

# Dual Augmented Lagrangian Algorithm 法による スパース正則化

富岡 亮太<sup>1</sup>, 鈴木 大慈<sup>1</sup>, 杉山 将<sup>2</sup>

<sup>1</sup> 東大 数理情報学専攻

<sup>2</sup> 東工大 計算工学専攻

tomioka@mist.i.u-tokyo.ac.jp

2010-4-13 @ 統計学輪講

# スパース回帰 (組み合わせ的)

- 入出力の組み

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m), \quad (\mathbf{x}_i \in \mathbb{R}^n).$$

- 仮定:

$$y_i = \langle \mathbf{w}, \mathbf{x}_i \rangle + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2).$$

- 経験誤差:

$$L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^m (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{w}\|^2.$$

- 問題:

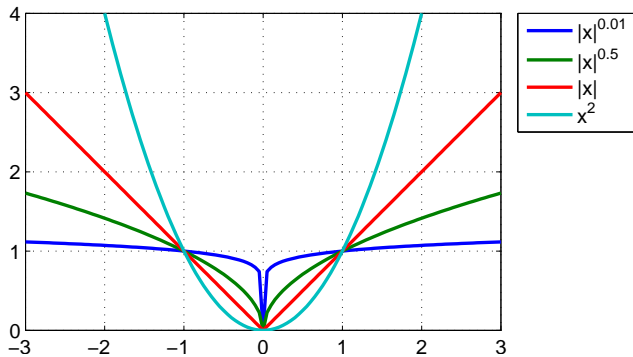
$$\text{minimize } L(\mathbf{w}), \quad \text{s.t. } \|\mathbf{w}\|_0 \leq C \quad (\text{NP 困難!!})$$

- 動機 1: 現実には多くの場合  $m < n$ .
- 動機 2: なるべく少ない数の変数で説明したい!  
( $\mathbf{w}$  の非ゼロ要素の数  $\|\mathbf{w}\|_0$  が少なければ少ないほどよい)

# スパース回帰 (連続)

- $p$ -ノルムの  $p$  乗 (のようなもの)

$$\|w\|_p^p = \sum_{j=1}^n |w_j|^p : \begin{cases} p \geq 1 \text{ ならば凸} \\ p < 1 \text{ ならば非凸} \end{cases}$$

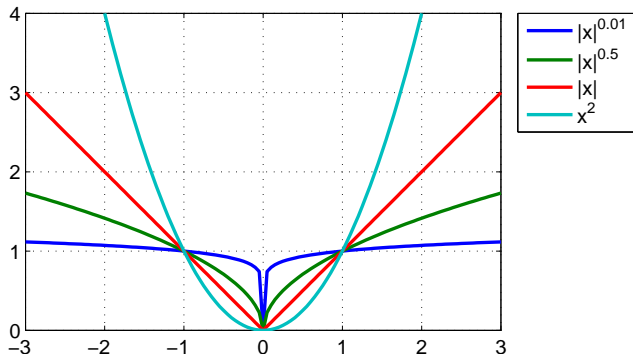


$\|w\|_1$  の正則化は凸の中ではもっとも  $\|w\|_0$  の正則化に近い!

# スパース回帰 (連続)

- $p$ -ノルムの  $p$  乗 (のようなもの)

$$\|w\|_p^p = \sum_{j=1}^n |w_j|^p : \begin{cases} p \geq 1 \text{ ならば凸} \\ p < 1 \text{ ならば非凸} \end{cases}$$



$\|w\|_1$  の正則化は凸の中ではもっとも  $\|w\|_0$  の正則化に近い!

## Lasso 回帰 (連続かつ凸)

## ● 問題 1:

$$\text{minimize } \|\mathbf{w}\|_1, \quad \text{s.t. } L(\mathbf{w}) \leq C.$$

## ● 問題 2:

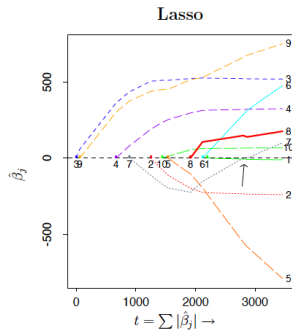
$$\text{minimize } L(\mathbf{w}), \quad \text{s.t. } \|\mathbf{w}\|_1 \leq C'.$$

## ● 問題 3:

$$\text{minimize } L(\mathbf{w}) + \lambda \|\mathbf{w}\|_1.$$

注意:

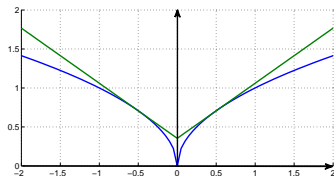
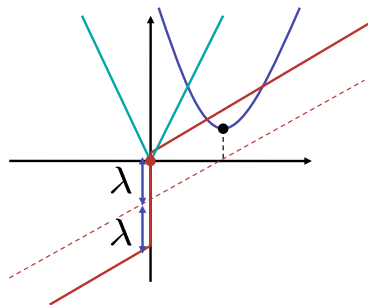
- 上の3つの問題はいずれも等価。
- 正則化項やロス項に単調な非線形変換をしても等価。
- この発表では問題3を扱う。



[From Efron et al. (2003)]

# なぜ $l_1$ -正則化か？

- 凸の中で最も  $\|\cdot\|_0$  に近い .
- 原点で微分不可能 (有限の  $\lambda$  でゼロに打ち切ることができる.)
- 凸でない正則化  
→ 繰り返し (重み付き)  $l_1$ -正則化問題を解けばよい .
- ベイズ周辺化尤度最大化  
→ (特殊な場合に) 繰り返し (重み付き)  $l_1$ -正則化問題を解けばよい .  
(Wipf&Nagarajan, 08)



# 問題設定

以下の最適化問題を効率良く解くためのアルゴリズムが求められている。

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

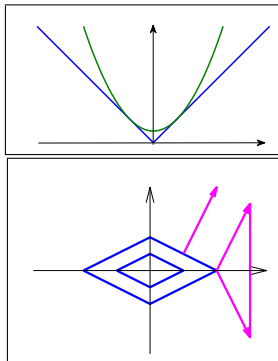
ただし，

- $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m$ : サンプル数,  $n$ : 未知変数の数) .
- $f_\ell$  は凸で 2 回微分可能 .
- $\phi_\lambda(\mathbf{w})$  は例えば,  $\phi_\lambda(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$  など, 凸だが, 微分不可能であってもよい . また,  $\eta\phi_\lambda = \phi_{\eta\lambda}$  を仮定 .
- 特定の  $f_\ell$  に依存したアルゴリズム (LARS など) ではもの足りない .
- No Free Lunch – 観測の数が変数の数より少ない場合 ( $m \ll n$ ) や  $\mathbf{A}$  のコンディションが悪い場合が応用上重要 .

# どこが難しいか？

今までの見方:  $\phi_\lambda(\mathbf{w})$  の微分不可能性が原因 .

- 正則化項を微分可能な関数で上から押さえる .
  - FOCUSS  
(Rao & Kreutz-Delgado, 99)
  - Majorization-Minimization  
(Figueiredo et al., 07)
- 微分不可能性を陽に考慮する .
  - Sub-gradient L-BFGS (Andrew & Gao, 07; Yu et al., 08)



我々の見方: **A** が変数の間にからみを導入するのが原因 .



# どこが難しいか？

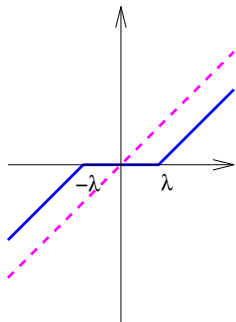
我々の見方:  $\mathbf{A}$  が変数の間からみを導入するのが原因.

$\mathbf{A} = \mathbf{I}_n$  (単位行列の場合)

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

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

解析的に解ける！



本発表では  $\phi_\lambda$  として、上の最小化が解析的に求められるもののみを扱う。

# $\phi_\lambda$ に関する proximation は解析的に計算できる

## 仮定

$\phi_\lambda$  に関する proximation (soft-thresholding)

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

は解析的に計算できる。

- グループラッソー

$$\phi_\lambda(\mathbf{w}) = \sum_{\mathbf{g} \in \mathcal{G}} \|\mathbf{w}_{\mathbf{g}}\|_2 : \text{グループ単位のノルムの和.}$$

- トレースノルム

$$\phi_\lambda(\mathbf{W}) = \sum_{j=1}^r \sigma_j(\mathbf{W}) : \text{特異値の和 ( } \mathbf{W} \text{ は行列).}$$

# Proximation の解釈 (脱線)

- Proximation は 集合に対する射影の一般化 .
- 集合  $C$  の定義関数  $\delta_C(\mathbf{x})$  を

$$\delta_C(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in C \text{ のとき,} \\ +\infty, & \text{その他,} \end{cases}$$

と定義する .

$\delta_C$  に関する proximation  $\equiv$  集合  $C$  への 2 乗距離の意味での射影 .

- Proximation による分解 :

$$\text{prox}_f(\mathbf{z}) + \text{prox}_{f^*}(\mathbf{z}) = \mathbf{z}.$$

ただし ,  $f$  と  $f^*$  は凸共役な関数の組み .

- 参照 : Moreau 1965, Rockafellar 1970.

# 発表の流れ

- ① イントロ
  - スパース正則化
  - どこが難しいか： **微分不可能性**  $\Rightarrow$  変数の間のからみ
- ② 最適化アルゴリズム
  - Iterative shrinkage-thresholding (IST)
  - **Dual Augmented Lagrangian** (提案手法)
- ③ 収束レート： **超 1 次収束**
  - **厳密な** 内部最小化の場合
  - **近似的な** 内部最小化の場合
- ④ 実験
  - OWLQN, SpaRSA, and FISTA との比較 .
- ⑤ まとめ

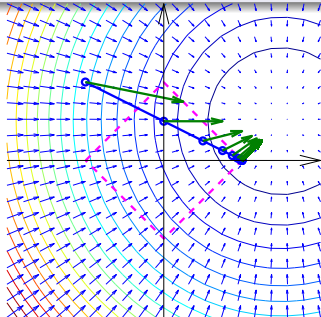
# Iterative Shrinkage/Thresholding (IST)

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

- ① 適当に初期解  $w^0$  を決める.
- ② 停止条件が満たされるまで反復 :

$$w^{t+1} \leftarrow \underbrace{ST_{\eta_t \lambda}}_{\text{縮小}} \left( \underbrace{w^t - \eta_t A^\top \nabla f_\ell(Aw^t)}_{\text{勾配ステップ}} \right).$$

- 利点: 実装が簡単 .
- 欠点: デザイン行列  $A$  の条件数が悪いと遅い .



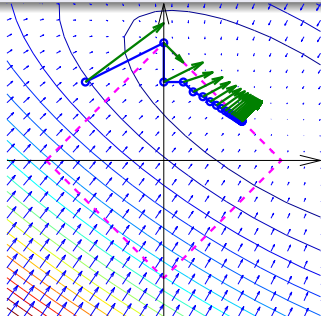
# Iterative Shrinkage/Thresholding (IST)

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

- 1 適当に初期解  $w^0$  を決める.
- 2 停止条件が満たされるまで反復 :

$$w^{t+1} \leftarrow \underbrace{ST_{\eta_t \lambda}}_{\text{縮小}} \left( \underbrace{w^t - \eta_t A^T \nabla f_\ell(Aw^t)}_{\text{勾配ステップ}} \right).$$

- 利点: 実装が簡単 .
- 欠点: デザイン行列  $A$  の条件数が悪いと遅い .



## Dual Augmented Lagrangian (DAL) 法 (提案手法)

## 主問題

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

## 双対問題

$$\begin{aligned} &\underset{\alpha, \mathbf{v}}{\text{maximize}} && -f_\ell^*(-\alpha) - \phi_\lambda^*(\mathbf{v}) \\ &\text{s.t.} && \mathbf{v} = \mathbf{A}^\top \alpha \end{aligned}$$

- $f_\ell^*$ ,  $\phi_\lambda^*$  は  $f_\ell$ ,  $\phi_\lambda$  の凸共役 :

$$f_\ell^*(\alpha) = \sup_{\mathbf{z} \in \mathbb{R}^m} (\langle \alpha, \mathbf{z} \rangle - f_\ell(\mathbf{z}))$$

$$\phi_\lambda^*(\mathbf{v}) = \sup_{\mathbf{w} \in \mathbb{R}^n} (\langle \mathbf{v}, \mathbf{w} \rangle - \phi_\lambda(\mathbf{w}))$$

- 主問題の最小値 = 双対問題の最大値 (強双対性)

## Dual Augmented Lagrangian (DAL) 法 (提案手法)

## 主問題

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

Proximal minimization:

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

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

- 解析がしやすい. 例えば  
 $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$ .
- 実用的でない (もとの問題と同程度に難しい!)

## 双対問題

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

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

⇔ Augmented Lagrangian  
(Tomioka & Sugiyama, 09):

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

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

- $\varphi_t(\alpha)$  の最小化は簡単 (なめらか).
- ステップサイズ  $\eta_t$  は増加.
- 同値性については Rockafellar 76 を参照.



## Dual Augmented Lagrangian (DAL) 法 (提案手法)

## 主問題

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

Proximal minimization:

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

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

- 解析がしやすい. 例えば  
 $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$ .
- 実用的でない (もとの問題と同程度に難しい!)

## 双対問題

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

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

⇔ Augmented Lagrangian  
(Tomioka & Sugiyama, 09):

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

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

- $\varphi_t(\alpha)$  の最小化は簡単 (なめらか).
- ステップサイズ  $\eta_t$  は増加.
- 同値性については Rockafellar 76 を参照.

## Dual Augmented Lagrangian (DAL) 法 (提案手法)

## 主問題

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

Proximal minimization:

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

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

- 解析がしやすい. 例えば  
 $f(\mathbf{w}^{t+1}) + \frac{1}{2\eta_t} \|\mathbf{w}^{t+1} - \mathbf{w}^t\|^2 \leq f(\mathbf{w}^t)$ .
- 実用的でない (もとの問題と同程度に難しい!)

## 双対問題

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

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

⇔ Augmented Lagrangian  
(Tomioka & Sugiyama, 09):

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

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

- $\varphi_t(\alpha)$  の最小化は簡単 (なめらか).
- ステップサイズ  $\eta_t$  は増加.
- 同値性については Rockafellar 76 を参照.

Dual Augmented Lagrangian 法 ( $l_1$ -正則化)

- ① 適当に初期解  $\mathbf{w}^1$  を決める .
- ② 停止条件が満たされるまで反復 :

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

ただし ,

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{argmin}} \left( \underbrace{f_\ell^*(-\boldsymbol{\alpha})}_{\text{損失関数 } f_\ell \text{ の凸共役}} + \frac{1}{2\eta_t} \|\text{ST}_{\eta_t \lambda}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha})\|_2^2 \right)$$

損失関数  $f_\ell$  の凸共役

Dual Augmented Lagrangian 法 ( $\ell_1$ -正則化の場合)

(1) この計算は解析的にできると仮定.

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

(2) この計算はスパースであればあるほど効率的. ( “アクティブな”  
未知変数の数  
に線形 )

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha} \in \mathbb{R}^m}{\text{argmin}} \left( \underbrace{f_\ell^*(-\boldsymbol{\alpha})}_{\substack{\mathbf{A} \text{ のスケーリン} \\ \text{グの影響を受け} \\ \text{ない}}} + \underbrace{\frac{1}{2\eta_t} \|\text{ST}_{\eta_t \lambda}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha})\|_2^2}_{\substack{\frac{\partial}{\partial \boldsymbol{\alpha}} : \mathbf{A} \text{ST}_{\eta_t \lambda}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}) \\ \frac{\partial}{\partial \alpha^2} : \eta_t \mathbf{A}_+ + \mathbf{A}_+^\top \\ (\mathbf{A}_+ \text{ は } \mathbf{A} \text{ の “アクティブな” 列から} \\ \text{なる部分行列; } \ell_1\text{-正則化の場合)}}} \right)$$

(3)  $\mathbf{A}$  のスケーリングの悪さの影響を受けにくい.

# IST と DAL の違い：いかに変数の間の絡みを取り除くか

目的関数  $f$  に関する Proximation は難しい：

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left( \underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{変数が絡みあっている}} + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- IST (既存手法)：線形にロス項を 近似：

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

→ 現在の点  $\mathbf{w}^t$  で最もタイト

- DAL (提案法)：線形なロス項の 下限

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

→ 次の点  $\mathbf{w}^{t+1}$  で最もタイト

# IST と DAL の違い：いかに変数の間の絡みを取り除くか

目的関数  $f$  に関する Proximation は難しい：

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left( \underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{変数が絡みあっている}} + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- IST (既存手法)：線形にロス項を 近似：

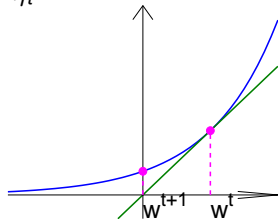
$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ 現在の点  $\mathbf{w}^t$  で最もタイト

- DAL (提案法)：線形なロス項の 下限

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

→ 次の点  $\mathbf{w}^{t+1}$  で最もタイト



# IST と DAL の違い：いかに変数の間の絡みを取り除くか

目的関数  $f$  に関する Proximation は難しい：

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left( \underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{変数が絡みあっている}} + \phi_{\lambda}(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- IST (既存手法)：線形にロス項を 近似：

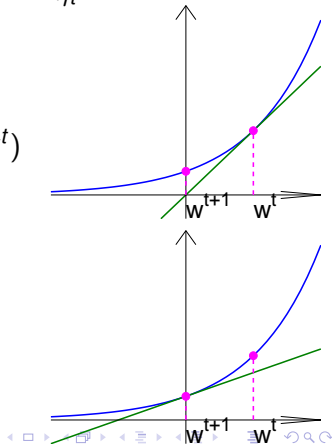
$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ 現在の点  $\mathbf{w}^t$  で最もタイト

- DAL (提案法)：線形なロス項の 下限

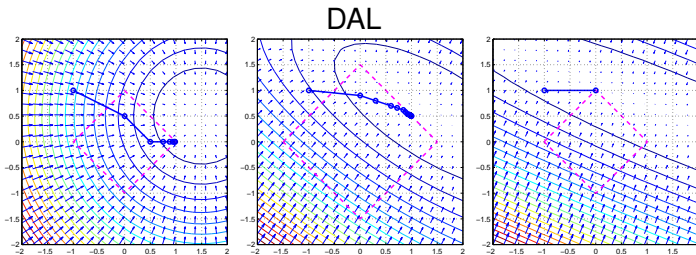
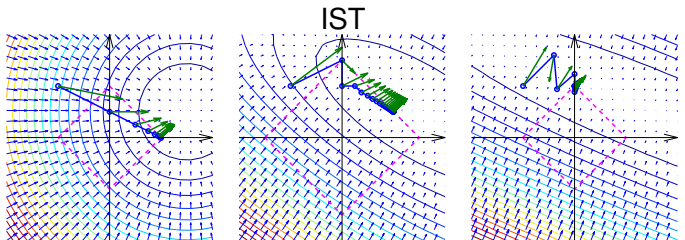
$$f_{\ell}(\mathbf{A}\mathbf{w}) = \max_{\alpha \in \mathbb{R}^m} \left( -f_{\ell}^*(-\alpha) - \mathbf{w}^{\top} \mathbf{A}^{\top} \alpha \right)$$

→ 次の点  $\mathbf{w}^{t+1}$  で最もタイト



# 数値例

デザイン行列  $\mathbf{A}$  のコンディションが悪くなるほど，DAL の方が有利．





# 定理 1 ( 厳密な最小化 )

## 定義

- $\mathbf{w}^t$  : 厳密な DAL 法 (  $\|\nabla\varphi_t(\alpha^t)\| = 0$  ) で得られる点列 .
- $\mathbf{w}^*$  : 目的関数  $f$  を最小化する点 .

## 仮定

正の定数  $\sigma$  が存在して

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

## 定理 1

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

$\Rightarrow \eta_t$  が増加するなら ,  $\mathbf{w}^t$  は  $\mathbf{w}^*$  に **超 1 次収束** する .

## 定理 2 (近似的最小化)

### 定義

- $w^t$ : 以下の停止基準による近似的な DAL 法で得られる点列 .

$$\|\nabla\varphi_t(\alpha^t)\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|w^{t+1} - w^t\| \quad \left( \begin{array}{l} 1/\gamma: \text{損失関数の微分} \\ \nabla f_\ell \text{ のリプシッツ定数.} \end{array} \right)$$

### 定理 2

定理 1 と同じ仮定のもとで

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

⇒  $\eta_t$  が増加するなら,  $w^t$  は  $w^*$  に超 1 次収束する .

- 収束レートは厳密な場合 ( $\|\nabla\varphi_t(\alpha^t)\| = 0$ ) より少し悪い .
- 同程度の収束レートは内部最小化をもう少し厳しくすることで達成可能  $\frac{\|\nabla\varphi_t(\alpha^t)\|}{\|w^{t+1} - w^t\|} \leq O(1/\eta_t)$ .

## 定理 2 (近似的最小化)

### 定義

- $\mathbf{w}^t$ : 以下の停止基準による近似的な DAL 法で得られる点列 .

$$\|\nabla\varphi_t(\alpha^t)\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|\mathbf{w}^{t+1} - \mathbf{w}^t\| \quad \left( \begin{array}{l} 1/\gamma: \text{損失関数の微分} \\ \nabla f_\ell \text{ のリプシッツ定数.} \end{array} \right)$$

### 定理 2

定理 1 と同じ仮定のもとで

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

$\Rightarrow \eta_t$  が増加するなら,  $\mathbf{w}^t$  は  $\mathbf{w}^*$  に超 1 次収束する .

- 収束レートは厳密な場合 ( $\|\nabla\varphi_t(\alpha^t)\| = 0$ ) より少し悪い .
- 同程度の収束レートは内部最小化をもう少し厳しくすることで達成可能  $\frac{\|\nabla\varphi_t(\alpha^t)\|}{\|\mathbf{w}^{t+1} - \mathbf{w}^t\|} \leq O(1/\eta_t)$ .

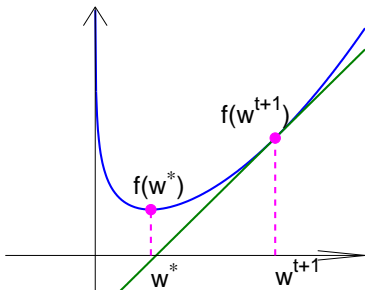
## 定理 1 の証明 (エッセンス)

$\mathbf{w}^{t+1}$  は,  $f(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2$  を最小化するので,

$$(\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t \in \partial f(\mathbf{w}^{t+1}) \quad (\text{劣微分に入る})$$

従って (Beck & Teboulle 09),

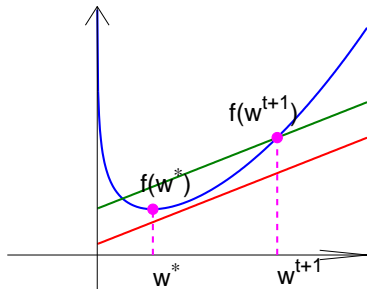
$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1})/\eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle.$$



## 定理 2 の証明 ( エッセンス )

$$f(\mathbf{w}^*) - f(\mathbf{w}^{t+1}) \geq \left\langle (\mathbf{w}^t - \mathbf{w}^{t+1}) / \eta_t, \mathbf{w}^* - \mathbf{w}^{t+1} \right\rangle - \underbrace{\frac{1}{2\gamma} \|\nabla \varphi_t(\alpha^t)\|^2}_{\text{近似最小化のコスト}} .$$

$1/\gamma$ : 損失関数の微分  $\nabla f_\ell$  のリプシッツ定数 .



## 他の手法との比較

	手法	説明	収束オーダー
1次	IST		$O(1/k)$
	FISTA	ISTの収束性を改善したもの	$O(1/k^2)$
	SpaRSA	曲率の情報を利用してISTのステップサイズを改善したもの	?
	OWLQN	劣微分を利用した擬似ニュートン法	?
高次	内点法	Koh, Kim, & Boyd 2007	$O(e^{-k})$
	DAL	提案法	$o(e^{-k})$

# 数値実験: $\ell_1$ -正則化付きロジスティック回帰

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

- FISTA

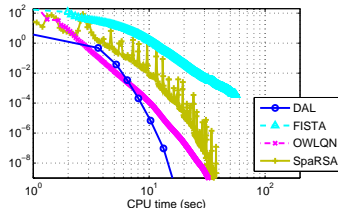
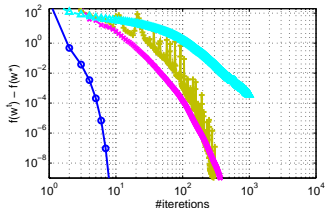
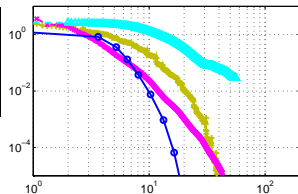
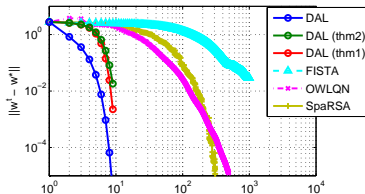
2 段階 IST (Beck & Teboulle 09)

- OWLQN

劣微分準ニュートン  
(Andrew & Gao 07)

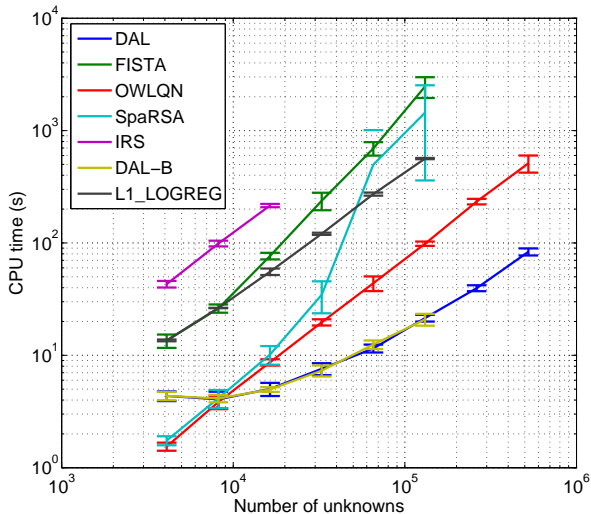
- SpaRSA

ステップサイズ改良  
IST (Wright et al. 09)



# 未知変数の数に対する計算時間の増加

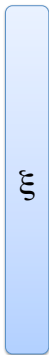
$m = 1,024$ .  $n = 4,096 - 524,288$





# 脳波解析における接続関係の推定

独立・  
非ガウス信号

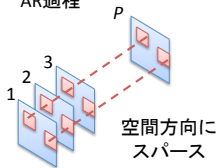


$$z(t) = \sum_{p=1}^P H^{(p)} z(t-p)$$

+  $\xi(t)$



時空間  
スパース  
AR過程



源信号

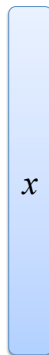


$$x(t) = M z(t)$$



空間方向  
即時的  
混合

観測信号



S. Haufe, R. Tomioka, G. Nolte, K-R Müller and M. Kawanabe, IEEE TBME 2010. accepted.

## EM アルゴリズム

- E-ステップ：隠れ信号  $\mathbf{z}(t)$  の推定 ( $\mathbf{B}$  に関する最尤法)
  - 混合行列の逆行列  $\mathbf{B} := \mathbf{M}^{-1}$  とする .
  - $\xi(t) = \mathbf{B}\mathbf{x}(t) - \sum_{p=1}^P \mathbf{H}^{(p)} \mathbf{B}\mathbf{x}(t-p)$  は *sech* 分布に従うと仮定 .

$$p(\xi) \propto \prod_{d=1}^D \frac{1}{e^{-\xi_d} + e^{\xi_d}}$$



- M-ステップ：AR 係数の推定 ( $\mathbf{H}^{(p)}$  に関するスパース推定)
  - 尤度： *sech* 分布
  - 正則化：(空間方向に) スパースな接続関係を見つけたい  
⇒ 時間方向にグループしたグループラッソー正則化 .

# 結果 (人工データ)

SCSA\_EM (提案手法) が混合行列  $M$  の推定および結合係数  $H^{(P)}$  の推定の両方の意味で優れている。

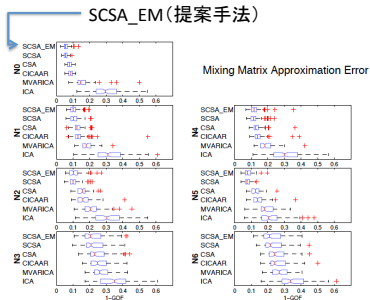


Fig. 3. Estimation errors of the mixing matrix according to the goodness-of-fit (GOF) criterion. Results are shown for the proposed (Sparsely-) Connected Sources Analysis variants (SCSA\_EM, SCBA, CSA) and three alternative approaches (CICAAR, MVARICA, ICA). Different subfigures depict the methods' performance in the noiseless case (N0), as well as in the presence of different types of noise (N1-N6, see TABLE I).

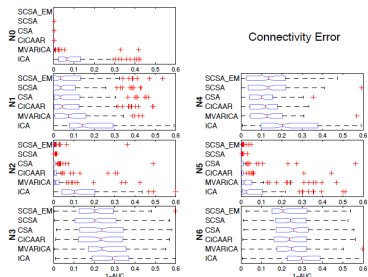


Fig. 5. Estimation errors regarding the source connectivity structure as measured by fitting an MVAR model subsequently to the demixed sources and testing the obtained coefficients for significant interaction. The performance measure reported is the area under the curve (AUC) score obtained by varying the significance level.

# Summary

- なぜスパース正則化は難しい? – 変数の間の絡み
  - 微分不可能性は悪くない.
  - 内部最小化はスパースであればあるほど速い.  
微分不可能性は効率的な内部最小化のために使える
- いかに変数の間の絡みを取り除くか
  - 線形近似の代わりに パラメトリックな下限を使う.
- 厳密 / 近似的内部最小化のもとでの超 1 次収束.
  - スパース正則化の状況に特殊化することで, 現実的・使いやすい条件のもとで, 既存の収束レートを改善.
- 実験結果も理論をサポート
  - 既存手法 OWLQN, SpaRSA, and FISTA より速い. かつより複雑な正則化に対応可能.

# Summary

- なぜスパース正則化は難しい? – 変数の間の絡み
  - 微分不可能性は悪くない.
  - 内部最小化はスパースであればあるほど速い.  
微分不可能性は効率的な内部最小化のために使える
- いかに変数の間の絡みを取り除くか
  - 線形近似の代わりに パラメトリックな下限を使う.
- 厳密 / 近似的内部最小化のもとでの超 1 次収束.
  - スパース正則化の状況に特殊化することで, 現実的・使いやすい条件のもとで, 既存の収束レートを改善.
- 実験結果も理論をサポート
  - 既存手法 OWLQN, SpaRSA, and FISTA より速い. かつより複雑な正則化に対応可能.

# Summary

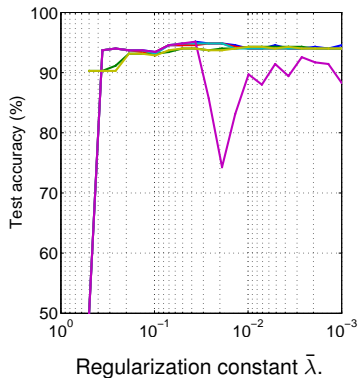
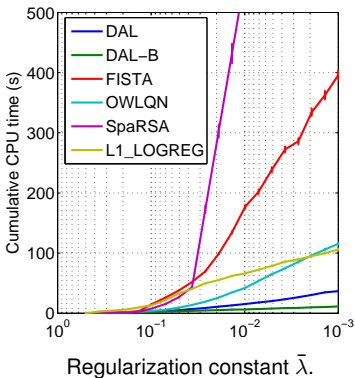
- なぜスパース正則化は難しい? – 変数の間の絡み
  - 微分不可能性は悪くない.
  - 内部最小化はスパースであればあるほど速い.  
微分不可能性は効率的な内部最小化のために使える
- いかに変数の間の絡みを取り除くか
  - 線形近似の代わりに パラメトリックな下限を使う.
- 厳密 / 近似的内部最小化のもとでの超 1 次収束.
  - スパース正則化の状況に特殊化することで, 現実的・使いやすい条件のもとで, 既存の収束レートを改善.
- 実験結果も理論をサポート
  - 既存手法 OWLQN, SpaRSA, and FISTA より速い. かつより複雑な正則化に対応可能.

# Summary

- なぜスパース正則化は難しい? – 変数の間の絡み
  - 微分不可能性は悪くない.
  - 内部最小化はスパースであればあるほど速い.  
微分不可能性は効率的な内部最小化のために使える
- いかに変数の間の絡みを取り除くか
  - 線形近似の代わりに パラメトリックな下限を使う.
- 厳密 / 近似的内部最小化のもとでの超 1 次収束.
  - スパース正則化の状況に特殊化することで, 現実的・使いやすい条件のもとで, 既存の収束レートを改善.
- 実験結果も理論をサポート
  - 既存手法 OWLQN, SpaRSA, and FISTA より速い. かつより複雑な正則化に対応可能.

# Benchmark datasets (dorothea)

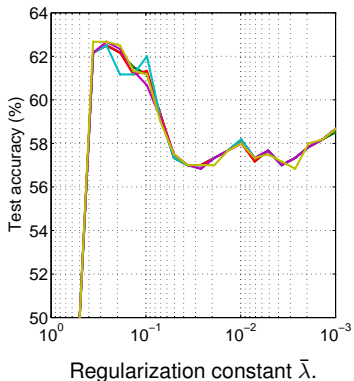
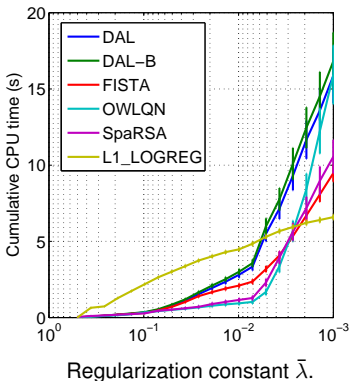
$m = 800, n = 100,000.$





# Benchmark datasets (madelon)

$m = 2,000, n = 500.$



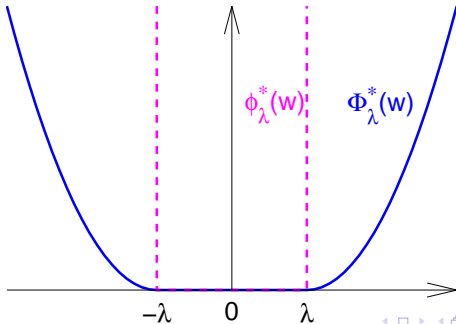
(1) Proximation wrt  $\phi_\lambda$  is analytic (though non-smooth):

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

(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 Rockafellar 76

## Assumption

The multifunction  $\nabla f^*$  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  $\mu$ .

# 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  $\mu$ .

# 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  $\mu$ .

# 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,

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