

# スパース正則化およびマルチカーネル学習のための最適化アルゴリズムとCV・PRへの応用

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

<sup>1</sup> 東京大学

<sup>2</sup> 東京工業大学

2009-08-31 @ PRMU/CVIM 仙台

<http://www.ibis.t.u-tokyo.ac.jp/ryotat/prmu09/>

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

## Lasso 回帰 (1/3)

- 入出力の組み  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$  .  $\mathbf{x}_i \in \mathbb{R}^n$  .
- 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 :

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

- **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 )

- 問題 1:

$$\text{minimize } \|\mathbf{w}\|_0, \quad \text{subject to } L(\mathbf{w}) \leq C.$$

- 問題 2:

$$\text{minimize } L(\mathbf{w}), \quad \text{subject to } \|\mathbf{w}\|_0 \leq C'.$$

どちらも NP 困難!!

## Lasso 回帰 (1/3)

- 入出力の組み  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$  .  $\mathbf{x}_i \in \mathbb{R}^n$  .
- 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 :

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

- **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 )

- 問題 1:

$$\text{minimize } \|\mathbf{w}\|_0, \quad \text{subject to } L(\mathbf{w}) \leq C.$$

- 問題 2:

$$\text{minimize } L(\mathbf{w}), \quad \text{subject to } \|\mathbf{w}\|_0 \leq C'.$$

どちらも NP 困難!!

## Lasso 回帰 (1/3)

- 入出力の組み  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$  .  $\mathbf{x}_i \in \mathbb{R}^n$  .
- 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 :

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

- **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 )

- 問題 1:

$$\text{minimize} \quad \|\mathbf{w}\|_0, \quad \text{subject to } L(\mathbf{w}) \leq C.$$

- 問題 2:

$$\text{minimize} \quad L(\mathbf{w}), \quad \text{subject to } \|\mathbf{w}\|_0 \leq C'.$$

どちらも NP 困難!!

# Lasso 回帰 (1/3)

- 入出力の組み  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$  .  $\mathbf{x}_i \in \mathbb{R}^n$  .
- 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 :

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

- **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 )

- 問題 1:

$$\text{minimize } \|\mathbf{w}\|_0, \quad \text{subject to } L(\mathbf{w}) \leq C.$$

- 問題 2:

$$\text{minimize } L(\mathbf{w}), \quad \text{subject to } \|\mathbf{w}\|_0 \leq C'.$$

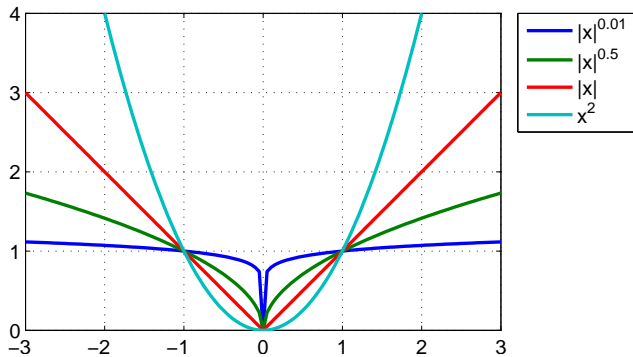
どちらも NP 困難!!



## Lasso 回帰 (2/3)

- $p$ -ノルムの  $p$  乗 ( のようなもの )

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

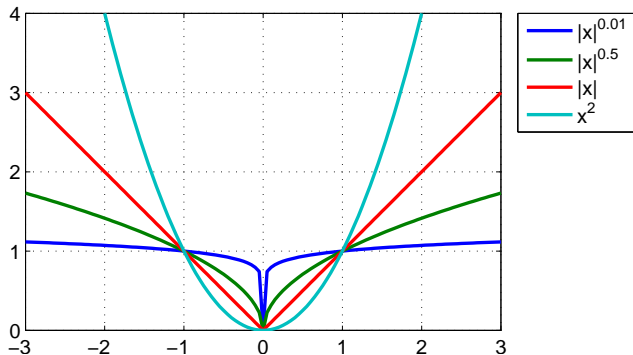


$\|w\|_1$  の正則化は凸の中ではもっとも  $\|w\|_0$  の正則化に近い!

## Lasso 回帰 (2/3)

- $p$ -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 (3/3)

## ● 問題 1:

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

## ● 問題 2:

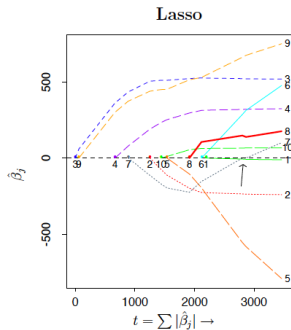
$$\text{minimize } L(\mathbf{w}), \quad \text{subject to } \|\mathbf{w}\|_1 \leq C'.$$

## ● 問題 3:

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

注意:

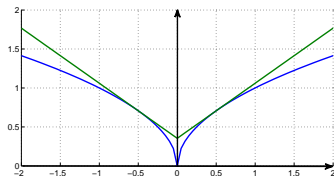
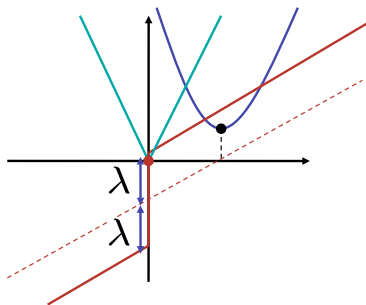
- 上の3つの問題はいずれも等価。
- 正則化項やロス項に単調な非線形変換をしても等価。
- この発表では問題3を扱う。



[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)

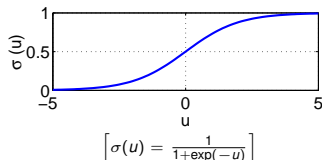


# 一般化

- 損失項を一般化... 例:  $l_1$ -ロジスティック回帰

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^m -\log P(y_i | \mathbf{x}_i; \mathbf{w}) + \lambda \|\mathbf{w}\|_1$$

$$\text{ただし, } P(y | \mathbf{x}; \mathbf{w}) = \sigma(y \langle \mathbf{w}, \mathbf{x} \rangle) \\ (y \in \{-1, +1\})$$



- 正則化項を一般化... 例: グループラッソー (Yuan&Lin,06)

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \sum_{g \in \mathcal{G}} \|\mathbf{w}_g\|_2$$

$$\text{ただし, } \mathcal{G} \text{ は } \{1, \dots, n\} \text{ の適当な分割で, } \mathbf{w} = \begin{pmatrix} (\mathbf{w}_{g_1}) \\ (\mathbf{w}_{g_2}) \\ \vdots \\ (\mathbf{w}_{g_q}) \end{pmatrix}, q = |\mathcal{G}|.$$

# さらに一般化：カーネルを導入

マルチカーネル学習 (Multiple Kernel Learning, MKL): (Lanckriet, Bach, et al., 04)

$\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_n$  を RKHS,  $K_1, K_2, \dots, K_n$  をそれに付随するカーネル関数とする.  $f = \underbrace{f_1}_{\in \mathcal{H}_1} + \underbrace{f_2}_{\in \mathcal{H}_2} + \dots + \underbrace{f_n}_{\in \mathcal{H}_n}$  で予測することで性能を上げよう!

$$\underset{f_j \in \mathcal{H}_j, b \in \mathbb{R}}{\text{minimize}} \quad L(f_1 + f_2 + \dots + f_n + b) + \lambda \sum_{j=1}^n \|f_j\|_{\mathcal{H}_j}$$

↓ (表現定理)

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad f_\ell \left( \sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1} \right) + \lambda \sum_{j=1}^n \|\alpha_j\|_{\mathcal{K}_j}$$

ただし,  $\|\alpha_j\|_{\mathcal{K}_j} = \sqrt{\alpha_j^\top \mathbf{K}_j \alpha_j}$ .

... カーネル行列で重み付けされたグループラッソー .

# さらに一般化：カーネルを導入

マルチカーネル学習 (Multiple Kernel Learning, MKL): (Lanckriet, Bach, et al., 04)

$\mathcal{H}_1, \mathcal{H}_2, \dots, \mathcal{H}_n$  を RKHS ,  $K_1, K_2, \dots, K_n$  をそれに付随するカーネル関数とする .  $f = \underbrace{f_1}_{\in \mathcal{H}_1} + \underbrace{f_2}_{\in \mathcal{H}_2} + \dots + \underbrace{f_n}_{\in \mathcal{H}_n}$  で予測することで性能を上げよう!

$$\underset{f_j \in \mathcal{H}_j, b \in \mathbb{R}}{\text{minimize}} \quad L(f_1 + f_2 + \dots + f_n + b) + \lambda \sum_{j=1}^n \|f_j\|_{\mathcal{H}_j}$$

↓ (表現定理)

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} \quad f_{\ell} \left( \sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1} \right) + \lambda \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}$$

ただし,  $\|\alpha_j\|_{\mathbf{K}_j} = \sqrt{\alpha_j^{\top} \mathbf{K}_j \alpha_j}$  .

… カーネル行列で重み付けされたグループラッソー .

# 絞り込み

損失項は（多くの場合）損失関数  $f_\ell$  とデザイン行列  $\mathbf{A}$  に分解可能．

- 2乗損失

$$f_\ell^Q(\mathbf{z}) = \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2, \quad \mathbf{A} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix}$$

$$f_\ell^Q(\mathbf{A}\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2$$

- ロジスティック損失

$$f_\ell^L(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-y_i z_i)), \quad \mathbf{A} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_m^\top \end{pmatrix}$$

$$f_\ell^L(\mathbf{A}\mathbf{w}) = \sum_{i=1}^m -\log \sigma(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$$



# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - **問題設定**
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# 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。

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

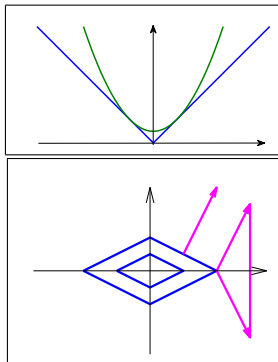
ただし,

- $\mathbf{A} \in \mathbb{R}^{m \times n}$  ( $m$ : 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:  $\mathbf{A}$  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？

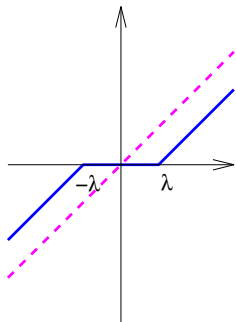
我々の見方:  $\mathbf{A}$  が変数の間にからみを導入するのが原因 .

$\mathbf{A} = \mathbf{I}_n$  (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.

# 先行研究

Iterative Shrinkage/Thresholding (IST) 法 (Figueiredo&Nowak, 03; Daubechies et al., 04,...):

## アルゴリズム

- 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 :

$$\mathbf{w}^{t+1} \leftarrow \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} (Q_{\eta^t}(\mathbf{w}; \mathbf{w}^t) + \phi_{\lambda}(\mathbf{w}))$$

ただし ,

$$Q_{\eta}(\mathbf{w}; \mathbf{w}^t) = \underbrace{L(\mathbf{w}^t) + \nabla L^{\top}(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t)}_{(1) \text{ 損失項を 1 次近似}} + \frac{1}{2\eta} \underbrace{\|\mathbf{w} - \mathbf{w}^t\|_2^2}_{(2) \text{ 前の反復からの距離}^2 \text{ にペナルティ}} .$$

注意 :  $Q_{\eta}(\mathbf{w}; \mathbf{w}^t)$  の最小化は普通の勾配ステップ (サイズ  $\eta$ ) を与える .

# 先行研究：IST法

$$\begin{aligned}
 & \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w})) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \operatorname{ST}_{\eta_t \lambda}(\tilde{\mathbf{w}}^t)
 \end{aligned}$$

この最小化は解析的にできると仮定

ただし， $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (勾配ステップ先)。  
 結局

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\operatorname{ST}_{\eta_t \lambda}}_{\text{縮小}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{勾配ステップ}} \right)$$

- 長所：簡単。
- 短所： $\mathbf{A}$ の悪スケーリングに弱い。

# 先行研究 : IST 法

$$\begin{aligned}
 & \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w})) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \text{ST}_{\eta_t \lambda}(\tilde{\mathbf{w}}^t)
 \end{aligned}$$

この最小化は解析的にできると仮定

ただし,  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (勾配ステップ先).

結局

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t \lambda}}_{\text{縮小}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{勾配ステップ}} \right)$$

- 長所 : 簡単 .
- 短所 :  $\mathbf{A}$  の悪スケーリングに弱い .

# 先行研究 : IST 法

$$\begin{aligned}
 & \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w})) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \text{ST}_{\eta_t\lambda}(\tilde{\mathbf{w}}^t)
 \end{aligned}$$

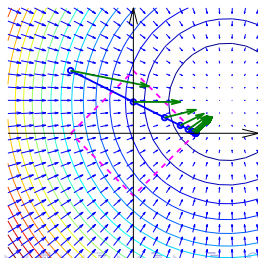
この最小化は解析的にできると仮定

ただし,  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (勾配ステップ先) .

結局

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t\lambda}}_{\text{縮小}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{勾配ステップ}} \right)$$

- 長所 : 簡単 .
- 短所 :  $\mathbf{A}$  の悪スケーリングに弱い .





# 先行研究：IST法

$$\begin{aligned}
 & \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (Q_{\eta_t}(\mathbf{w}; \mathbf{w}^t) + \phi_\lambda(\mathbf{w})) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \text{const.} + \nabla L^\top(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \frac{1}{2\eta_t} \|\mathbf{w} - \tilde{\mathbf{w}}^t\|_2^2 + \phi_\lambda(\mathbf{w}) \right) =: \text{ST}_{\eta_t\lambda}(\tilde{\mathbf{w}}^t)
 \end{aligned}$$

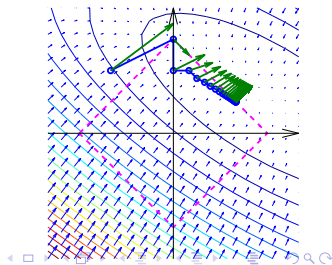
この最小化は解析的にできると仮定

ただし,  $\tilde{\mathbf{w}}^t = \mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)$  (勾配ステップ先) .

結局

$$\mathbf{w}^{t+1} \leftarrow \underbrace{\text{ST}_{\eta_t\lambda}}_{\text{縮小}} \left( \underbrace{\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)}_{\text{勾配ステップ}} \right)$$

- 長所：簡単 .
- 短所： $\mathbf{A}$  の悪スケーリングに弱い .



# ここまでのまとめ

最適化問題：

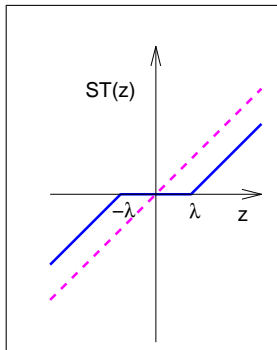
$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}).$$

を解きたい。ただし，

- $f_\ell$  は凸で 2 階微分可能。
- $\phi_\lambda(\mathbf{w})$  は例えば  $\phi_\lambda(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$  で，以下の最小化が陽に求まるもの。

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

- $\phi_\lambda$  の微分不可能性  $\Rightarrow$  ST による打ち切り． **スパース性が計算効率を高める！**
- $\mathbf{A}$  のスケージングにロバストにしたい．



# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

## Proximal Minimization (Rockafellar, 1976)

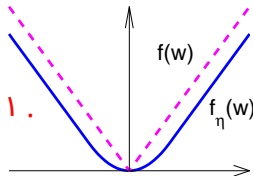
- 1 適当に初期解  $\mathbf{w}^1$  を決める .
- 2 停止条件が満たされるまで反復 :

$$\mathbf{w}^{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( \underbrace{f_\ell(\mathbf{A}\mathbf{w})}_{\text{線形近似しない}} + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \underbrace{\|\mathbf{w} - \mathbf{w}^t\|_2^2}_{\substack{\text{前の反復からの} \\ \text{距離}^2 \text{ にペナルティ}}} \right)$$

$$\bullet f_\eta(\mathbf{w}) = \min_{\tilde{\mathbf{w}} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\tilde{\mathbf{w}}) + \phi_\lambda(\tilde{\mathbf{w}}) + \frac{1}{2\eta} \|\tilde{\mathbf{w}} - \mathbf{w}\|_2^2 \right)$$

とおくと ,

- 事実 1:  $f_\eta(\mathbf{w}) \leq f(\mathbf{w}) = f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w})$  .
- 事実 2:  $f_\eta(\mathbf{w}^*) = f(\mathbf{w}^*)$  .
- このままではもとの最適化問題と同程度に難しい .



## IST 法 (既存手法)

- 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 :

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

## Dual Augmented Lagrangian 法 (提案法)

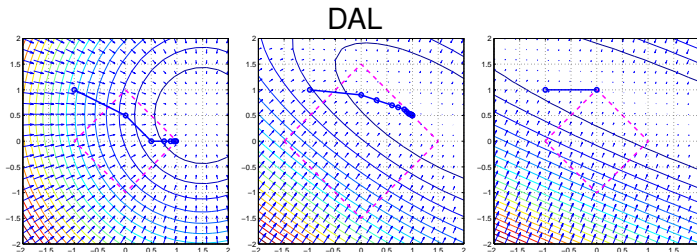
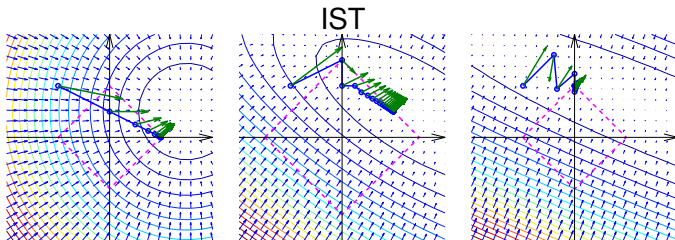
- 1 適当に初期解  $\mathbf{w}^1$  を決める .
- 2 停止条件が満たされるまで反復 :

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

ただし ,

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

## 数値例



## Dual Augmented Lagrangian 法

- 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 :

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

ただし ,

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

## 凸共役 ( Legendre 変換 )

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} \left( \mathbf{y}^\top \mathbf{x} - f(\mathbf{x}) \right)$$



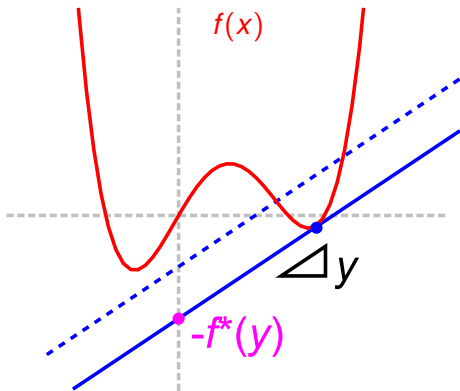
# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# Legendre 変換

関数  $f(x)$  を関数  $f^*(y)$  に移す変換 (フーリエ変換のようなもの)

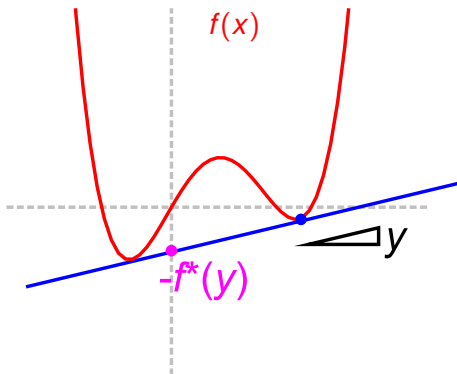
$$f^*(y) = \sup_x (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$$



# Legendre 変換

関数  $f(x)$  を関数  $f^*(y)$  に移す変換 (フーリエ変換のようなもの)

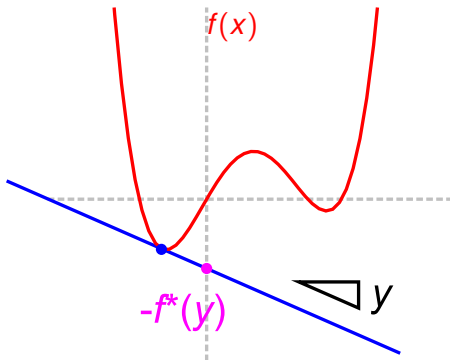
$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$$



# Legendre 変換

関数  $f(x)$  を関数  $f^*(y)$  に移す変換 (フーリエ変換のようなもの)

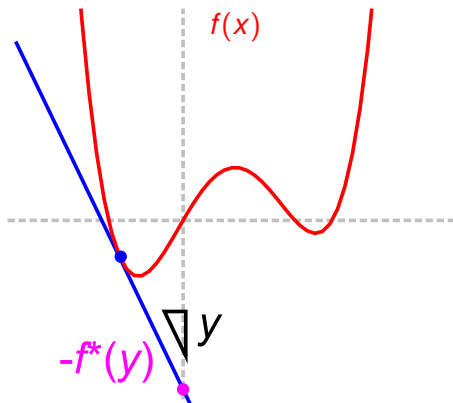
$$f^*(y) = \sup_x \left( y^\top x - f(x) \right)$$



# Legendre 変換

関数  $f(x)$  を関数  $f^*(y)$  に移す変換 (フーリエ変換のようなもの)

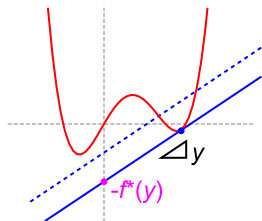
$$f^*(y) = \sup_x (y^\top x - f(x))$$



## Legendre 変換

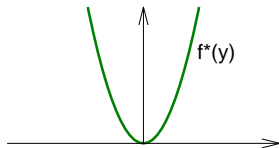
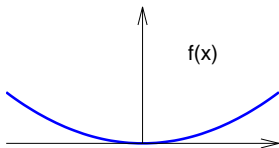
関数  $f(x)$  を関数  $f^*(y)$  に移す変換  
(フーリエ変換のようなもの)

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} \left( \mathbf{y}^\top \mathbf{x} - f(\mathbf{x}) \right)$$



例 1 (2乗誤差関数):

$$f(x) = \frac{x^2}{2\sigma^2} \Rightarrow f^*(y) = \sup_x \left( xy - \frac{x^2}{2\sigma^2} \right) = \left( (\sigma^2 y)y - \frac{(\sigma^2 y)^2}{2\sigma^2} \right) = \frac{\sigma^2 y^2}{2}$$



## Legendre 変換の性質

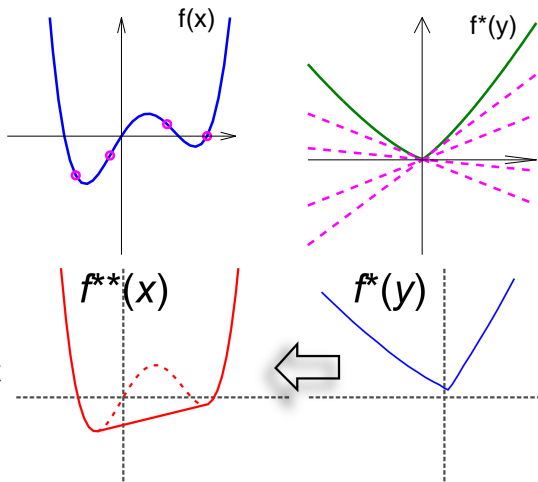
- $f^*(y)$  はいつも凸 .

$$\because f^*(y) \geq \underbrace{xy - f(x)}_{\text{線形な下限}}$$

- $f$  が凸なら  $f^{**}(x) = f(x)$  .

$$\because f(x) \geq xy - f^*(y)$$

(ただし等号が成り立つとは限らない)



Legendre 変換:  $f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$ 

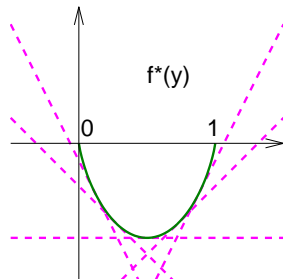
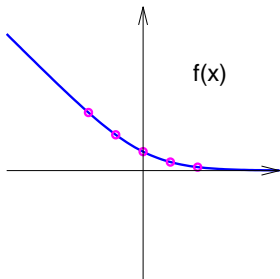
例 2 (ロジスティック損失関数):

$$f(x) = \log(1 + \exp(-x))$$

$$\Rightarrow f^*(-y) = \sup_x (-xy - \log(1 + \exp(-x)))$$

$$= \left( -(\log \frac{1-y}{y})y - \log(1 + \frac{y}{1-y}) \right)$$

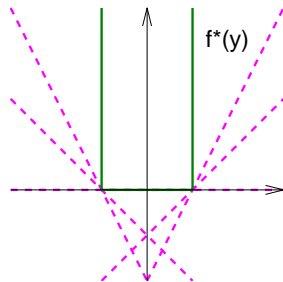
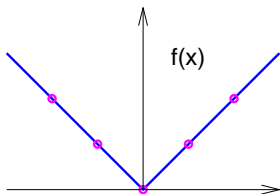
$$= y \log(y) + (1-y) \log(1-y) \quad (\text{エントロピーの符号反転})$$





Legendre 変換:  $f^*(\mathbf{y}) = \sup_{\mathbf{x}} (\mathbf{y}^\top \mathbf{x} - f(\mathbf{x}))$ 例 3 ( $\ell_1$  正則化関数):

$$f(x) = |x| \Rightarrow f^*(y) = \sup_x (xy - |x|) = \begin{cases} 0 & (|y| \leq 1), \\ +\infty & (\text{otherwise}). \end{cases}$$



# 双対を使う

## 主問題の目的関数を双対を使って表現

$$f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) = \max_{\mathbf{v} \in \mathbb{R}^n, \alpha \in \mathbb{R}^m} \left( \underbrace{\mathbf{w}^\top (\mathbf{v} - \mathbf{A}^\top \alpha) - (f_\ell^*(-\alpha) + \phi_\lambda^*(\mathbf{v}))}_{\mathbf{w} \text{ に関し線形な下限}} \right)$$

$$\begin{aligned} \therefore & \max_{\mathbf{v} \in \mathbb{R}^n, \alpha \in \mathbb{R}^m} \left( \mathbf{w}^\top (\mathbf{v} - \mathbf{A}^\top \alpha) - f_\ell^*(-\alpha) - \phi_\lambda^*(\mathbf{v}) \right) \\ &= \max_{\alpha \in \mathbb{R}^m} \left( -\mathbf{w}^\top \mathbf{A}^\top \alpha - f_\ell^*(-\alpha) \right) + \max_{\mathbf{v} \in \mathbb{R}^n} \left( \mathbf{w}^\top \mathbf{v} - \phi_\lambda^*(\mathbf{v}) \right) \\ &= f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) \end{aligned}$$

## Proximal minimization の式に代入

$$\begin{aligned}
 \mathbf{w}^{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \max_{\mathbf{v} \in \mathbb{R}^n, \boldsymbol{\alpha} \in \mathbb{R}^m} \left( -f_\ell^*(-\boldsymbol{\alpha}) - \phi_\lambda^*(\mathbf{v}) + \mathbf{w}^\top (\mathbf{v} - \mathbf{A}^\top \boldsymbol{\alpha}) \right) \right. \\
 &\quad \left. + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right\}
 \end{aligned}$$

min と max の順番を交換し，あとは計算，計算 ...

## Proximal minimization の式に代入

$$\begin{aligned}
 \mathbf{w}^{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \max_{\mathbf{v} \in \mathbb{R}^n, \boldsymbol{\alpha} \in \mathbb{R}^m} \left( -f_\ell^*(-\boldsymbol{\alpha}) - \phi_\lambda^*(\mathbf{v}) + \mathbf{w}^\top (\mathbf{v} - \mathbf{A}^\top \boldsymbol{\alpha}) \right) \right. \\
 &\quad \left. + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right\}
 \end{aligned}$$

min と max の順番を交換し，あとは計算，計算 ...

## Proximal minimization の式に代入

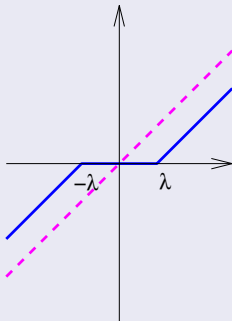
$$\begin{aligned}
 \mathbf{w}^{t+1} &\leftarrow \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left( f_\ell(\mathbf{A}\mathbf{w}) + \phi_\lambda(\mathbf{w}) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right) \\
 &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \left\{ \max_{\mathbf{v} \in \mathbb{R}^n, \boldsymbol{\alpha} \in \mathbb{R}^m} \left( -f_\ell^*(-\boldsymbol{\alpha}) - \phi_\lambda^*(\mathbf{v}) + \mathbf{w}^\top (\mathbf{v} - \mathbf{A}^\top \boldsymbol{\alpha}) \right) \right. \\
 &\quad \left. + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right\}
 \end{aligned}$$

min と max の順番を交換し，あとは計算，計算 ...

## Dual Augmented Lagrangian 法

(1) この計算は解析的にできると仮定 .

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



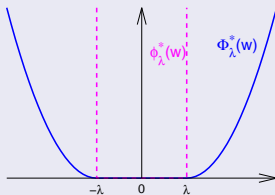
## Dual Augmented Lagrangian 法

(1) この計算は解析的にできると仮定 .

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

(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 )

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



## Dual Augmented Lagrangian 法

(1) この計算は解析的にできると仮定 .

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

(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 .



# Outline

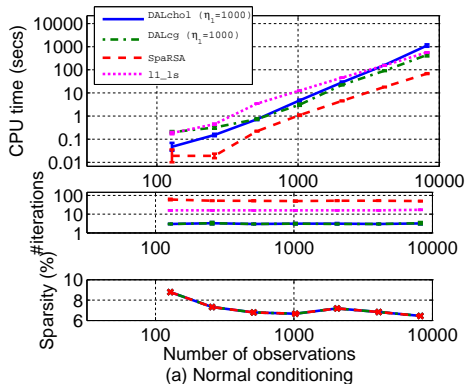
- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - **実験評価**
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# 実験 (設定)

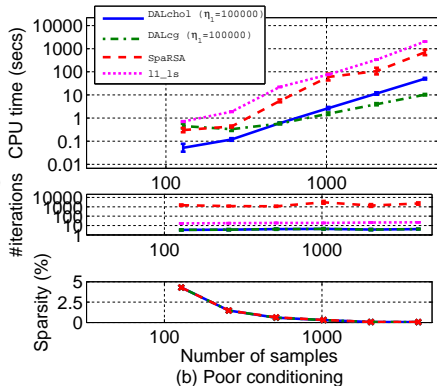
- 問題: LASSO (2乗損失 +  $\ell_1$  正則化)
- 対抗手法:
  - `l1_ls` (内点法)
  - `SpaRSA` (ステップサイズ改良 IST)
- ランダムデザイン行列  $A \in \mathbb{R}^{m \times n}$  ( $m$ : サンプル数  $n$ : 未知変数の数)
  - $A = \text{randn}(m, n)$ ; (良条件数)
  - $A = U * \text{diag}(1 ./ (1:m)) * V'$ ; (悪条件数)
- 2つの状況
  - 中規模 ( $n = 4m, n < 10000$ )
  - 大規模 ( $m = 1024, n < 1e+6$ )

## 結果 (中規模)

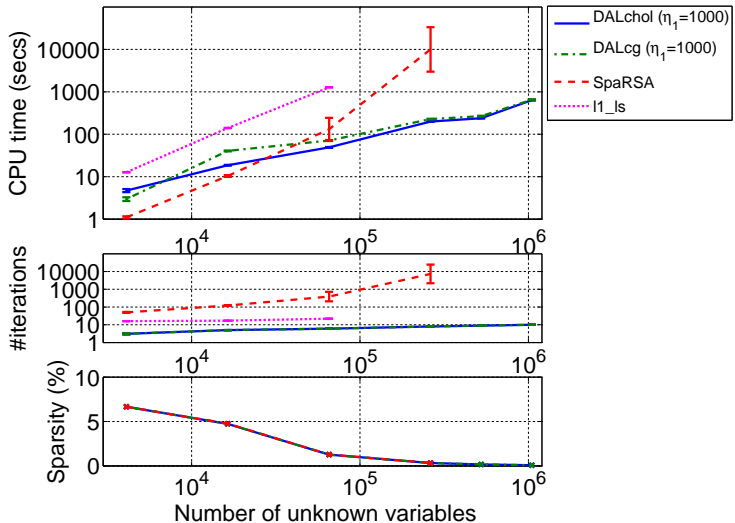
良条件数



悪条件数



## 結果 (大規模)



# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - **マルチカーネル学習**
- 3 デモンストレーション
- 4 まとめ

# 目的

本来やりたいこと :

$$\underset{f \in \mathcal{H}, b \in \mathbb{R}, \mathbf{d} \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^m \xi_i + \frac{\lambda}{2} \|f\|_{\mathcal{H}(\mathbf{d})}^2$$

$$\text{subject to} \quad y_i(f(x_i) + b) \geq 1 - \xi_i \quad (i = 1, \dots, m)$$

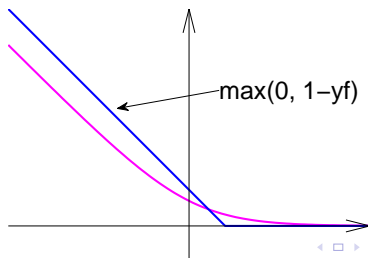
$$\mathbf{K}(\mathbf{d}) = \sum_{i=1}^n d_j \mathbf{K}_j, \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

## 表現定理を適用

$$\underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, d \in \mathbb{R}^n}{\text{minimize}} \quad \sum_{i=1}^m \xi_i + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d) \alpha$$

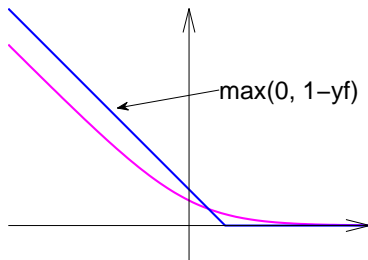
$$\text{subject to} \quad y_i ((\mathbf{K}(d) \alpha)_i + b) \geq 1 - \xi_i \quad (i = 1, \dots, m)$$

$$\mathbf{K}(d) = \sum_{j=1}^n d_j \mathbf{K}_j, \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$



# 損失関数を定義

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, d \in \mathbb{R}^n}{\text{minimize}} && L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha \\ & \text{subject to} && \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$





## 正則化項を緩和

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} && L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha \\ & \text{subject to} && \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_i \geq 0, \quad \sum_i d_i \leq 1. \end{aligned}$$

補助変数  $\alpha_j$  ( $j = 1, \dots, n$ ) を導入して正則化項を緩和

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{subject to} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

ラグランジュ乗数  $\beta$  を導入し,

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right)$$

を最小化.  $\alpha_j = d_j \beta$ ,  $\beta = \alpha$  を得る.

## 正則化項を緩和

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} && L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha \\ & \text{subject to} && \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_i \geq 0, \quad \sum_i d_i \leq 1. \end{aligned}$$

補助変数  $\alpha_j$  ( $j = 1, \dots, n$ ) を導入して正則化項を緩和

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{subject to} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

ラグランジュ乗数  $\beta$  を導入し,

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right)$$

を最小化.  $\alpha_j = d_j \beta$ ,  $\beta = \alpha$  を得る.

## 正則化項を緩和

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} && L(\mathbf{K}(d)\alpha + b\mathbf{1}) + \frac{\lambda}{2} \alpha^\top \mathbf{K}(d)\alpha \\ & \text{subject to} && \mathbf{K}(d) = \sum_{i=1}^n d_i \mathbf{K}_i, \quad d_i \geq 0, \quad \sum_i d_i \leq 1. \end{aligned}$$

補助変数  $\alpha_j$  ( $j = 1, \dots, n$ ) を導入して正則化項を緩和

$$\alpha^\top \mathbf{K}(d)\alpha = \min_{\alpha_j \in \mathbb{R}^m} \left( \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \right) \quad \text{subject to} \quad \sum_{j=1}^n \mathbf{K}_j \alpha_j = \mathbf{K}(d)\alpha$$

ラグランジュ乗数  $\beta$  を導入し,

$$\frac{1}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} + \beta^\top \left( \mathbf{K}(d)\alpha - \sum_{j=1}^n \mathbf{K}_j \alpha_j \right)$$

を最小化.  $\alpha_j = d_j \beta$ ,  $\beta = \alpha$  を得る.

# 上限を最小化

$$\begin{aligned} & \underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} && L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} \\ & \text{subject to} && d_j \geq 0, \quad \sum_j d_j \leq 1. \end{aligned}$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\| \kappa_j}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\| \kappa_j}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\| \kappa_j \right)^2 \end{aligned}$$

ノルムの線形和

# 上限を最小化

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{subject to} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ ノルムの線形和 .

# 上限を最小化

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{subject to} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ ノルムの線形和 .

# 上限を最小化

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{subject to} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ ノルムの線形和 .

# 上限を最小化

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{subject to} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ ノルムの線形和 .



# 上限を最小化

$$\underset{\alpha_j \in \mathbb{R}^m, b \in \mathbb{R}, \gamma \in \mathbb{R}^n}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j}$$

$$\text{subject to} \quad d_j \geq 0, \quad \sum_j d_j \leq 1.$$

$$\begin{aligned} \sum_{j=1}^n \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j} &= \sum_{j=1}^n d_j \frac{\alpha_j^\top \mathbf{K}_j \alpha_j}{d_j^2} = \sum_{j=1}^n d_j \left( \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \\ &\geq \left( \sum_{j=1}^n d_j \frac{\|\alpha_j\|_{\mathbf{K}_j}}{d_j} \right)^2 \quad \left( \sum_j d_j = 1 \text{ より Jensen の不等式} \right) \\ &= \left( \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \right)^2 \end{aligned}$$

↑ ノルムの線形和 .

## 2つの定式化の同値性

ノルム線形和の2乗で正則化 (一般的な定式化)

$$\underset{\alpha_j \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \frac{\lambda}{2} \left(\sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}\right)^2 \quad (\text{A})$$

ノルム線形和で正則化 (我々の定式化)

$$\Leftrightarrow \underset{\alpha_j \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} \quad L\left(\sum_{j=1}^n \mathbf{K}_j \alpha_j + b \mathbf{1}\right) + \lambda' \sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j} \quad (\text{B})$$

(A) の最適性 :  $\nabla_{\alpha_j} L + \lambda \left(\sum_{j=1}^n \|\alpha_j\|_{\mathbf{K}_j}\right) \partial_{\alpha_j} \|\alpha_j\|_{\mathbf{K}_j} \ni 0$

(B) の最適性 :  $\nabla_{\alpha_j} L + \lambda' \partial_{\alpha_j} \|\alpha_j\|_{\mathbf{K}_j} \ni 0$

# SpicyMKL

DAL + MKL = **SpicyMKL (Sparse Iterative MKL)**

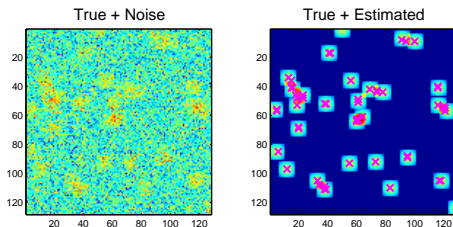
- 基本的には DAL と同じ .
- バイアス項を扱う必要がある .
- ヒンジロス は微分できないので特別に扱う必要がある ( 今回の実験は ロジスティック損失 ) .
- Soft-thresholding が ( 変数単位ではなく ) カーネル単位でかかる .

$$ST_{\lambda}(\alpha_j) = \begin{cases} 0 & (\|\alpha_j\|_{\mathbf{K}_j} \leq \lambda) \\ \left( \|\alpha_j\|_{\mathbf{K}_j} - \lambda \right) \frac{\alpha_j}{\|\alpha_j\|_{\mathbf{K}_j}} & (\text{otherwise}) \end{cases}$$

# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 **デモンストレーション**
- 4 まとめ

## Demo1 – デコンボリューション



- 画像は 128x128 .
- フィルタは  $\sigma = 5$  のガウシアンぼかし .
- コマンド :

```
[xx,stat]=dalsqll(zeros(m*n,1),H,Y(:),lambda,'eta',500,'solver','cg');
```

↑  
 2乗ロス+L1正則化

↗  
 初期値

↘  
 畳み込み行列

↑  
 入力画像

↑  
 正則化定数

{  
 ペナルティーの強さ  
 (の初期値)

{  
 インナーループの最  
 適化にCG法を使う

# Demo2 – バイオインフォマティクス

- 多発性硬化症に対する  $\beta$  インターフェロン療法の効果を検証 .
- 53 人の患者の 70 遺伝子の発現データが投薬開始から最長 2 年間に渡って集められた (t=0, 3, 6, 9, 12, 18, 24ヶ月後)
- 2 値分類問題 (効果的 / 効果なし) → **ロジスティック損失を使う** .
- 2 つの設定
  - 時系列情報を扱うグループラッソーの問題 .
  - 遺伝子の組みを探す MKL の問題 .

## Demo2.1 – グループラッソー

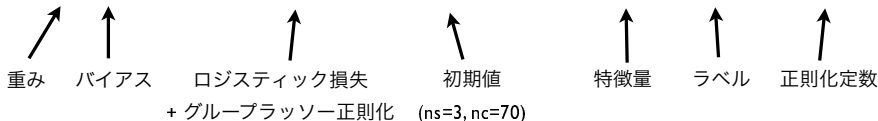
- 各遺伝子ごとに，
  - ① (時間方向の) 平均発現量
  - ② 時間方向 1 階差分の平均
  - ③ 時間方向 2 階差分の平均
 を計算 (3 × 70 次元特徴)

- グループラッソー: 
$$\underset{\mathbf{w} \in \mathbb{R}^{3 \times 70}, b \in \mathbb{R}}{\text{minimize}} \sum_{i=1}^m \ell^L(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) + \lambda \sum_{j=1}^{70} \|\mathbf{w}_j\|_2$$

- Soft-thresholding: 
$$\text{ST}_\lambda(\mathbf{w}_j) = \max(0, \|\mathbf{w}_j\|_2 - \lambda) \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|_2}$$

- コマンド

```
[ww,bb,stat]=dallrgl(zeros(ns,nc),F(:,,:),Y(:,), lambda);
```



## Demo2.2 – MKL

- 時刻 0 (治療開始時) のデータだけを利用 .
- Baranzini らが見つけた遺伝子の 3 つ組 9 つにそれぞれ 2 次の多項式カーネル  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^\top \mathbf{x}_j)^2$  を導入 .

```
opt=struct('loss','logit'); ←ロジスティック損失を指定
```

```
[alpha,d,b,actset]=SpicyMKL(K,Y,lambda,opt);
```

サンプル重み    カーネル重み    バイアス    アクティブセット    カーネル (m x m x n)    ラベル    正則化定数



## Demo3 – 画像認識

- Caltech101 (Fei-Fei et al., 2004) の中から **anchor, ant, cannon, chair, cup** の 5 クラスを利用 .
- 10 通りの 2 クラス分類問題 .
- **カーネル数 1,760** = 特徴抽出法 ( 4 通り ) × 領域分割 ( 22 通り ) × カーネル関数 ( 20 通り )
  - 特徴抽出: van de Sande らのコードを利用 . hsvsift, sift ( スケール自動 ) , sift ( スケール 4px 固定 ) , sift ( スケール 8px 固定 ) の 4 通り .
  - 領域分割と統合: 画像全体 , 4 分割 , 16 分割し , それぞれの領域で visual words の出現頻度を計算 , さらに , それらを spatial pyramid で統合したもの ( 計 22 通り ) .
  - カーネル関数: ガウシアンカーネルと  $\chi^2$  カーネルをそれぞれ 10 通りのハイパーパラメータで用意 .

# Outline

- 1 イントロ - スパース正則化とは
  - 具体例
  - 問題設定
- 2 Dual Augmented Lagrangian 法 (提案法)
  - Proximal minimization からのアプローチ
  - Legendre 変換
  - 実験評価
  - マルチカーネル学習
- 3 デモンストレーション
- 4 まとめ

# まとめ

- $l_1$ -正則化：凸最適化だからといって終わりではない．まだまだ工夫の必要 / 余地がある．
- DAL：スパース性を計算の面でも積極的に使う．
- Legendre 変換：微分を取って線形化 → Legendre 変換で線形化  
… 困ったら下限を作ってみる．
- MKL：「最適化問題」を信用しない．同じ問題を表現する方法は無数にある．

## 謝辞

電通大の柳内先生には画像認識に関して詳細にアドバイス頂き，感謝しています．

## 陰勾配法としての提案法

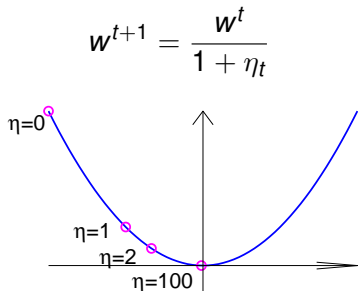
$$\mathbf{w}^{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w}} \left( f(\mathbf{w}) + \frac{1}{2\eta^t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 \right)$$

より,

$$\partial f(\mathbf{w}^{t+1}) + \frac{1}{\eta_t} (\mathbf{w}^{t+1} - \mathbf{w}^t) \ni 0$$

整理すると,

$$\mathbf{w}^{t+1} - \mathbf{w}^t \in -\eta_t \underbrace{\partial f(\mathbf{w}^{t+1})}_{\text{遷移先での勾配}}$$



# Convolution

Inf-convolution:

$$(f \circ g)(x) = \inf_y (f(x - y) + g(y))$$

## 畳み込みと Legendre 変換

$$(f \circ g)^*(\alpha) = f^*(\alpha) + g^*(\alpha)$$

$$\begin{aligned} \therefore (f \circ g)^*(\alpha) &= \sup_x \left( \alpha x - \inf_y (f(x - y) + g(y)) \right) \\ &= \sup_x \sup_y (\alpha x - f(x - y) - g(y)) \\ &= f^*(\alpha) + \sup_y (\alpha y - g(y)) \\ &= f^*(\alpha) + g^*(\alpha) \end{aligned}$$

# Proximity Operation

$f(\mathbf{x})$  は凸で下半連続 .  $f^*(\mathbf{y}) = \sup_{\mathbf{x}}(\mathbf{x}^\top \mathbf{y} - f(\mathbf{x}))$  とする .

Proximity operator (凸集合への射影の一般化)

$$\text{prox}_f(\mathbf{z}) = \inf_{\mathbf{x}} \left( \frac{1}{2} \|\mathbf{z} - \mathbf{x}\|_2^2 + f(\mathbf{x}) \right)$$

Moreau's decomposition (Moreau, 65; Combettes&Wajs, 05)

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