

機械学習における連続最適化の新しいトレンド

富岡 亮太¹

共同研究者: 鈴木 大慈¹ 杉山 将²

¹ 東京大学

² 東京工業大学

2011-01-20 @ NEC

Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting
- 6 まとめ
- 7 Appendix

標準形は好き？

例) 線形計画 (LP)

主問題

$$(P) \quad \min \quad \mathbf{c}^\top \mathbf{x},$$

$$\text{s.t.} \quad \mathbf{Ax} = \mathbf{b}, \mathbf{x} \geq 0.$$

双対問題

$$(D) \quad \max \quad \mathbf{b}^\top \mathbf{y},$$

$$\text{s.t.} \quad \mathbf{A}^\top \mathbf{y} \leq \mathbf{c}.$$

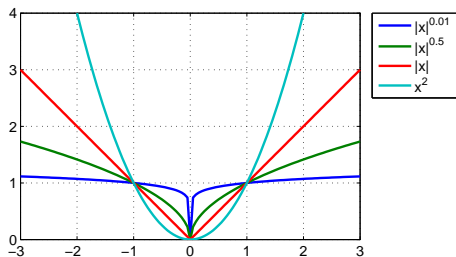
2次計画 (QP), 2次錐計画 (SOCP), 半正定値計画 (SDP), etc...

- **良い点** : 既存の (ある程度) 高性能なソルバーが利用できる .
- **悪い点** : 問題のモデル化に制限, 問題の書き換えが必要 .

モデル化に制限が生じる例： ℓ_p -正則化

$$\underset{\mathbf{w}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \sum_{j=1}^n |w_j|^p.$$

- $p = 2$: 正則化最小二乗法
⇒ 逆行列で一発.
- $1 < p < 2$: ?
- $p = 1$: Lasso ⇒ 2次計画.

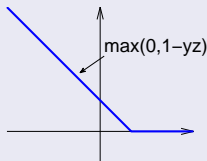


それでも標準形が好まれる理由（これを解決します）

① 微分不可能性: 微分不可能性 < 制約付き最適化

SVM (微分不可能)

$$\min_{\mathbf{w}} C \sum_{i=1}^m \ell_H(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) + \frac{1}{2} \|\mathbf{w}\|^2$$



SVM (2次計画)

$$\begin{aligned} \min_{\mathbf{w}} \quad & C \sum_{i=1}^m \xi_i + \frac{1}{2} \|\mathbf{w}\|^2, \\ \text{s.t.} \quad & y_i \langle \mathbf{x}_i, \mathbf{w} \rangle \geq 1 - \xi_i \\ & (i = 1, \dots, m). \end{aligned}$$

(制約の数 = サンプル数)

② 実装がめんどくさい

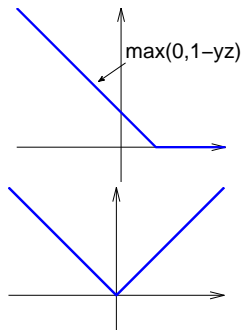
微分不可能性を持つ凸最適化問題の例

- SVM（ロスが微分不可能）

$$\underset{\mathbf{w}}{\text{minimize}} \quad C \sum_{i=1}^m \ell_H(y_i \langle \mathbf{x}_i, \mathbf{w} \rangle) + \frac{1}{2} \|\mathbf{w}\|^2$$

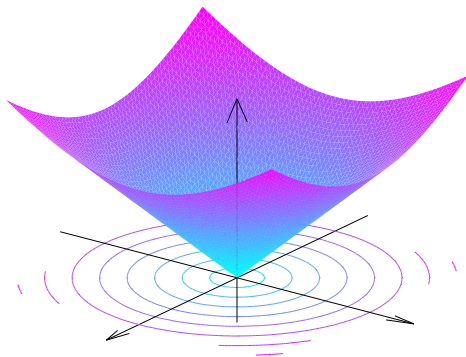
- Lasso（正則化項が微分不可能）

$$\underset{\mathbf{w}}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \sum_{j=1}^n |w_j|$$



マルチタスク学習 (Evgeniou et al 05)

$$\underset{\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_{12}}{\text{minimize}} \underbrace{L(\mathbf{w}_1 + \mathbf{w}_{12})}_{\text{タスク 1 のロス}} + \underbrace{L(\mathbf{w}_2 + \mathbf{w}_{12})}_{\text{タスク 2 のロス}} + \lambda(\|\mathbf{w}_1\| + \|\mathbf{w}_2\| + \|\mathbf{w}_{12}\|)$$



マルチカーネル学習 (Bach et al 05; Suzuki & Tomioka 09)

- M 個の情報源があり，それぞれについて入力サンプル x_i と x_j の内積が $k_m(x_i, x_j)$ ($m = 1, \dots, M$) と計算できるとき，それらをいかに組み合わせさせて予測するか？
- 最適化問題

$$\underset{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_M}{\text{minimize}} \quad C \sum_{i=1}^N \ell_H \left(y_i \sum_{m=1}^M f_m(x_i) \right) + \lambda \sum_{m=1}^M \|\mathbf{f}_m\|$$

- それぞれの情報源について，表現定理より予測関数

$$f_m(x_i) = \sum_{j=1}^N k_m(x_i, x_j) \alpha_j^{(m)}.$$

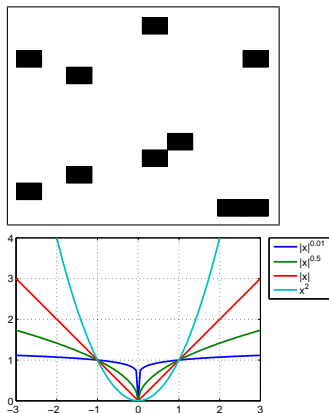
- 正則化項だけでなくロスも微分不可能．

行列穴埋め

- 部分的に観測された行列の未観測部分を低ランク性の仮定を使って予測する問題． \Rightarrow 協調フィルタリング．

$$\underset{\mathbf{X}}{\text{minimize}} \quad \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2 + \lambda \underbrace{\sum_{j=1}^r \sigma_j(\mathbf{X})}_{\text{特異値の線形和}} .$$

- Ω は見えている部分だけを取り出す操作．



行列の空間上の判別問題

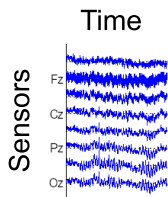
判別ラベル

予測関数

$$\mathbf{y} \in \{-1, +1\} \quad \Leftarrow \quad f(\mathbf{X}) = \langle \mathbf{W}, \mathbf{X} \rangle + b$$

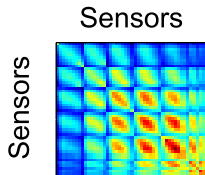
- Multivariate Time Series

$$\mathbf{X} =$$

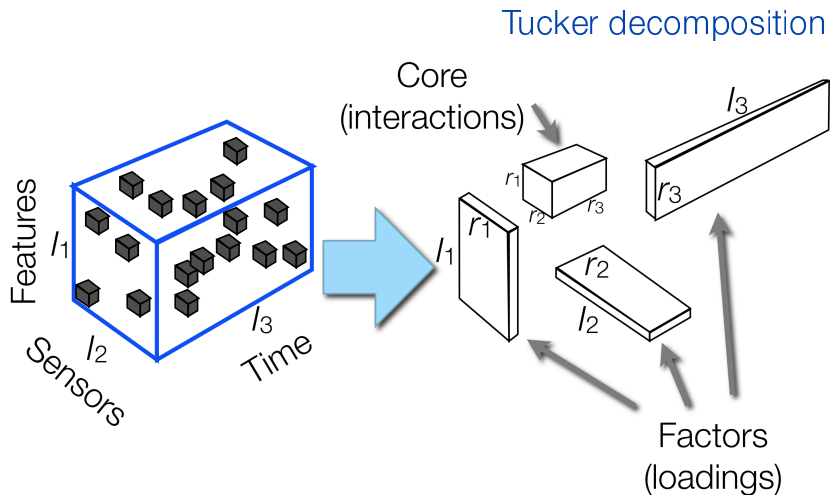


- Second order statistics

$$\mathbf{X} =$$



テンソルの穴埋め



アジェンダ

(割と) 簡単に微分不可能な最適化問題を解くテクニック

- Prox 作用素 (proximity operator)
- 凸関数の共役 (Legendre 変換)
- Operator splitting

Outline

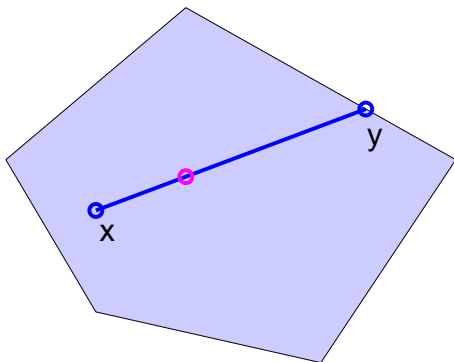
- 1 イントロ
- 2 準備**
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting
- 6 まとめ
- 7 Appendix

凸集合

\mathbb{R}^n の部分集合 V は凸集合

$\Leftrightarrow V$ の任意の 2 点 \mathbf{x}, \mathbf{y} を結ぶ線分が V に入っている . つまり

$$\forall \mathbf{x}, \mathbf{y} \in V, \forall \lambda \in [0, 1], \lambda \mathbf{x} + (1 - \lambda) \mathbf{y} \in V.$$



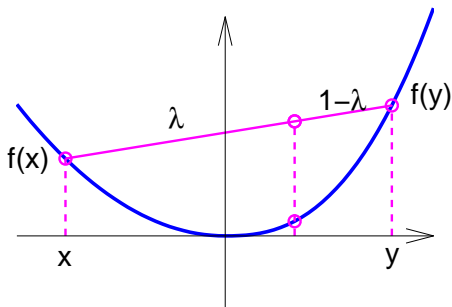
凸関数

関数 $f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ が (拡張値) 凸関数

\Leftrightarrow 関数 f の“弦”が つねに関数自身より上にある．つまり

$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \forall \lambda \in [0, 1], \quad f((1 - \lambda)\mathbf{x} + \lambda\mathbf{y}) \leq (1 - \lambda)f(\mathbf{x}) + \lambda f(\mathbf{y})$$

$<$ が成り立つ場合を **strictly convex** と言う．



拡張値凸関数 - なぜ $+\infty$ を許すか

- $f(\mathbf{x}) = +\infty$ を許すと定義域 (or 制約) を含めて扱える .
例) $1/x$ は $x > 0$ で凸 ,

$$f(x) = \begin{cases} 1/x & (x > 0), \\ +\infty & (\text{otherwise}). \end{cases}$$

と定義すればよい .

定義 (定義関数)

凸集合 C の定義関数 δ_C を以下のように定義する .

$$\delta_C(\mathbf{x}) = \begin{cases} 0 & (\mathbf{x} \in C), \\ +\infty & (\text{otherwise}). \end{cases}$$

このように定義した δ_C は凸関数 .

凸最適化問題

- 拡張値凸関数 f として

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}).$$

- 制約付き問題

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) + \delta_C(\mathbf{x}).$$

- 正則化付き最小化問題

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{L(\mathbf{x}) + \phi_\lambda(\mathbf{x})}_{=: f(\mathbf{x})}.$$

Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素**
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting
- 6 まとめ
- 7 Appendix

Proximal Minimization (Rockafellar 76)

- ① \mathbf{x}^0 を適当に初期化する .
- ② 以下を繰り返す .

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x}} \left(f(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right)$$

- **良い点**: f が凸関数であれば必ず収束 . しかも速い (超 1 次)

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

(Tomioka et al. 11)

- **悪い点**: そもそも $f(\mathbf{x})$ の最小化が難しい時 , 上の最小化を各ステップやるのはつらい .

勾配法

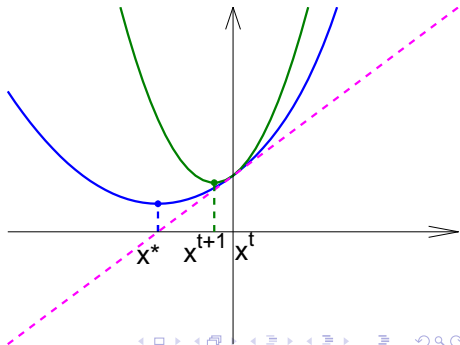
各ステップで $f(\mathbf{x})$ を \mathbf{x}^t の周りで線形化して Proximal Minimization する .

$$\begin{aligned}\mathbf{x}^{t+1} &= \operatorname{argmin}_{\mathbf{x}} \left(\nabla f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right) \\ &= \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)\end{aligned}$$

- 安定性の条件: $\eta_t \leq 1/L(f)$.
- $L(f)$ は微分 ∇f のリプシッツ定数:

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq L(f) \|\mathbf{y} - \mathbf{x}\|.$$

- f が 2 階微分可能な場合は $L(f) =$ ヘシアン の 最大固有値 の 上限 .



0G: 射影付き勾配法 (Bertsekas 99; Nesterov 03)

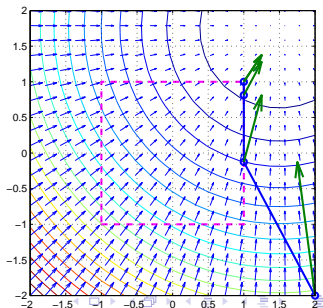
目的関数を線形化, δ_C は束縛条件の定義関数,

$$\begin{aligned}\mathbf{x}^{t+1} &= \operatorname{argmin}_{\mathbf{x}} \left(\nabla f(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) + \delta_C(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right) \\ &= \operatorname{argmin}_{\mathbf{x}} \left(\delta_C(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - (\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t))\|^2 \right) \\ &= \operatorname{proj}_C(\mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t)).\end{aligned}$$

- 安定性の条件は制約のない場合と同様: $\eta_t \leq 1/L(f)$.
- 収束の速さ

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{L(f) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2k}$$

- もちろん射影 proj_C が効率的に計算できることが必要.

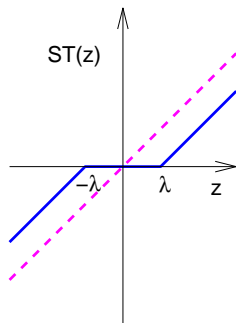


Proximal Operator: 射影の一般化

$$\text{prox}_{\phi_\lambda}(\mathbf{z}) = \underset{\mathbf{x}}{\text{argmin}} \left(\phi_\lambda(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2 \right)$$

- 凸集合への射影: $\text{proj}_C(\mathbf{z}) = \text{prox}_{\delta_C}(\mathbf{z})$.
- Soft-Threshold ($\phi_\lambda(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$)

$$\begin{aligned} \text{prox}_\lambda(\mathbf{z}) &= \underset{\mathbf{x}}{\text{argmin}} \left(\lambda \|\mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2 \right) \\ &= \begin{cases} z_j + \lambda & (z_j < -\lambda), \\ 0 & (-\lambda \leq z_j \leq \lambda), \\ z_j - \lambda & (z_j > \lambda). \end{cases} \end{aligned}$$



- 何らかの意味で**分離可能**な ϕ_λ は Prox が簡単に計算できる .
- 微分不可能でも解析的に計算できる .

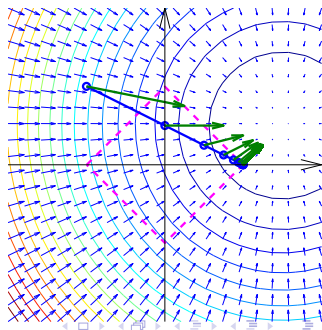
1G: Iterative Shrinkage Thresholding (IST)

$$\begin{aligned}
 \mathbf{x}^{t+1} &= \operatorname{argmin}_{\mathbf{x}} \left(\nabla L(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) + \phi_{\lambda}(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right) \\
 &= \operatorname{argmin}_{\mathbf{x}} \left(\phi_{\lambda}(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - (\mathbf{x}^t - \eta_t \nabla L(\mathbf{x}^t))\|^2 \right) \\
 &= \operatorname{prox}_{\lambda\eta_t}(\mathbf{x}^t - \eta_t \nabla L(\mathbf{x}^t)).
 \end{aligned}$$

- 安定性の条件，収束性は射影付き勾配法と同じ．(Beck & Teboulle 09)

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{L(f) \|\mathbf{x}_0 - \mathbf{x}^*\|^2}{2k}$$

- Prox 作用素 $\operatorname{prox}_{\lambda}$ が計算できれば実装簡単．
- Forward-Backward Splitting と呼ばれる (Lions & Mercier 76)



IST を用いた行列穴埋め (Mazumder et al. 10)

ロス:

$$L(\mathbf{X}) = \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2.$$

微分:

$$\nabla L(\mathbf{X}) = \Omega^\top (\Omega(\mathbf{X} - \mathbf{Y}))$$

反復式:

$$\mathbf{X}^{t+1} = \text{prox}_{\lambda\eta_t} \left((\mathbf{I} - \eta_t \Omega^\top \Omega)(\mathbf{X}^t) + \eta_t \Omega^\top \Omega(\mathbf{Y}^t) \right)$$

- $\eta_t = 1$ の時, 予測値で穴埋め, 観測値で上書き, Soft-Threshold かける, の繰り返し.

正則化項:

$$\phi_\lambda(\mathbf{X}) = \lambda \sum_{j=1}^r \sigma_j(\mathbf{X}) \quad (\text{特異値の線形和}).$$

Prox 作用素 (Singular Value Thresholding):

$$\text{prox}_\lambda(\mathbf{Z}) = \mathbf{U} \max(\mathbf{S} - \lambda \mathbf{I}, 0) \mathbf{V}^\top.$$

FISTA: IST の加速版 (Beck & Teboulle 09; Nesterov 07)

① \mathbf{x}^0 を適当に初期化, $\mathbf{y}^1 = \mathbf{x}^0$, $s_1 = 1$ とする.

② \mathbf{x}^t を更新:

$$\mathbf{x}^t = \text{prox}_{\lambda\eta_t}(\mathbf{y}^t - \eta_t \nabla L(\mathbf{y}^t)).$$

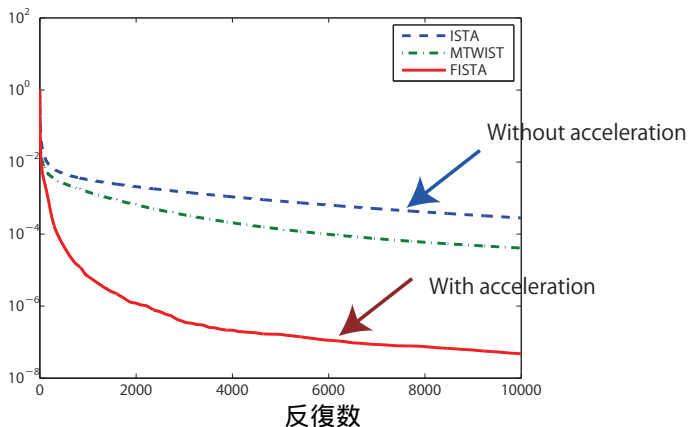
③ \mathbf{y}^t を更新:

$$\mathbf{y}^{t+1} = \mathbf{x}^t + \left(\frac{s_t - 1}{s_{t+1}} \right) (\mathbf{x}^t - \mathbf{x}^{t-1}),$$

$$\text{ただし, } s_{t+1} = (1 + \sqrt{1 + 4s_t^2})/2.$$

- ステップあたりの計算量は IST とほぼ同じ. 収束性は $O(1/k^2)$.
- どこの点で勾配ステップを取るべきか予測している感じ?

加速の効果



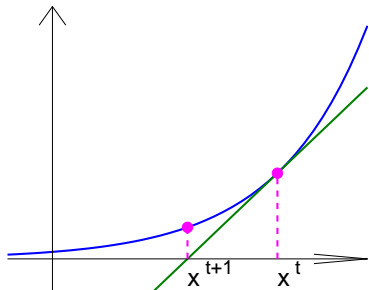
Beck & Teboulle 2009 SIAM J. IMAGING SCIENCES

Vol. 2, No. 1, pp. 183-202 より

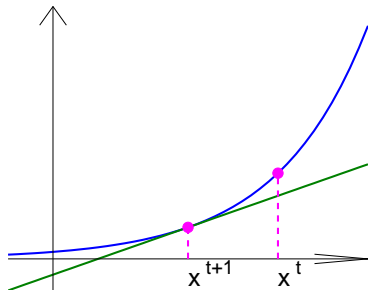
3G?: もっとタイトな下限を作るアイデア

IST:

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x}} \left(\nabla L(\mathbf{x}^t)(\mathbf{x} - \mathbf{x}^t) + \phi_\lambda(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right)$$



IST: 現在の点 \mathbf{x}^t で下限を作る .



DAL: 次の点 \mathbf{x}^{t+1} で下限を作る .

次の点で下限を作るには？

下限の傾き \mathbf{y} をパラメータとして最適化する:

$$\begin{aligned}\mathbf{x}^{t+1} &= \operatorname{argmin}_{\mathbf{x}} \left(L(\mathbf{x}) + \phi_{\lambda}(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right) \\ &= \operatorname{argmin}_{\mathbf{x}} \left(\max_{\mathbf{y}} (\langle \mathbf{x}, \mathbf{y} \rangle - L^*(\mathbf{y})) + \phi_{\lambda}(\mathbf{x}) + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}^t\|^2 \right)\end{aligned}$$

↓ Min と Max の順番を入れ替えて計算する

DAL (Tomioka et al 11)

$$\mathbf{x}^{t+1} = \operatorname{prox}_{\lambda\eta_t} (\mathbf{x}^t - \eta_t \mathbf{y}^t),$$

ただし,

$$\mathbf{y}^t = \operatorname{argmax}_{\mathbf{y}} \left(-L^*(\mathbf{y}) - \frac{1}{\eta_t} \Phi_{\lambda\eta_t}^* (\mathbf{x}^t - \eta_t \mathbf{y}) \right)$$

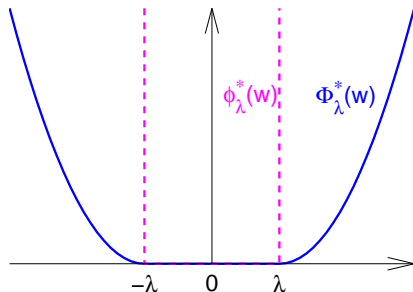
DAL の利点 (ℓ_1 -正則化の場合)

(1) Prox 作用素は解析的に計算可能

$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t \lambda} (\mathbf{x}^t - \eta_t \mathbf{y}^t)$$

(2) 内部最適化は微分可能

$$\mathbf{y}^t = \underset{\mathbf{y}}{\text{argmax}} \left(\underbrace{-L^*(\mathbf{y})}_{\text{微分可能}} - \frac{1}{2\eta_t} \underbrace{\|\text{prox}_{\lambda \eta_t}(\mathbf{x}^t - \eta_t \mathbf{y})\|^2}_{\text{非ゼロ成分の数に比例}} \right)$$



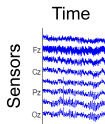
DAL を用いた行列の空間上の判別問題

判別ラベル

予測関数

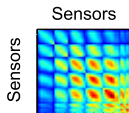
- Multivariate Time Series

$$X =$$



- Second order statistics

$$X =$$



$$\mathbf{y} \in \{-1, +1\} \quad \Leftarrow \quad f(\mathbf{X}) = \langle \mathbf{W}, \mathbf{X} \rangle + b$$

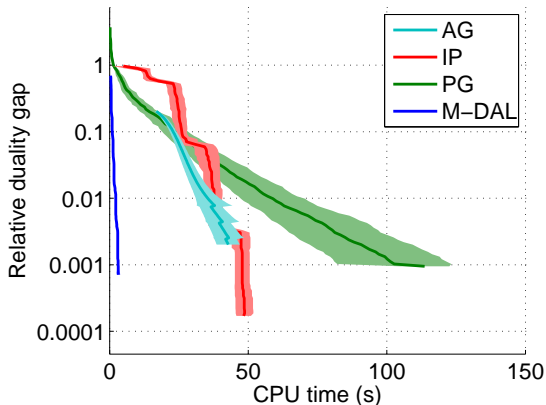
最適化問題:

$$\underset{\mathbf{W}}{\text{minimize}} \quad \sum_{i=1}^m \ell_{\text{LR}}(y_i f(\mathbf{X}_i)) + \lambda \underbrace{\|\mathbf{W}\|_{\text{tr}}}_{\text{低ランク化}}$$

ロジスティック損失:

$$\ell_{\text{LR}}(z_i) = \log(1 + \exp(-z_i))$$

DAL を用いた行列の空間上の判別問題 (Tomioka et al 10)



- AG: 加速付き IST — IP: 内点法
— PG: 射影勾配法 — M-DAL: 提案法

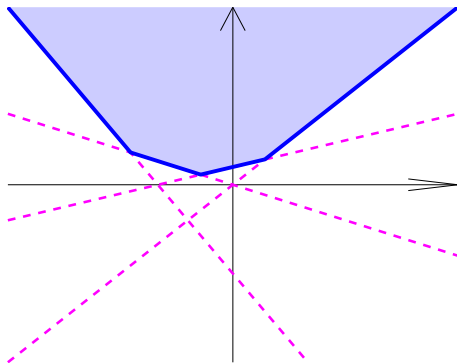
Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換**
- 5 手法 3: Operator Splitting
- 6 まとめ
- 7 Appendix

凸関数の作り方

線形関数の point-wise maximum は凸関数 .

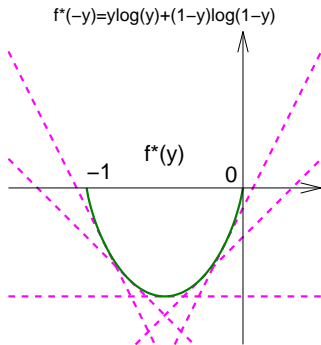
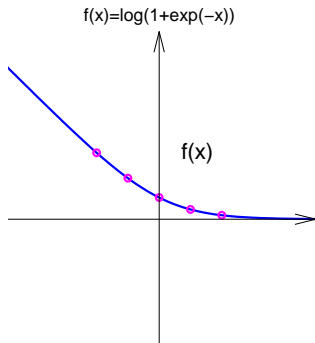
$$f(\mathbf{x}) = \max_{i=1, \dots, k} (\langle \mathbf{a}_i, \mathbf{x} \rangle - b_i).$$



凸共役 (Fenchel-Legendre 変換)

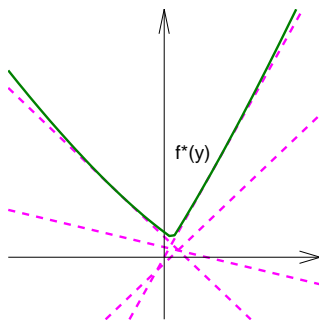
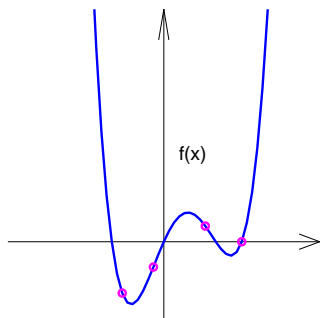
$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathbb{R}^n} (\langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x})).$$

のように定義される f^* を関数 f の凸共役と呼ぶ。



注意: $f(x)$ は凸でなくてもよい

$f(x)$ が凸関数であってもなくても $f^*(y)$ は凸関数 .



凸共役の性質

定義より

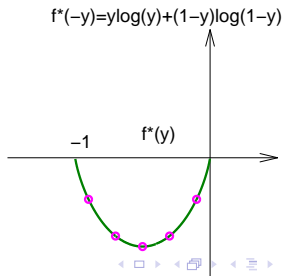
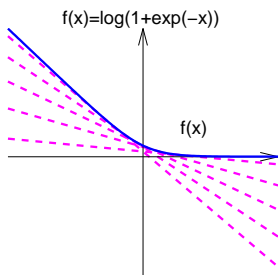
$$\forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n, \quad f(\mathbf{x}) + f^*(\mathbf{y}) \geq \langle \mathbf{y}, \mathbf{x} \rangle.$$

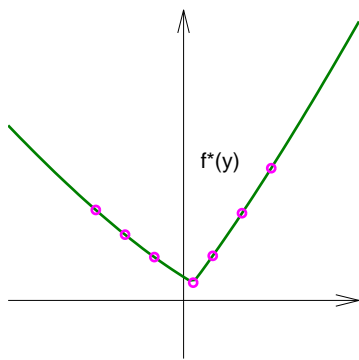
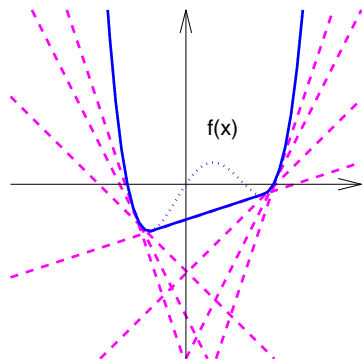
\mathbf{x} がある \mathbf{y} に関して $\langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x})$ を最大化するならば

$$f^*(\mathbf{y}) = \langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x}).$$

つまり,

$$f(\mathbf{x}) = f^{**}(\mathbf{x}) = \sup_{\mathbf{y}} (\langle \mathbf{y}, \mathbf{x} \rangle - f^*(\mathbf{y})).$$



凸でない f の場合 $f \neq f^{**}$ 

- f^* が微分可能 $\Leftrightarrow f$ が strictly convex.

Uzawa's method (Dual ascent; 双対分解)

最適化問題（あえて制約付き問題として書く）:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m}{\text{minimize}} && f(\mathbf{z}) + g(\mathbf{x}), \\ & \text{subject to} && \mathbf{z} = \mathbf{A}\mathbf{x}. \end{aligned}$$

ラグランジアン:

$$L(\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{z}) + g(\mathbf{x}) + \boldsymbol{\alpha}^\top (\mathbf{z} - \mathbf{A}\mathbf{x}).$$

ラグランジアンの最小化は凸共役の計算

$$\begin{aligned} -f^*(-\boldsymbol{\alpha}) &= \min_{\mathbf{z}} (f(\mathbf{z}) + \langle \boldsymbol{\alpha}^t, \mathbf{z} \rangle), \\ -g^*(\mathbf{A}^\top \boldsymbol{\alpha}) &= \min_{\mathbf{x}} (g(\mathbf{x}) - \langle \mathbf{A}^\top \boldsymbol{\alpha}^t, \mathbf{x} \rangle). \end{aligned}$$

Uzawa's method (Uzawa 58; Bertsekas 99)

ラグランジアンを \mathbf{x}, \mathbf{z} に関して最小化:

$$\begin{cases} \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z}} (f(\mathbf{z}) + \langle \boldsymbol{\alpha}^t, \mathbf{z} \rangle), \\ \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x}} (g(\mathbf{x}) - \langle \mathbf{A}^\top \boldsymbol{\alpha}^t, \mathbf{x} \rangle). \end{cases}$$

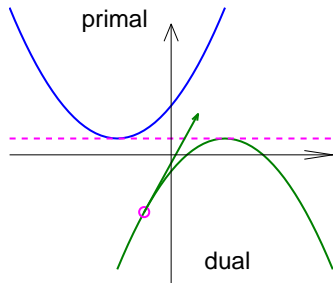
ラグランジュ乗数 $\boldsymbol{\alpha}^t$ を更新:

$$\boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta_t (\mathbf{z}^{t+1} - \mathbf{A}\mathbf{x}^{t+1}).$$

- ラグランジュ乗数の更新は双対での勾配法 (dual ascent) に対応
- **良い点**: シンプル!
- **悪い点**: 凸共役 f^*, g^* が微分できないとき劣勾配法 (sub-gradient ascent) 収束するの?



宇澤 弘文



Uzawa's method を用いた行列穴埋め (Cai et al. 08)

$$\begin{aligned} & \underset{\mathbf{X}}{\text{minimize}} && \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{\text{Strictly convex}} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{\text{Strictly convex}}, \\ & \text{subject to} && \Omega(\mathbf{X}) = \mathbf{z}. \end{aligned}$$

$$\Downarrow$$

ラグランジアン:

$$L(\mathbf{X}, \mathbf{z}, \alpha) = \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{=f(\mathbf{z})} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{=g(\mathbf{x})} + \alpha^\top (\mathbf{z} - \Omega(\mathbf{X})).$$

Uzawa's method:

$$\begin{cases} \mathbf{X}^{t+1} = \text{prox}_\tau (\Omega^\top(\alpha^t)) & (\text{Singular-Value Thresholding}) \\ \mathbf{z}^{t+1} = \mathbf{y} - \lambda \alpha^t \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \Omega(\mathbf{X}^{t+1})) \end{cases}$$

Uzawa's method を用いた行列穴埋め (Cai et al. 08)

$$\begin{aligned} & \underset{\mathbf{X}}{\text{minimize}} && \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{\text{Strictly convex}} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{\text{Strictly convex}}, \\ & \text{subject to} && \Omega(\mathbf{X}) = \mathbf{z}. \end{aligned}$$

$$\Downarrow$$

ラグランジアン:

$$L(\mathbf{X}, \mathbf{z}, \alpha) = \underbrace{\frac{1}{2\lambda} \|\mathbf{z} - \mathbf{y}\|^2}_{=f(\mathbf{z})} + \underbrace{\left(\tau \|\mathbf{X}\|_{\text{tr}} + \frac{1}{2} \|\mathbf{X}\|^2 \right)}_{=g(\mathbf{x})} + \alpha^\top (\mathbf{z} - \Omega(\mathbf{X})).$$

Uzawa's method:

$$\begin{cases} \mathbf{X}^{t+1} = \text{prox}_\tau (\Omega^\top(\alpha^t)) & (\text{Singular-Value Thresholding}) \\ \mathbf{z}^{t+1} = \mathbf{y} - \lambda \alpha^t \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \Omega(\mathbf{X}^{t+1})) \end{cases}$$

拡張ラグランジュ法 (Augmented Lagrangian Method)

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{z}) + g(\mathbf{x}) + \alpha^{\top}(\mathbf{z} - \mathbf{Ax}) + \frac{\eta}{2}\|\mathbf{z} - \mathbf{Ax}\|^2.$$

拡張ラグランジュ法:

$$\left\{ \begin{array}{l} \text{拡張ラグランジアンを } \mathbf{x}, \mathbf{z} \text{ に関して最小化:} \\ (\mathbf{x}^{t+1}, \mathbf{z}^{t+1}) = \underset{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m}{\operatorname{argmin}} L_{\eta_t}(\mathbf{x}, \mathbf{z}, \alpha^t). \\ \\ \text{ラグランジュ乗数を更新:} \\ \alpha^{t+1} = \alpha^t + \eta_t(\mathbf{z}^{t+1} - \mathbf{Ax}^{t+1}). \end{array} \right.$$

- **良い点:** ペナルティ項を追加 (双対は必ず微分可能になる)
- **悪い点:** \mathbf{x} と \mathbf{z} の間に絡みが発生! (別々に最小化できない)

実は

- 拡張ラグランジュ法をまじめにやると双対側で Proximal Minimization

$$\alpha^{t+1} = \operatorname{argmin}_{\alpha \in \mathbb{R}^m} \left(\underbrace{f^*(-\alpha) + g^*(\mathbf{A}^\top \alpha)}_{\text{双対目的関数の符号反転}} + \frac{1}{2\eta_t} \|\alpha - \alpha^t\|^2 \right)$$

をやっているのと等価 .

- DAL (Dual Augmented Lagrangian) はその逆 .
- 拡張ラグランジュを不真面目に解くと双対側で色々な手法が出てくる .

対応関係

	主問題	双対
厳密	拡張ラグランジアン	Proximal Minimization
	Proximal Minimization	DAL
近似	Alternating Minimization Algorithm (Tseng 91)	Forward-Backward Splitting (IST と等価)
	Alternating Direction Method of Multipliers (Gabay & Mercier 76)	Douglas-Rachford Splitting (Lions & Mercier 76)

Alternating Direction Method of Multipliers (ADMM; Gabay & Mercier 76)

$$\left\{ \begin{array}{l} \text{拡張ラグランジアンを } \mathbf{x} \text{ に関して最小化:} \\ \mathbf{x}^{t+1} = \underset{\mathbf{x} \in \mathbb{R}^n}{\operatorname{argmin}} L_{\eta_t}(\mathbf{x}, \mathbf{z}^t, \boldsymbol{\alpha}^t). \\ \text{拡張ラグランジアンを } \mathbf{z} \text{ に関して最小化:} \\ \mathbf{z}^{t+1} = \underset{\mathbf{z} \in \mathbb{R}^m}{\operatorname{argmin}} L_{\eta_t}(\mathbf{x}^{t+1}, \mathbf{z}, \boldsymbol{\alpha}^t). \\ \text{ラグランジュ乗数を更新:} \\ \boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta_t(\mathbf{z}^{t+1} - \mathbf{A}\mathbf{x}^{t+1}). \end{array} \right.$$

- 今更新した \mathbf{x}^{t+1} が \mathbf{z}^{t+1} の計算に入っているところがポイント。

ADMM (Gabay & Mercier 76)

書き直すと

$$\begin{cases} \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left(g(\mathbf{x}) + \frac{\eta_t}{2} \|\mathbf{z}^t - \mathbf{A}\mathbf{x} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right) \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} \left(f(\mathbf{z}) + \frac{\eta_t}{2} \|\mathbf{z} - \mathbf{A}\mathbf{x}^{t+1} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right) \\ \boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta_t (\mathbf{z}^{t+1} - \mathbf{A}\mathbf{x}^{t+1}). \end{cases}$$

- \mathbf{z} に関する最小化は Prox 作用素 prox_f (簡単) .
- \mathbf{x} に関する最小化は行列 \mathbf{A} が変数を絡ませるのでちょっと難しい .
- 1 反復あたりのコストが同じなら FBS より明らかに速い (理論的には不明)
- 双対側での Douglas Rachford Splitting と等価 \Rightarrow **ステップサイズ η によらず** ADMM は安定 . (Lions & Mercier 76; Eckstein & Bertsekas 92)

ADMM (Gabay & Mercier 76)

書き直すと

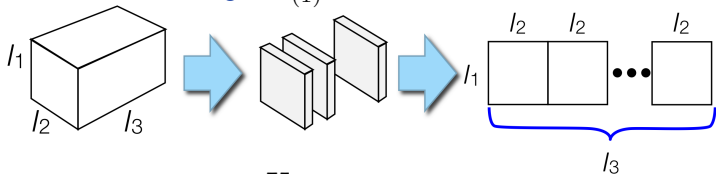
$$\begin{cases} \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left(g(\mathbf{x}) + \frac{\eta_t}{2} \|\mathbf{z}^t - \mathbf{Ax} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right) \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} \left(f(\mathbf{z}) + \frac{\eta_t}{2} \|\mathbf{z} - \mathbf{Ax}^{t+1} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right) \\ \boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta_t (\mathbf{z}^{t+1} - \mathbf{Ax}^{t+1}). \end{cases}$$

- \mathbf{z} に関する最小化は Prox 作用素 prox_f (簡単) .
- \mathbf{x} に関する最小化は行列 \mathbf{A} が変数を絡ませるのでちょっと難しい .
- 1 反復あたりのコストが同じなら FBS より明らかに速い (理論的には不明)
- 双対側での Douglas Rachford Splitting と等価 \Rightarrow **ステップサイズ η によらず** ADMM は安定 . (Lions & Mercier 76; Eckstein & Bertsekas 92)

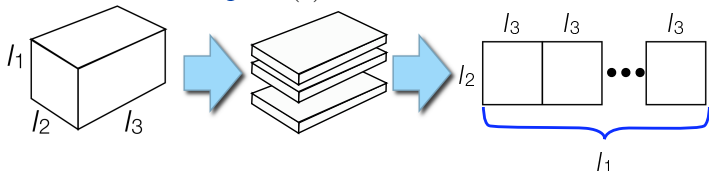
テンソルの穴埋め問題への ADMM の適用

- 凸最適化の適用のポイント: テンソルの行列化 (Matricization)
- テンソルが Tucker 分解の意味で低ランク
 \Leftrightarrow そのテンソルの行列化は (行列の意味で) 低ランク

Mode-1 unfolding $\mathbf{X}_{(1)}$



Mode-2 unfolding $\mathbf{X}_{(2)}$



テンソルの穴埋め問題への ADMM の適用

数学的な定式化:

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K \in \mathbb{R}^N}{\text{minimize}} && \frac{1}{2\lambda} \|\Omega \mathbf{x} - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \underbrace{\|\mathbf{z}_k\|_{\text{tr}}}_{\text{低ランク化}}, \\ & \text{subject to} && \mathbf{P}_k \mathbf{x} = \mathbf{z}_k \quad (k = 1, \dots, K), \end{aligned}$$

- \mathbf{x} は推定すべきテンソルをベクトルとして書いたもの。
- $\mathbf{y} \in \mathbb{R}^M$ は観測 ($M \ll N = n_1 n_2 \cdots n_K$)
- ベクトル化, 行列化は要素の並び替え (線形変換) に過ぎない。
- $\mathbf{P}_k^\top \mathbf{P}_k = \mathbf{I}$ (置換は直交変換)。
- すべてのモードが同時に低ランクになるように正則化。

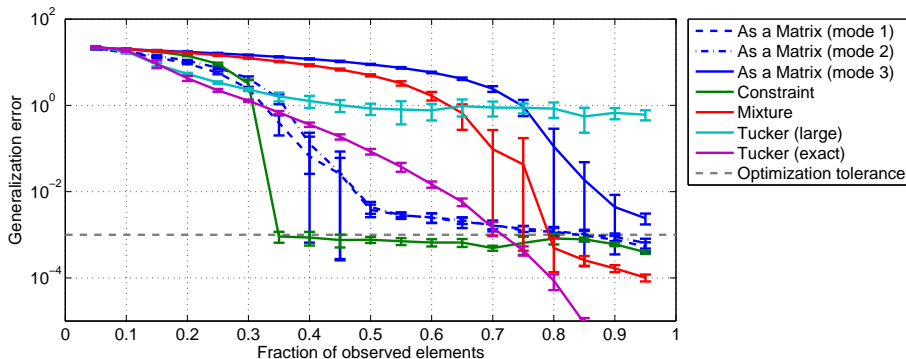
テンソルの穴埋め問題への ADMM の適用

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \{\mathbf{Z}_k\}_{k=1}^K, \{\boldsymbol{\alpha}_k\}_{k=1}^K) = \frac{1}{2\lambda} \|\boldsymbol{\Omega}\mathbf{x} - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_{\text{tr}} \\ + \sum_{k=1}^K \left(\boldsymbol{\alpha}_k^{\top} (\mathbf{P}_k \mathbf{x} - \mathbf{z}_k) + \frac{\eta}{2} \|\mathbf{P}_k \mathbf{x} - \mathbf{z}_k\|^2 \right).$$

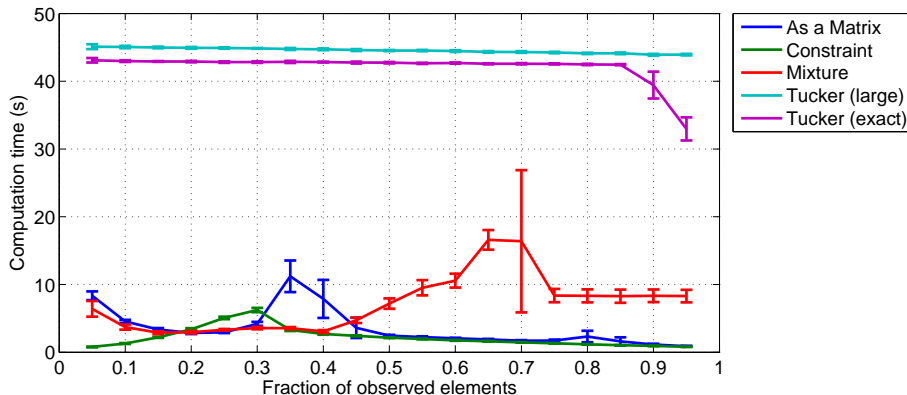
- \mathbf{x} に関する最小化 \mathbf{P}_k が直交行列なので解析的に $O(N)$ で計算可能 .
- \mathbf{Z}_k (\mathbf{z}_k を行列として並べたもの) に関する最小化はトレースノルムに関する Prox 作用素 .
- ラグランジュ乗数ベクトルは制約の数 (モードの数) だけ必要 .

テンソル結果 1: 予測精度



- 提案手法 Constraint は 35% くらい見ればほぼ完璧に予測可能．
ランクを前もって決める必要なし．
- 既存手法 Tucker(EM アルゴリズム) はランクが合っていれば OK．
ランクが間違っていると汎化誤差が収束しない．

テンソル結果 2: 計算速度



- しかも凸最適化は速い！ (Tomioka et al. 10)

Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting**
- 6 まとめ
- 7 Appendix

凸最適化 \Leftrightarrow 方程式のゼロ点

最小化問題

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x}).$$

方程式のゼロ点

Find \mathbf{x} such that

$$(T_f + T_g)(\mathbf{x}) \ni \mathbf{0},$$

where

$$T_f = \partial f,$$

$$T_g = \partial g \quad (\text{劣微分作用素}).$$

例 (Prox 作用素):

$$\mathbf{x} = \text{prox}_f(\mathbf{z}) = \underset{\mathbf{x}' \in \mathbb{R}^n}{\text{argmin}} \left(f(\mathbf{x}') + \frac{1}{2} \|\mathbf{x}' - \mathbf{z}\|^2 \right)$$

$$\Leftrightarrow T_f(\mathbf{x}) + (\mathbf{x} - \mathbf{z}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{prox}_f(\mathbf{z}) = (I + T_f)^{-1}(\mathbf{z})$$

凸最適化 \Leftrightarrow 方程式のゼロ点

最小化問題

$$\min_{\mathbf{x}} f(\mathbf{x}) + g(\mathbf{x}).$$

方程式のゼロ点

Find \mathbf{x} such that

$$(T_f + T_g)(\mathbf{x}) \ni \mathbf{0},$$

where

$$T_f = \partial f,$$

$$T_g = \partial g \quad (\text{劣微分作用素}).$$

例 (Prox 作用素):

$$\mathbf{x} = \text{prox}_f(\mathbf{z}) = \underset{\mathbf{x}' \in \mathbb{R}^n}{\text{argmin}} \left(f(\mathbf{x}') + \frac{1}{2} \|\mathbf{x}' - \mathbf{z}\|^2 \right)$$

$$\Leftrightarrow T_f(\mathbf{x}) + (\mathbf{x} - \mathbf{z}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{prox}_f(\mathbf{z}) = (I + T_f)^{-1}(\mathbf{z})$$

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:

$$\Downarrow$$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:

$$\Downarrow$$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:

$$\Downarrow$$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:



$$\mathbf{x}^{t+1} = \text{prox}_{\eta_t f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:



$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:

$$\Downarrow$$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Forward-Backward Splitting (\Leftrightarrow IST) (Lions & Mercier 76)

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad T_f(\mathbf{x}) + T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \eta T_f(\mathbf{x}) + \eta T_g(\mathbf{x}) \ni \mathbf{0}$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} + \eta T_f(\mathbf{x}) \ni \mathbf{x} - \eta T_g(\mathbf{x})$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad (I + \eta T_f)(\mathbf{x}) = (\mathbf{x} - \eta T_g(\mathbf{x}))$$

$$\Leftrightarrow \text{find } \mathbf{x} \quad \mathbf{x} = \text{prox}_{\eta f}(\mathbf{x} - \eta T_g(\mathbf{x}))$$

反復式:

$$\Downarrow$$

$$\mathbf{x}^{t+1} = \text{prox}_{\eta f}(\mathbf{x}^t - \eta_t \nabla g(\mathbf{x}^t))$$

- Douglas Rachford Splitting も同様の方法で導出できる (複雑すぎるので省略)

Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting
- 6 まとめ**
- 7 Appendix

メッセージ

- ブラックボックス最適化から中身を考慮した最適化へ
- 微分不可能でも怖くない
- Prox 作用素や凸共役を計算しよう

Outline

- 1 イントロ
- 2 準備
- 3 手法 1: Prox 作用素
- 4 手法 2: Legendre 変換
- 5 手法 3: Operator Splitting
- 6 まとめ
- 7 Appendix**

凸共役の性質 (contd.)

- 凸共役 (Legendre 変換) と畳み込み:

$$(f + g)^*(\mathbf{y}) = \inf_{\alpha} (f^*(\mathbf{y} - \alpha) + g^*(\alpha)).$$

- 最小化と凸共役:

$$\inf_{\mathbf{x}} f(\mathbf{x}) = -\sup_{\mathbf{x}} (\langle \mathbf{0}, \mathbf{x} \rangle - f(\mathbf{x})) = -f^*(\mathbf{0}).$$

- Fenchel 双対:

$$\inf_{\mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = -(f + g)^*(\mathbf{0}) = \sup_{\alpha} (-f^*(-\alpha) - g^*(\alpha)).$$

凸共役の性質 (contd.)

- 凸共役 (Legendre 変換) と畳み込み:

$$(f + g)^*(\mathbf{y}) = \inf_{\alpha} (f^*(\mathbf{y} - \alpha) + g^*(\alpha)).$$

- 最小化と凸共役:

$$\inf_{\mathbf{x}} f(\mathbf{x}) = -\sup_{\mathbf{x}} (\langle \mathbf{0}, \mathbf{x} \rangle - f(\mathbf{x})) = -f^*(\mathbf{0}).$$

- Fenchel 双対:

$$\inf_{\mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = -(f + g)^*(\mathbf{0}) = \sup_{\alpha} (-f^*(-\alpha) - g^*(\alpha)).$$

凸共役の性質 (contd.)

- 凸共役 (Legendre 変換) と畳み込み:

$$(f + g)^*(\mathbf{y}) = \inf_{\alpha} (f^*(\mathbf{y} - \alpha) + g^*(\alpha)).$$

- 最小化と凸共役:

$$\inf_{\mathbf{x}} f(\mathbf{x}) = -\sup_{\mathbf{x}} (\langle \mathbf{0}, \mathbf{x} \rangle - f(\mathbf{x})) = -f^*(\mathbf{0}).$$

- Fenchel 双対:

$$\inf_{\mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = -(f + g)^*(\mathbf{0}) = \sup_{\alpha} (-f^*(-\alpha) - g^*(\alpha)).$$

Fenchel 双対定理

$$\inf_{\mathbf{x}} (f(\mathbf{Ax}) + g(\mathbf{x})) = \sup_{\alpha} \left(-f^*(-\alpha) - g^*(\mathbf{A}^T \alpha) \right).$$

- 等式制約付き最適化問題

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^n, \mathbf{z} \in \mathbb{R}^m}{\text{minimize}} && f(\mathbf{z}) + g(\mathbf{x}), \\ & \text{subject to} && \mathbf{Ax} = \mathbf{z} \end{aligned}$$

の双対問題として導出できる .

(<http://www.ibis.t.u-tokyo.ac.jp/RyotaTomioka/Notes/DerivingDual>)

- 凸共役の一覧表があれば機械的に双対問題を作れる点が美味しい .

Fenchel 双対の例 (1): SVM

主問題

$$\min_{\mathbf{w}} f(\mathbf{y} \circ \mathbf{X}\mathbf{w}) + \phi_{\lambda}(\mathbf{x})$$

$$\begin{cases} f(\mathbf{z}) = \sum_{i=1}^m \max(0, 1 - z_i), \\ \phi_{\lambda}(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x}\|^2. \end{cases}$$

双対問題

$$\max_{\alpha} -f^*(-\alpha) - \phi_{\lambda}^*(\mathbf{X}^{\top}(\alpha \circ \mathbf{y}))$$

$$\begin{cases} f^*(-\alpha) = \begin{cases} \sum_{i=1}^m -\alpha_i & (0 \leq \alpha \leq 1), \\ +\infty & (\text{otherwise}), \end{cases} \\ \phi_{\lambda}^*(\mathbf{v}) = \frac{1}{2\lambda} \|\mathbf{v}\|^2. \end{cases}$$

Fenchel 双対の例 (2): 正則化付きロジスティック回帰 (Jaakkola & Haussler 99)

主問題

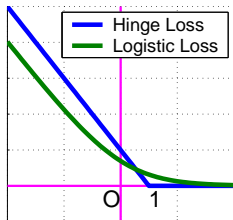
$$\min_{\mathbf{w}} f(\mathbf{y} \circ \mathbf{X}\mathbf{w}) + \phi_{\lambda}(\mathbf{x})$$

$$\begin{cases} f(\mathbf{z}) = \sum_{i=1}^m \log(1 + \exp(-z_i)), \\ \phi_{\lambda}(\mathbf{x}) = \frac{\lambda}{2} \|\mathbf{x}\|^2. \end{cases}$$

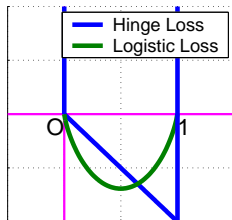
双対問題

$$\max_{\alpha} -f^*(-\alpha) - \phi_{\lambda}^*(\mathbf{X}^{\top}(\alpha \circ \mathbf{y}))$$

$$\begin{cases} f^*(-\alpha) = \sum_{i=1}^m \alpha_i \log(\alpha_i) \\ \quad + (1 - \alpha_i) \log(1 - \alpha_i), \\ \phi_{\lambda}^*(\mathbf{v}) = \frac{1}{2\lambda} \|\mathbf{v}\|^2, \end{cases}$$



(a) primal losses



(b) dual losses

Fenchel 双対の例 (3): 線形計画

主問題

$$(P) \quad \min \quad \mathbf{c}^\top \mathbf{x},$$

$$\text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq 0.$$



双対問題

$$(D) \quad \max \quad \mathbf{b}^\top \mathbf{y},$$

$$\text{s.t.} \quad \mathbf{A}^\top \mathbf{y} \leq \mathbf{c}.$$



主問題'

$$(P') \quad \min_{\mathbf{x}} \quad f(\mathbf{A}\mathbf{x}) + g(\mathbf{x})$$

$$\begin{cases} f(\mathbf{z}) = \begin{cases} 0 & (\mathbf{z} = \mathbf{b}), \\ +\infty & (\text{otherwise}), \end{cases} \\ g(\mathbf{x}) = \begin{cases} \mathbf{c}^\top \mathbf{x} & (\mathbf{x} \geq 0), \\ +\infty & (\text{otherwise}). \end{cases} \end{cases}$$

双対問題'

$$(D') \quad \max_{\mathbf{y}} \quad -f(-\mathbf{y}) - g(\mathbf{A}^\top \mathbf{y})$$

$$\begin{cases} f^*(\mathbf{y}) = \mathbf{b}^\top \mathbf{y}, \\ g^*(\boldsymbol{\alpha}) = \begin{cases} 0 & (\boldsymbol{\alpha} \leq \mathbf{c}), \\ +\infty & (\text{otherwise}). \end{cases} \end{cases}$$

凸じゃない最適化は？

- 一般的に対処するのは困難
- 基本的な方針は凸最適化を繰り返し解くこと
 - EM
 - Difference of Convex (DC) programming
 - CCCP

文献

全体的なまとめ

- Tomioka, Suzuki, & Sugiyama (2011) Augmented Lagrangian Methods for Learning, Selecting, and Combining Features. In Sra, Nowozin, Wright., editors, *Optimization for Machine Learning*, MIT Press.
- Combettes & Pesquet (2010) Proximal splitting methods in signal processing. In *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Springer-Verlag.
- Boyd, Parikh, Peleato, & Eckstein (2010) Distributed optimization and statistical learning via the alternating direction method of multipliers.

教科書

- Rockafellar (1970) *Convex Analysis*. Princeton University Press.
- Bertsekas (1999) *Nonlinear Programming*. Athena Scientific.
- Nesterov (2003) *Introductory Lectures on Convex Optimization: A Basic Course*. Springer.

文献

Proximal Point Algorithm/Augmented Lagrangian

- Arrow, Hurwicz, & Uzawa (1958) Studies in Linear and Non-Linear Programming. Stanford University Press.
- Moreau (1965) Proximité et dualité dans un espace Hilbertien. Bul letin de la S. M. F.
- Rockafellar (1976) Monotone operators and the proximal point algorithm. SIAM J Control Optim 14, 877–898.
- Bertsekas (1982) Constrained Optimization and Lagrange Multiplier Methods. Academic Press.
- Tomioka, Suzuki, & Sugiyama (2011) Super-Linear Convergence of Dual Augmented Lagrangian Algorithm for Sparse Learning. Arxiv:0911.4046.

IST/FISTA

- Nesterov (2007) Gradient Methods for Minimizing Composite Objective Function.
- Beck & Teboulle (2009) A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems. SIAM J Imag Sci 2, 183–202.

文献

Operator Splitting

- Lions & Mercier (1976) Splitting Algorithms for the Sum of Two Nonlinear Operators. SIAM J Numer Anal 16, 964–979.
- Gabay & Mercier (1976) A dual algorithm for the solution of nonlinear variational problems via finite element approximation. Comput Math Appl 2, 17–40.
- Eckstein & Bertsekas (1992) On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators.

マルチタスク / マルチカーネル

- Evgeniou, Micchelli, & Pontil (2005) Learning Multiple Tasks with Kernel Methods. JMLR 6, 615–637.
- Bach, Thibaux, & Jordan (2005) Computing regularization paths for learning multiple kernels. Advances in NIPS, 73–80.
- Suzuki & Tomioka (2009) SpicyMKL. Arxiv:0909.5026.

文献

行列 / テンソル

- Srebro, Rennie, & Jaakkola (2005) Maximum-Margin Matrix Factorization. *Advances in NIPS 17*, 1329–1336.
- Cai, Candès, & Shen (2008) A singular value thresholding algorithm for matrix completion.
- Tomioka, Suzuki, Sugiyama, & Kashima (2010) A Fast Augmented Lagrangian Algorithm for Learning Low-Rank Matrices. In *ICML 2010*.
- Tomioka, Hayashi, & Kashima (2010) On the extension of trace norm to tensors. *ArXiv:1010.0789*.
- Mazumder, Hastie, & Tibshirani (2010) Spectral Regularization Algorithms for Learning Large Incomplete Matrices. *JMLR 11*, 2287–2322.