j) \\ &= \begin{cases} y_j - \lambda & (\lambda \leq y_j), \\ 0 & (-\lambda \leq y_j \leq \lambda), \\ y_j + \lambda & (y_j \leq -\lambda). \end{cases} \end{aligned}$$

min is obtained analytically!



We focus on ϕ_λ for which the above min can be obtained analytically

Proximation wrt ϕ_λ can be computed analytically

Assumption

Proximation wrt ϕ_λ (soft-thresholding):

$$\text{ST}_\lambda(\mathbf{y}) = \underset{\mathbf{w} \in \mathbb{R}^n}{\text{argmin}} \left(\phi_\lambda(\mathbf{w}) + \frac{1}{2} \|\mathbf{y} - \mathbf{w}\|_2^2 \right)$$

can be computed analytically.

Outline

- 1 Introduction
 - Sparse regularized learning.
 - Why is it difficult? *not the non-differentiability*
- 2 Methods
 - Iterative shrinkage-thresholding (IST)
 - **Dual Augmented Lagrangian** (proposed)
- 3 Theoretical results: *super-linear convergence*
 - **Exact** inner minimization
 - **Approximate** inner minimization
- 4 Empirical results
 - Comparison against OWLQN, SpaRSA, and FISTA.
- 5 Summary

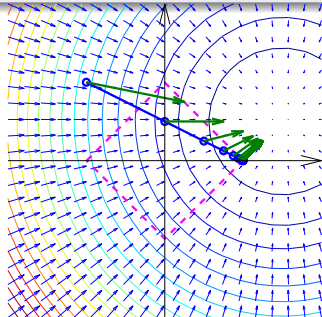
Iterative Shrinkage/Thresholding (IST)

Algorithm (Figueiredo&Nowak, 03; Daubechies et al., 04;...)

- 1 Choose an initial solution \mathbf{w}^0 .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t \lambda}}_{\text{shrink}} \left(\underbrace{\mathbf{w}^t - \eta_t \mathbf{A}^\top \nabla f_\ell(\mathbf{A} \mathbf{w}^t)}_{\text{gradient step}} \right).$$

- **Pro**: easy to implement.
- **Con**: bad for poorly conditioned \mathbf{A} .
- Also known as:
 - Forward-Backward Splitting [Combettes & Wajs, 05]
 - Thresholded Landweber Iteration [Daubechies et al., 04]



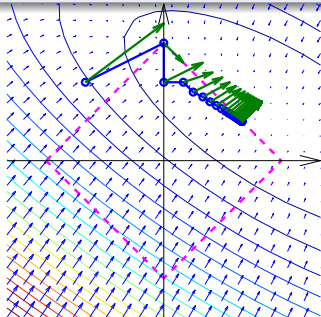
Iterative Shrinkage/Thresholding (IST)

Algorithm (Figueiredo&Nowak, 03; Daubechies et al., 04;...)

- 1 Choose an initial solution \mathbf{w}^0 .
- 2 Repeat until some stopping criterion is satisfied:

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t \lambda}}_{\text{shrink}} \left(\underbrace{\mathbf{w}^t - \eta_t \mathbf{A}^\top \nabla f_\ell(\mathbf{A} \mathbf{w}^t)}_{\text{gradient step}} \right).$$

- **Pro**: easy to implement.
- **Con**: bad for poorly conditioned \mathbf{A} .
- Also known as:
 - Forward-Backward Splitting [Combettes & Wajs, 05]
 - Thresholded Landweber Iteration [Daubechies et al., 04]



Dual Augmented Lagrangian (DAL) method

Primal problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad \underbrace{f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w})}_{f(\mathbf{w})}$$

Proximal minimization:

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

$$(\eta_0 \leq \eta_1 \leq \dots)$$

- Easy to analyze.
- $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$.
- Not practical! (as difficult as the original problem)

Dual problem

$$\underset{\boldsymbol{\alpha}, \mathbf{v}}{\text{maximize}} \quad -f_\ell^*(-\boldsymbol{\alpha}) - \phi_\lambda^*(\mathbf{v})$$

$$\text{s.t.} \quad \mathbf{v} = \mathbf{A}^\top \boldsymbol{\alpha}$$

\Leftrightarrow Augmented Lagrangian
(Tomioka & Sugiyama, 09):

$$\mathbf{w}^{t+1} = \text{ST}_{\lambda\eta_t}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}^t)$$

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha}}{\text{argmin}} \varphi_t(\boldsymbol{\alpha})$$

- Minimization of $\varphi_t(\boldsymbol{\alpha})$ is easy (smooth).
- Step-size η_t is increased.
- See Rockafellar 76 for the equivalence.

Dual Augmented Lagrangian (DAL) method

Primal problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad \underbrace{f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w})}_{f(\mathbf{w})}$$

Proximal minimization:

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

$$(\eta_0 \leq \eta_1 \leq \dots)$$

- Easy to analyze.
- $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$.
- Not practical! (as difficult as the original problem)

Dual problem

$$\underset{\alpha, \mathbf{v}}{\text{maximize}} \quad -f_\ell^*(-\alpha) - \phi_\lambda^*(\mathbf{v})$$

$$\text{s.t.} \quad \mathbf{v} = \mathbf{A}^\top \alpha$$

\Leftrightarrow Augmented Lagrangian
(Tomioka & Sugiyama, 09):

$$\mathbf{w}^{t+1} = \text{ST}_{\lambda\eta_t}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \alpha^t)$$

$$\alpha^t = \underset{\alpha}{\text{argmin}} \varphi_t(\alpha)$$

- Minimization of $\varphi_t(\alpha)$ is easy (smooth).
- Step-size η_t is increased.
- See Rockafellar 76 for the equivalence.

Dual Augmented Lagrangian (DAL) method

Primal problem

$$\underset{\mathbf{w}}{\text{minimize}} \quad \underbrace{f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w})}_{f(\mathbf{w})}$$

Proximal minimization:

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\text{argmin}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

$$(\eta_0 \leq \eta_1 \leq \dots)$$

- Easy to analyze.
- $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$.
- Not practical! (as difficult as the original problem)

Dual problem

$$\underset{\alpha, \mathbf{v}}{\text{maximize}} \quad -f_\ell^*(-\alpha) - \phi_\lambda^*(\mathbf{v})$$

$$\text{s.t.} \quad \mathbf{v} = \mathbf{A}^\top \alpha$$

\Leftrightarrow Augmented Lagrangian
(Tomioka & Sugiyama, 09):

$$\mathbf{w}^{t+1} = \text{ST}_{\lambda\eta_t}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \alpha^t)$$

$$\alpha^t = \underset{\alpha}{\text{argmin}} \varphi_t(\alpha)$$

- Minimization of $\varphi_t(\alpha)$ is easy (smooth).
- Step-size η_t is increased.
- See Rockafellar 76 for the equivalence.

Difference: How do we get rid of the couplings?

Proximation wrt f is hard:

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{variables are coupled}} + \underbrace{\phi_{\lambda}(\mathbf{w})}_{f(\mathbf{w})} + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right).$$

- IST: linearly approximates the loss term:

$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ tightest at the current point \mathbf{w}^t

- DAL (proposed): linearly lower-bounds the loss term:

$$f_{\ell}(\mathbf{A}\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m} \left(-f_{\ell}^*(-\alpha) - \mathbf{w}^{\top} \mathbf{A}^{\top} \alpha \right)$$

→ tightest at the next point \mathbf{w}^{t+1}

Difference: How do we get rid of the couplings?

Proximation wrt f is hard:

$$\mathbf{w}^{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} \left(\underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{variables are coupled}} + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right).$$

- IST: linearly approximates the loss term:

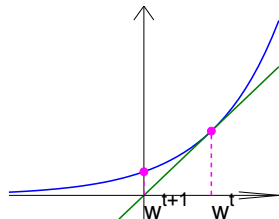
$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ tightest at the current point \mathbf{w}^t

- DAL (proposed): linearly lower-bounds the loss term:

$$f_{\ell}(\mathbf{A}\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m} \left(-f_{\ell}^*(-\alpha) - \mathbf{w}^{\top} \mathbf{A}^{\top} \alpha \right)$$

→ tightest at the next point \mathbf{w}^{t+1}



Difference: How do we get rid of the couplings?

Proximation wrt f is hard:

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{variables are coupled}} + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right).$$

- IST: linearly approximates the loss term:

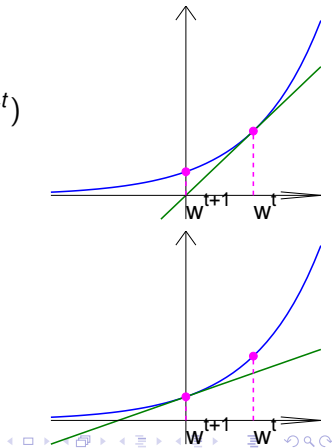
$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ tightest at the current point \mathbf{w}^t

- DAL (proposed): linearly lower-bounds the loss term:

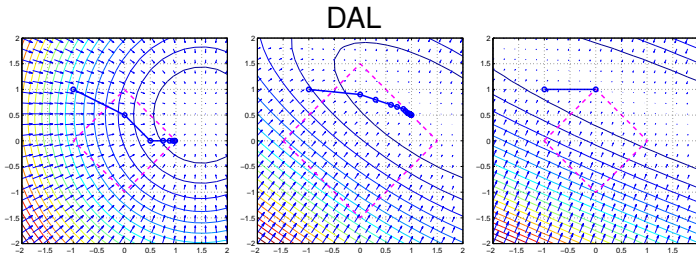
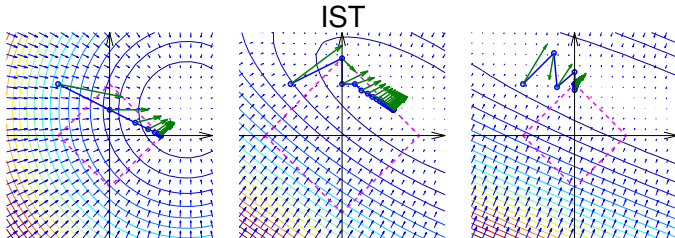
$$f_{\ell}(\mathbf{A}\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m} \left(-f_{\ell}^*(-\alpha) - \mathbf{w}^{\top} \mathbf{A}^{\top} \alpha \right)$$

→ tightest at the next point \mathbf{w}^{t+1}



Numerical examples

DAL is better when \mathbf{A} is poorly conditioned.



Theorem 1 (exact minimization)

Definition

- \mathbf{w}^t : sequence generated by the DAL algorithm with $\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\| = 0$ (exact minimization).
- \mathbf{w}^* : the unique minimizer of the objective f .

Assumption

There is a constant σ such that

$$f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*) \geq \sigma \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2 \quad (t = 0, 1, 2, \dots).$$

Theorem 1

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\| \leq \frac{1}{1 + \sigma\eta_t} \|\mathbf{w}^t - \mathbf{w}^*\|.$$

I.e., \mathbf{w}^t converges super-linearly to \mathbf{w}^* if η_t is increasing.

Theorem 2 (approximate minimization)

Definition

- \mathbf{w}^t : sequence generated by the DAL algorithm with

$$\|\nabla\varphi_t(\alpha^t)\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|\mathbf{w}^{t+1} - \mathbf{w}^t\| \quad \left(\begin{array}{l} 1/\gamma: \text{ Lipschitz con-} \\ \text{stant of } \nabla f_\ell. \end{array} \right)$$

Theorem 2

Under the same assumption as in Theorem 1,

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\| \leq \frac{1}{\sqrt{1 + 2\sigma\eta_t}} \|\mathbf{w}^t - \mathbf{w}^*\|.$$

I.e., \mathbf{w}^t converges super-linearly to \mathbf{w}^* if η_t is increasing.

Note

- Convergence is slower than the exact case ($\|\nabla\varphi_t(\alpha^t)\| = 0$).
- A faster rate can be obtained if we choose $\frac{\|\nabla\varphi_t(\alpha^t)\|}{\|\mathbf{w}^{t+1} - \mathbf{w}^t\|} \leq O(1/\eta_t)$.

Theorem 2 (approximate minimization)

Definition

- \mathbf{w}^t : sequence generated by the DAL algorithm with

$$\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|\mathbf{w}^{t+1} - \mathbf{w}^t\| \quad \left(\begin{array}{l} 1/\gamma: \text{ Lipschitz constant of } \nabla f_\ell. \end{array} \right)$$

Theorem 2

Under the same assumption as in Theorem 1,

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\| \leq \frac{1}{\sqrt{1 + 2\sigma\eta_t}} \|\mathbf{w}^t - \mathbf{w}^*\|.$$

I.e., \mathbf{w}^t converges super-linearly to \mathbf{w}^* if η_t is increasing.

Note

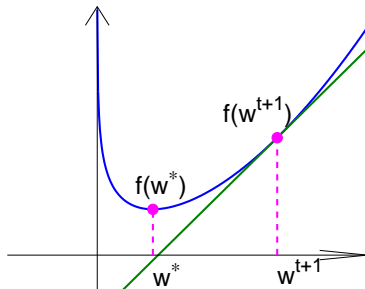
- Convergence is slower than the exact case ($\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\| = 0$).
- A faster rate can be obtained if we choose $\frac{\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\|}{\|\mathbf{w}^{t+1} - \mathbf{w}^t\|} \leq O(1/\eta_t)$.

Proof (in essence) of Theorem 1

Since $\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$,
 $(\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t \in \partial f(\mathbf{w}^{t+1})$ (is a subgradient of f). I.e.,

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle.$$

(inspired by Beck & Teboulle 09)

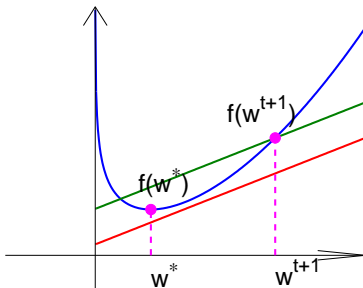


Proof (in essence) of Theorem 2

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle - \underbrace{\frac{1}{2\gamma} \|\nabla\varphi_t(\alpha^t)\|^2}_{\text{cost of approximate minimization}}.$$

cost of approximate
minimization

$1/\gamma$: Lipschitz constant of ∇f_ℓ .



Empirical results: ℓ_1 -logistic regression

#samples=1,024, #unknowns=16,384.

- FISTA

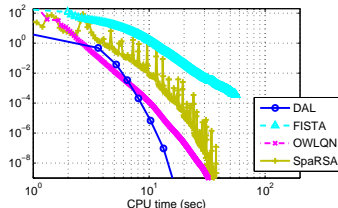
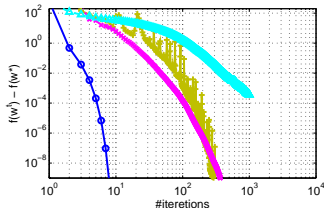
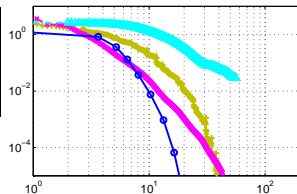
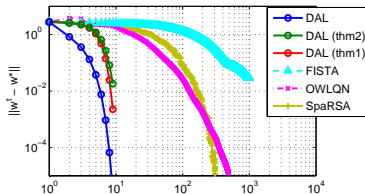
Two-step IST (Beck & Teboulle 09)

- OWLQN

Orthant-wise L-BFGS (Andrew & Gao 07)

- SpaRSA

Step-size improved IST (Wright et al. 09)



Summary

- Why is sparse learning difficult to optimize? – *couplings*
 - Non-differentiability is not bad.
 - Cost of inner minimization $O(m^2 n^+)$ (n^+ : number of active variables). **Sparsity makes inner minimization efficient.**
- How do we get rid of the couplings?
 - Use **linear parametric lower bound** instead of linear approximation.
- Super-linear convergence for exact/approximate inner minimization.
 - Improved a classic result in optimization by specializing the setting to **sparse learning**; i.e., proximation wrt ϕ_λ can be performed analytically.
- Empirical results are promising.
 - Faster than OWLQN, SpaRSA, and FISTA with the potential to be generalized further.

Summary

- Why is sparse learning difficult to optimize? – *couplings*
 - Non-differentiability is not bad.
 - Cost of inner minimization $O(m^2 n^+)$ (n^+ : number of active variables). **Sparsity makes inner minimization efficient.**
- How do we get rid of the couplings?
 - Use **linear parametric lower bound** instead of linear approximation.
- Super-linear convergence for exact/approximate inner minimization.
 - Improved a classic result in optimization by specializing the setting to **sparse learning**; i.e., proximation wrt ϕ_λ can be performed analytically.
- Empirical results are promising.
 - Faster than OWLQN, SpaRSA, and FISTA with the potential to be generalized further.

Summary

- Why is sparse learning difficult to optimize? – *couplings*
 - Non-differentiability is not bad.
 - Cost of inner minimization $O(m^2 n^+)$ (n^+ : number of active variables). **Sparsity makes inner minimization efficient.**
- How do we get rid of the couplings?
 - Use **linear parametric lower bound** instead of linear approximation.
- Super-linear convergence for exact/approximate inner minimization.
 - Improved a classic result in optimization by specializing the setting to **sparse learning**; i.e., proximation wrt ϕ_λ can be performed analytically.
- Empirical results are promising.
 - Faster than OWLQN, SpaRSA, and FISTA with the potential to be generalized further.

Summary

- Why is sparse learning difficult to optimize? – *couplings*
 - Non-differentiability is not bad.
 - Cost of inner minimization $O(m^2 n^+)$ (n^+ : number of active variables). **Sparsity makes inner minimization efficient.**
- How do we get rid of the couplings?
 - Use **linear parametric lower bound** instead of linear approximation.
- Super-linear convergence for exact/approximate inner minimization.
 - Improved a classic result in optimization by specializing the setting to **sparse learning**; i.e., proximation wrt ϕ_λ can be performed analytically.
- Empirical results are promising.
 - Faster than OWLQN, SpaRSA, and FISTA with the potential to be generalized further.

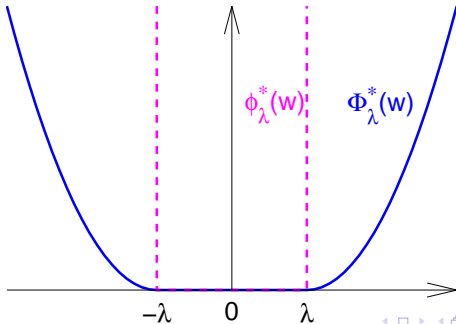
(1) Proximation wrt ϕ_λ is analytic (though non-smooth):

$$\mathbf{w}^{t+1} = \text{ST}_{\eta_t \lambda} \left(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}^t \right)$$

(2) Inner minimization is smooth:

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{argmin}} \left(\underbrace{f_\ell^*(-\boldsymbol{\alpha})}_{\text{independent of } \mathbf{A}.} + \frac{1}{2\eta_t} \underbrace{\|\text{ST}_{\eta_t \lambda}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha})\|_2^2}_{= \Phi_\lambda^*(\cdot)} \right)$$

(linear to the number of
active variables)



Comparison to other algorithms

- DAL (this talk)
 $\|\mathbf{w}^k - \mathbf{w}^*\| = O(\exp(-k))$
- SpaRSA (Step-size improved IST)
Convergence shown but no rate given. (Wright et al. 09)
- OWLQN (Orthant-wise L-BFGS)
Convergence shown but no rate given. (Andrew & Gao 07)
- IST (Iterative Soft-thresholding)
 $f(\mathbf{w}^k) - f(\mathbf{w}^*) = O(1/k)$ (Beck & Teboulle 09)
- FISTA (Two-step IST)
 $f(\mathbf{w}^k) - f(\mathbf{w}^*) = O(1/k^2)$ (Beck & Teboulle 09)

Comparison to Rockafellar 76

Assumption

The multifunction ∇f^* is (locally) Lipschitz continuous at the origin:

$$\|\nabla f^*(\beta) - \nabla f^*(0)\| \leq L\|\beta\| \quad (\|\beta\| \leq \tau)$$

\Rightarrow Implies our assumption with $\sigma = \frac{1}{2} \min(1/L, \tau/\|\mathbf{w}^0 - \mathbf{w}^*\|)$.

Convergence (exact minimization) – comparable to Thm 1

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\| \leq \frac{1}{\sqrt{1 + (\eta_t/L)^2}} \|\mathbf{w}^t - \mathbf{w}^*\|$$

Convergence (approximate minimization) – much worse than Thm 2

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\| \leq \frac{\mu_t + \epsilon_t}{1 - \epsilon_t} \|\mathbf{w}^t - \mathbf{w}^*\| \quad \left(\mu_t = \frac{1}{\sqrt{1 + (\eta_t/L)^2}} \right)$$

(assuming $\|\nabla \varphi_t\| \leq \epsilon_t \sqrt{\gamma/\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|$)

Outline of proof of Theorem 1

- ① Since $\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$,
 $(\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t$ is a subgradient of f at \mathbf{w}^{t+1} . I.e.,

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle.$$

- ② For any $\mu > 0$,

$$\|\mathbf{w}^* - \mathbf{w}^{t+1}\| \|\mathbf{w}^t - \mathbf{w}^*\| \leq \frac{\mu}{2} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|^2 + \frac{1}{2\mu} \|\mathbf{w}^t - \mathbf{w}^*\|^2.$$

- ③ Combining 1 & 2 and using $f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*) \geq \sigma \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2$,

$$\frac{1}{2} \|\mathbf{w}^t - \mathbf{w}^*\|^2 \geq \left((1 + \sigma\eta_t)\mu - \frac{\mu^2}{2} \right) \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2.$$

- ④ Maximize RHS wrt μ .

Outline of proof of Theorem 1

- ① Since $\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$,
 $(\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t$ is a subgradient of f at \mathbf{w}^{t+1} . I.e.,

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle.$$

- ② For any $\mu > 0$,

$$\|\mathbf{w}^* - \mathbf{w}^{t+1}\| \|\mathbf{w}^t - \mathbf{w}^*\| \leq \frac{\mu}{2} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|^2 + \frac{1}{2\mu} \|\mathbf{w}^t - \mathbf{w}^*\|^2.$$

- ③ Combining 1 & 2 and using $f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*) \geq \sigma \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2$,

$$\frac{1}{2} \|\mathbf{w}^t - \mathbf{w}^*\|^2 \geq \left((1 + \sigma\eta_t)\mu - \frac{\mu^2}{2} \right) \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2.$$

- ④ Maximize RHS wrt μ .

Outline of proof of Theorem 1

- ① Since $\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$,
 $(\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t$ is a subgradient of f at \mathbf{w}^{t+1} . I.e.,

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle.$$

- ② For any $\mu > 0$,

$$\|\mathbf{w}^* - \mathbf{w}^{t+1}\| \|\mathbf{w}^t - \mathbf{w}^*\| \leq \frac{\mu}{2} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|^2 + \frac{1}{2\mu} \|\mathbf{w}^t - \mathbf{w}^*\|^2.$$

- ③ Combining 1 & 2 and using $f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*) \geq \sigma \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2$,

$$\frac{1}{2} \|\mathbf{w}^t - \mathbf{w}^*\|^2 \geq \left((1 + \sigma\eta_t)\mu - \frac{\mu^2}{2} \right) \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2.$$

- ④ Maximize RHS wrt μ .

Outline of proof of Theorem 2

1 Let $\delta^t := \nabla \varphi_t(\alpha^t)$,

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle - \frac{1}{2\gamma} \|\delta^t\|^2.$$

2 By assumption

$$f(\mathbf{w}^{t+1}) - f(\mathbf{w}^*) \geq \sigma \|\mathbf{w}^{t+1} - \mathbf{w}^*\|^2,$$

$$\|\delta^t\|^2 \leq \frac{\gamma}{\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2.$$

3 Combining 1 & 2,

$$\frac{1}{2} \|\mathbf{w}^* - \mathbf{w}^t\|^2 \geq \left(\sigma \eta_t + \frac{1}{2} \right) \|\mathbf{w}^* - \mathbf{w}^{t+1}\|^2.$$

EEG problem – P300 visual speller dataset (subject A)

- Number of samples $m = 2550$.
- 6 class classification.
- $\mathbf{w} \in \mathbb{R}^{37 \times 64}$.
- Trace-norm regularization.

