

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: 鹿島久嗣¹、杉山将²、鈴木大慈¹、林浩平³

¹ 東京大学 ² 東京工業大学 ³ 奈良先端科学技術大学院大学

2011-10-25 @ RAMP 2011

最適化業界-機械学習業界の間の需給ミスマッチ

- 最適化業界

- ▶ 最適化のことよく分からなくても使えるツールボックスが必要
- ▶ ワンストップサービス — CVX (Grant & Boyd)
- ▶ 連続最適化なら内点法 (80年代~)

- 機械学習業界

- ▶ モデルが変わってもすぐ実装を変更できる方がよい.
- ▶ なるべく簡単な手法が好ましい.
- ▶ 並列化できるとなおよい.

⇒ 古い手法 (60-70年代) がどうやら熱い.

最適化業界-機械学習業界の間の需給ミスマッチ

● 最適化業界

- ▶ 最適化のことよく分からなくても使えるツールボックスが必要
- ▶ ワンストップサービス — CVX (Grant & Boyd)
- ▶ 連続最適化なら内点法 (80年代~)

● 機械学習業界

- ▶ モデルが変わってもすぐ実装を変更できる方がよい.
- ▶ なるべく簡単な手法が好ましい.
- ▶ 並列化できるとなおよい.

⇒ 古い手法 (60-70年代) がどうやら熱い.

- ▶ (Accelerated) Proximal gradient methods

最適化業界-機械学習業界の間の需給ミスマッチ

● 最適化業界

- ▶ 最適化のことよく分からなくても使えるツールボックスが必要
- ▶ ワンストップサービス — CVX (Grant & Boyd)
- ▶ 連続最適化なら内点法 (80年代~)

● 機械学習業界

- ▶ モデルが変わってもすぐ実装を変更できる方がよい.
- ▶ なるべく簡単な手法が好ましい.
- ▶ 並列化できるとなおよい.

⇒ 古い手法 (60-70年代) がどうやら熱い.

- ▶ (Accelerated) Proximal gradient methods
- ▶ Dual decomposition (Uzawa's method)

最適化業界-機械学習業界の間の需給ミスマッチ

● 最適化業界

- ▶ 最適化のことよく分からなくても使えるツールボックスが必要
- ▶ ワンストップサービス — CVX (Grant & Boyd)
- ▶ 連続最適化なら内点法 (80年代~)

● 機械学習業界

- ▶ モデルが変わってもすぐ実装を変更できる方がよい.
- ▶ なるべく簡単な手法が好ましい.
- ▶ 並列化できるとなおよい.

⇒ 古い手法 (60-70年代) がどうやら熱い.

- ▶ (Accelerated) Proximal gradient methods
- ▶ Dual decomposition (Uzawa's method)
- ▶ Alternating Direction Method of Multipliers (ADMM)

機械学習における連続最適化の古いトレンド?

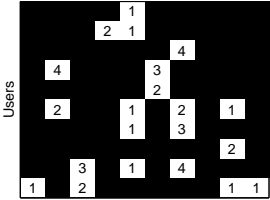
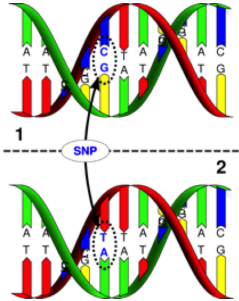
なぜこれらの手法がいま注目されるか — スパース推定

- 高次元データ (サンプル数 \ll 次元)

- ▶ バイオインフォマティクス (遺伝子発現, SNP 解析, etc)
- ▶ テキストマイニング (系列ラベリング, 係り受け解析)
- ▶ イメージング (MRI) — 圧縮センシング

- 構造があるデータ

- ▶ 協調フィルタリング — **低ランク構造**
- ▶ グラフィカルモデル推定 — **グラフ構造**



例 1: SNP (一塩基多型) 解析

x_i : 入力 (SNP), $y_i = 1$: 病気, $y_i = -1$: 健康

目的: ゲノムの個人差 x_i と病気になるかならないか y_i の関係を知りたい.

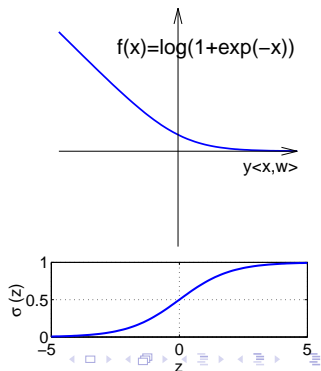
ロジスティック回帰: 2 値分類規則の学習法 ($y_i \in \{-1, +1\}$)

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \underbrace{\sum_{i=1}^m \log(1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle))}_{\text{data-fit}} + \underbrace{\lambda \|\mathbf{w}\|_1}_{\text{Regularization}}$$

- 例えば, SNP の数 $n = 500,000$, 被験者の数 $m = 5,000$
- 事後確率最大化 (MAP) 法の一つ.
ロジスティック損失関数:

$$\log(1 + e^{-yz}) = -\log P(Y = y|z)$$

$$\text{where } P(Y = +1|z) = \frac{e^z}{1+e^z}.$$



例 2: 圧縮センシング [Candes, Romberg, & Tao 06]

低次元 (ノイズ入り) 観測からの信号 (MRI 画像) 復元

$$\underset{\mathbf{w} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{\Omega} \mathbf{w}\|_2^2 + \lambda \|\mathbf{\Phi} \mathbf{w}\|_1$$

- \mathbf{y} : ノイズ入り観測信号
- \mathbf{w} : 原信号
- $\mathbf{\Omega}: \mathbb{R}^n \rightarrow \mathbb{R}^m$: 観測行列 (ランダム, フーリエ変換)
- $\mathbf{\Phi}$: 原信号がスパースとなる基底への変換行列
- $\mathbf{\Phi}^{-1}$ が存在すれば, より簡単な問題

$$\underset{\tilde{\mathbf{w}} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{A} \tilde{\mathbf{w}}\|_2^2 + \lambda \|\tilde{\mathbf{w}}\|_1,$$

(ただし $\mathbf{A} = \mathbf{\Omega} \mathbf{\Phi}^{-1}$) を解けばよい.

例 3: 低ランク行列の推定 [Fazel+ 01; Srebro+ 05]

行列 \mathbf{X} を部分的な (ノイズ入り) 観測 \mathbf{Y} から復元したい:

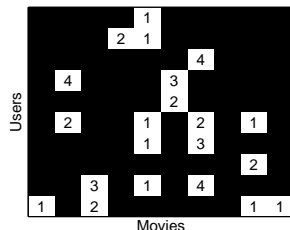
$$\underset{\mathbf{X}}{\text{minimize}} \quad \frac{1}{2} \|\Omega(\mathbf{X} - \mathbf{Y})\|^2 + \lambda \|\mathbf{X}\|_{S_1}$$

where $\|\mathbf{X}\|_{S_1} := \sum_{j=1}^r \sigma_j(\mathbf{X})$ (Schatten 1-norm)

特異値の線形和

⇒ 特異値の意味でスパース

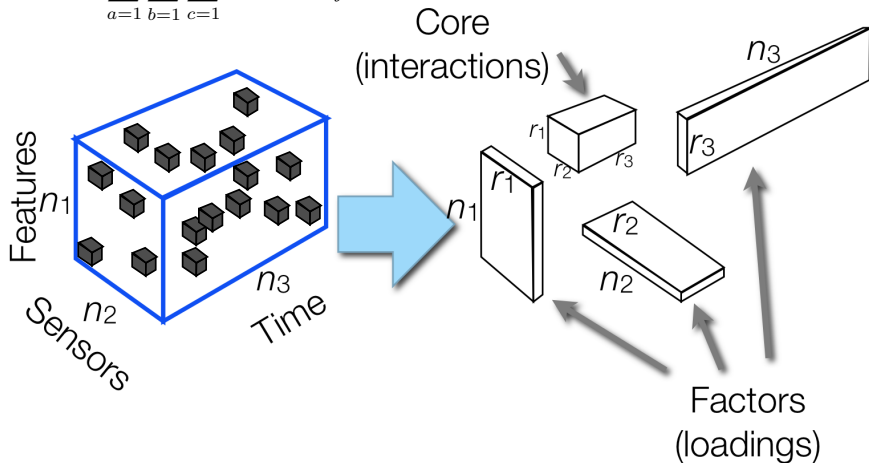
⇒ 低ランク



例 4: 低ランクテンソルの補完 [Tucker 66]

$$X_{ijk} = \sum_{a=1}^{r_1} \sum_{b=1}^{r_2} \sum_{c=1}^{r_3} C_{abc} U_{ia}^{(1)} U_{jb}^{(2)} U_{kc}^{(c)}$$

Tucker decomposition



単純スパース推定問題と構造付きスパース推定問題

● 単純スパース推定問題

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

- ▶ SNP 解析
- ▶ 圧縮センシングで Φ^{-1} が存在する場合 (ウェーブレット)
- ▶ 協調フィルタリング (行列穴埋め)

● 構造付きスパース推定問題

$$\underset{\mathbf{w}}{\text{minimize}} \quad L(\mathbf{w}) + \lambda \|\Phi \mathbf{w}\|_1$$

- ▶ 圧縮センシングで Φ^{-1} が存在しない場合 (Total variation)
- ▶ テンソルの Tucker 分解

今日の内容

- 単純スパース推定問題のための最適化手法
 - ▶ (加速付き) 近接勾配法 (proximal gradient method)
 - ▶ Dual Augmented Lagrangian (DAL)
- 構造付きスパース推定問題のための最適化手法
 - ▶ Alternating Direction Method of Multipliers (ADMM)

単純スパース推定問題のための最適化手法

- (加速付き) 近接勾配法 (proximal gradient method)
- Dual Augmented Lagrangian (DAL)

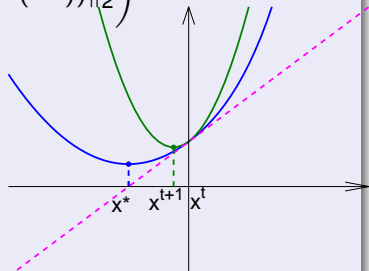
近接勾配法 (proximal gradient method)

最小化問題

$$\underset{\mathbf{w}}{\text{minimize}} \quad \underbrace{L(\mathbf{w})}_{\text{微分可能}} + \underbrace{\lambda \|\mathbf{w}\|_1}_{\text{微分不可能}}$$

線形化/最小化

$$\begin{aligned} \mathbf{w}^{t+1} &= \underset{\mathbf{w}}{\text{argmin}} \left(\nabla L(\mathbf{w}^t)(\mathbf{w} - \mathbf{w}^t) + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|_2^2 + \lambda \|\mathbf{w}\|_1 \right) \\ &= \underset{\mathbf{w}}{\text{argmin}} \left(\lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - (\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t))\|_2^2 \right) \\ &= \text{prox}_{\lambda\eta_t}(\mathbf{w}^t - \eta_t \nabla L(\mathbf{w}^t)). \end{aligned}$$

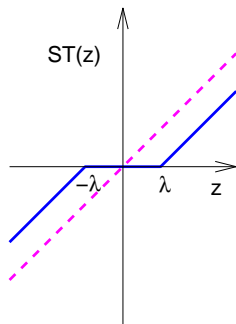


Proximal operator: 射影の一般化

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

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

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



- 何らかの意味で分離可能な関数 r 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 .

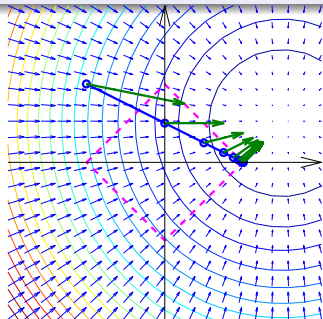
近接勾配法 (proximal gradient method)

近接勾配法 (Lions & Mercier 79; Figueiredo&Nowak 03; Daubechies 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 :

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

- 利点: 実装が簡単 .
- 欠点: 損失項 L のヘシアン の条件数が悪いと遅い .
- 別名: Forward-Backward splitting, Iterative Shrinkage/Thresholding



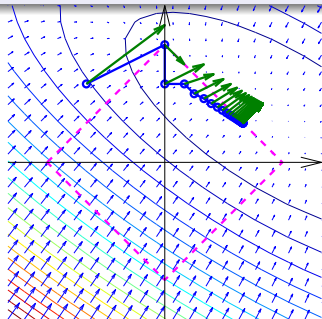
近接勾配法 (proximal gradient method)

近接勾配法 (Lions & Mercier 79; Figueiredo&Nowak 03; Daubechies 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 :

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

- 利点: 実装が簡単 .
- 欠点: 損失項 L のヘシアン の条件数が悪いと遅い .
- 別名: Forward-Backward splitting, Iterative Shrinkage/Thresholding



近接勾配法の収束レートと加速

- 損失項 L が強凸, かつリプシッツ定数

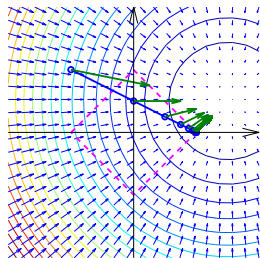
$$\|\nabla L(\mathbf{x}) - \nabla L(\mathbf{y})\| \leq H\|\mathbf{x} - \mathbf{y}\|$$

が存在すれば **1 次収束** (勾配法と同じ)

- 強凸でない場合, ステップサイズ $\eta_t \leq 1/H$ とすることで, 多項式レート

$$f(\mathbf{x}^k) - f(\mathbf{x}^*) \leq \frac{H\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2k}$$

- 下限 $O(1/k^2)$ を達成するための加速法も提案されている . (Nesterov 07; Beck & Teboulle 09)



Dual Augmented Lagrangian (DAL) [Tomioka & Sugiyama 09]

- 1次ブラックボックスモデルでは下限が達成されている。
- もう少し機械学習における問題の構造を考慮したい。
 - ① 損失項は $L(\mathbf{w}) = f_\ell(\mathbf{A}\mathbf{w})$ と分解できる。 f_ℓ : ロス関数, $\mathbf{A} \in \mathbb{R}^{m \times n}$: データ行列

$$\text{(例 1)} \quad L(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{w}\|_2^2 \quad (\text{2乗ロス回帰})$$

$$\text{(例 2)} \quad L(\mathbf{w}) = \sum_{i=1}^m \log(1 + \exp(-y_i \langle \mathbf{x}_i, \mathbf{w} \rangle)) \quad (\text{ロジスティック回帰})$$

- ② スパースな解に興味があるので、解がスパースであるほど効率的な解法が望ましい (データ行列は必ずしもスパースではない)

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

主問題

$$\min_{\mathbf{w}} \underbrace{f_{\ell}(\mathbf{A}\mathbf{w}) + \lambda \|\mathbf{w}\|_1}_{f(\mathbf{w})}$$

双対問題

$$\begin{aligned} \max_{\boldsymbol{\alpha}, \mathbf{v}} \quad & -f_{\ell}^*(-\boldsymbol{\alpha}) - (\lambda \|\cdot\|_1)^*(\mathbf{v}) \\ \text{s.t.} \quad & \mathbf{v} = \mathbf{A}^{\top} \boldsymbol{\alpha} \end{aligned}$$

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

主問題

$$\min_{\mathbf{w}} \underbrace{f_{\ell}(\mathbf{A}\mathbf{w}) + \lambda \|\mathbf{w}\|_1}_{f(\mathbf{w})}$$

Proximal minimization

[Rockafellar 76]:

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

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

- 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

$$\begin{aligned} \max_{\boldsymbol{\alpha}, \mathbf{v}} \quad & -f_{\ell}^*(-\boldsymbol{\alpha}) - (\lambda \|\cdot\|_1)^*(\mathbf{v}) \\ \text{s.t.} \quad & \mathbf{v} = \mathbf{A}^{\top} \boldsymbol{\alpha} \end{aligned}$$

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

主問題

$$\min_{\mathbf{w}} \underbrace{f_{\ell}(\mathbf{A}\mathbf{w}) + \lambda \|\mathbf{w}\|_1}_{f(\mathbf{w})}$$

Proximal minimization

[Rockafellar 76]:

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

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

- 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

$$\begin{aligned} \max_{\alpha, \mathbf{v}} \quad & -f_{\ell}^*(-\alpha) - (\lambda \|\cdot\|_1)^*(\mathbf{v}) \\ \text{s.t.} \quad & \mathbf{v} = \mathbf{A}^{\top} \alpha \end{aligned}$$

\Leftrightarrow Augmented Lagrangian

[Powell 69; Hestenes 69]:

$$\begin{aligned} \mathbf{w}^{t+1} &= \operatorname{prox}_{\lambda\eta_t}(\mathbf{w}^t + \eta_t \mathbf{A}^{\top} \alpha^t) \\ \alpha^t &= \operatorname{argmin}_{\alpha} \varphi_t(\alpha) \end{aligned}$$

- $\varphi_t(\alpha)$ の最小化は簡単 (なめらか).
- ステップサイズ η_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}^0 を決める .
- ② 停止条件が満たされるまで反復 :

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

ただし ,

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

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

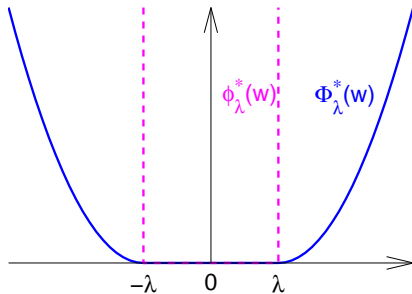
$$\mathbf{w}^{t+1} = \text{prox}_{\eta_t \lambda} \left(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}^t \right)$$

(2) 内部最適化は微分可能

$$\boldsymbol{\alpha}^t = \underset{\boldsymbol{\alpha}}{\text{argmin}} \left(\underbrace{f_\ell^*(-\boldsymbol{\alpha})}_{\text{微分可能. } \mathbf{A} \text{ のスケーリングの影響を受けない}} + \frac{1}{2\eta_t} \underbrace{\|\text{prox}_{\lambda\eta_t}(\mathbf{w}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha})\|^2}_{\text{非ゼロ成分の数に比例}} \right)$$

微分可能. \mathbf{A} の
スケーリングの
影響を受けない

非ゼロ成分の数に比例



近接勾配法と DAL の違い : いかに変数の間の絡みを除くか

目的関数 f に関する Proximation は難しい :

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{f(\mathbf{w})}_{f_l(\mathbf{A}\mathbf{w})} + \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

変数が絡みあっている

近接勾配法とDALの違い：いかに変数の間の絡みを除くか

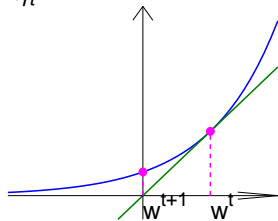
目的関数 f に関する Proximation は難しい：

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{変数が絡みあっている}} + \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- 近接勾配法（既存）：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



近接勾配法とDALの違い：いかに変数の間の絡みを除くか

目的関数 f に関する Proximation は難しい：

$$\mathbf{w}^{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\underbrace{f_{\ell}(\mathbf{A}\mathbf{w})}_{\text{変数が絡みあっている}} + \lambda \|\mathbf{w}\|_1 + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}^t\|^2 \right)$$

- 近接勾配法（既存）：線形にロス項を 近似：

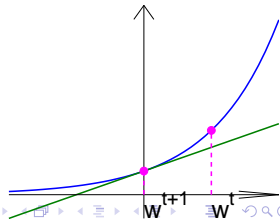
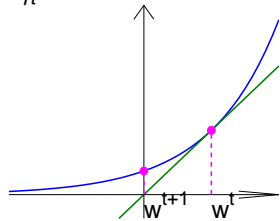
$$f_{\ell}(\mathbf{A}\mathbf{w}) \simeq f_{\ell}(\mathbf{A}\mathbf{w}^t) + (\mathbf{w} - \mathbf{w}^t)^{\top} \mathbf{A}^{\top} \nabla f_{\ell}(\mathbf{A}\mathbf{w}^t)$$

→ 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

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

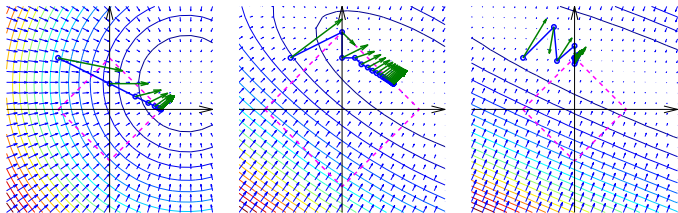
→ 次の点 \mathbf{w}^{t+1} で最もタイト



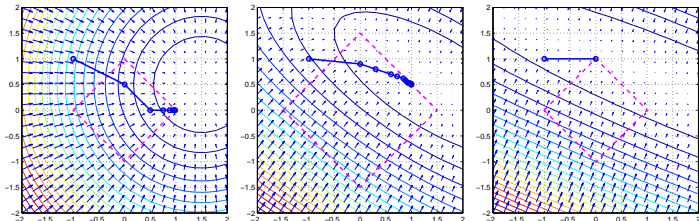
数値例

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．

近接勾配法



DAL



定理 1 (厳密な最小化)

定義

- \mathbf{w}^t : 厳密な DAL 法 ($\|\nabla\varphi_t(\boldsymbol{\alpha}^t)\| = 0$) 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 σ 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

- \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 .

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

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 .

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

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

⇒ η_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 $\frac{\|\nabla\varphi_t(\alpha^t)\|}{\|\mathbf{w}^{t+1} - \mathbf{w}^t\|} \leq O(1/\eta_t)$.

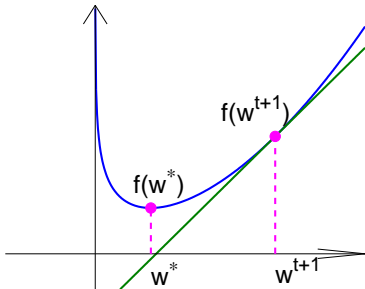
定理 1 の証明 (エッセンス)

\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),

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

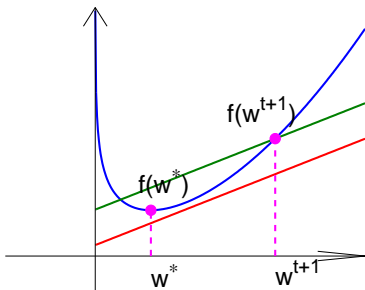


定理 2 の証明 (エッセンス)

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

近似最小化のコスト

$1/\gamma$: 損失関数の微分 ∇f_ℓ のリプシッツ定数 .



構造付きスパース推定問題のための最適化手法

- Alternating Direction Method of Multipliers (ADMM)

拡張ラグランジュ法 [Powell 69; Hestenes 69]

最小化問題

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{z}}{\text{minimize}} && f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1, \\ & \text{s.t.} && \mathbf{z} = \Phi \mathbf{x} \end{aligned}$$

拡張ラグランジアン

$$L_\eta(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

拡張ラグランジュ法

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

x と z の間に絡みが発生! (別々に最小化できない)

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

& Mercier 76)

拡張ラグランジアン

$$L_\eta(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \alpha^\top (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

拡張ラグランジアンを \mathbf{x} に関して最小化:

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} L_{\eta_t}(\mathbf{x}, \mathbf{z}^t, \alpha^t).$$

拡張ラグランジアンを \mathbf{z} に関して最小化:

$$\mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} L_{\eta_t}(\mathbf{x}^{t+1}, \mathbf{z}, \alpha^t).$$

ラグランジュ乗数を更新:

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

- 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)

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \alpha^{\top} (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

書き直すと

$$\begin{cases} \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} L_{\eta_t}(\mathbf{x}, \mathbf{z}^t, \alpha^t). \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} L_{\eta_t}(\mathbf{x}^{t+1}, \mathbf{z}, \alpha^t). \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \Phi \mathbf{x}^{t+1}). \end{cases}$$

ADMM (Gabay & Mercier 76)

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \alpha^{\top} (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

書き直すと

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

ADMM (Gabay & Mercier 76)

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \mathbf{z}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \boldsymbol{\alpha}^{\top} (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

書き直すと

$$\begin{cases} \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left(f(\mathbf{x}) + \frac{\eta_t}{2} \|\mathbf{z}^t - \Phi \mathbf{x} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right). \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} \left(\lambda \|\mathbf{z}\|_1 + \frac{\eta_t}{2} \|\mathbf{z} - \Phi \mathbf{x}^{t+1} + \boldsymbol{\alpha}^t / \eta_t\|^2 \right). \\ \boldsymbol{\alpha}^{t+1} = \boldsymbol{\alpha}^t + \eta_t (\mathbf{z}^{t+1} - \Phi \mathbf{x}^{t+1}). \end{cases}$$

ADMM (Gabay & Mercier 76)

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \mathbf{z}, \alpha) = f(\mathbf{x}) + \lambda \|\mathbf{z}\|_1 + \alpha^{\top} (\mathbf{z} - \Phi \mathbf{x}) + \frac{\eta}{2} \|\mathbf{z} - \Phi \mathbf{x}\|^2.$$

書き直すと

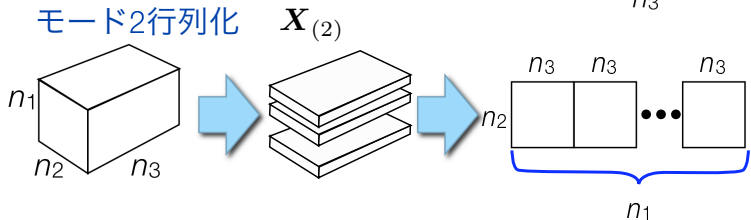
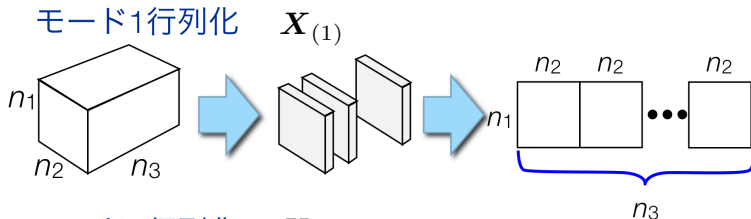
$$\begin{cases} \mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n} \left(f(\mathbf{x}) + \frac{\eta_t}{2} \|\mathbf{z}^t - \Phi \mathbf{x} + \alpha^t / \eta_t\|^2 \right). \\ \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathbb{R}^m} \left(\lambda \|\mathbf{z}\|_1 + \frac{\eta_t}{2} \|\mathbf{z} - \Phi \mathbf{x}^{t+1} + \alpha^t / \eta_t\|^2 \right). \\ \alpha^{t+1} = \alpha^t + \eta_t (\mathbf{z}^{t+1} - \Phi \mathbf{x}^{t+1}). \end{cases}$$

- \mathbf{z} に関する最小化は Prox 作用素 $\operatorname{prox}_{\lambda/\eta_t}$ (簡単) .
- \mathbf{x} に関する最小化は行列 Φ 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近接勾配法より経験的に速い (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凸最適化の適用 [Liu+09,

Signoretto +10, Tomioka+10, Gandy+11]

- 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
⇔ そのテンソルの行列化は (行列の意味で) 低ランク



テンソルの穴埋め問題への ADMM の適用

数学的な定式化:

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_K \in \mathbb{R}^N}{\text{minimize}} && \frac{1}{2\lambda} \|\Omega \mathbf{x} - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \underbrace{\|\mathbf{z}_k\|_{S_1}}_{\text{低ランク化}}, \\ & \text{s.t.} && \mathbf{P}_k \mathbf{x} = \mathbf{z}_k \quad (k = 1, \dots, K), \end{aligned}$$

- \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 ($M \ll N = n_1 n_2 \cdots n_K$)
- \mathbf{P}_k はモード k 行列化の操作を行列で表現したもの .
- $\mathbf{P}_k^\top \mathbf{P}_k = \mathbf{I}$ (行列化は直交変換) .
- すべてのモードが同時に低ランクになるように正則化 .

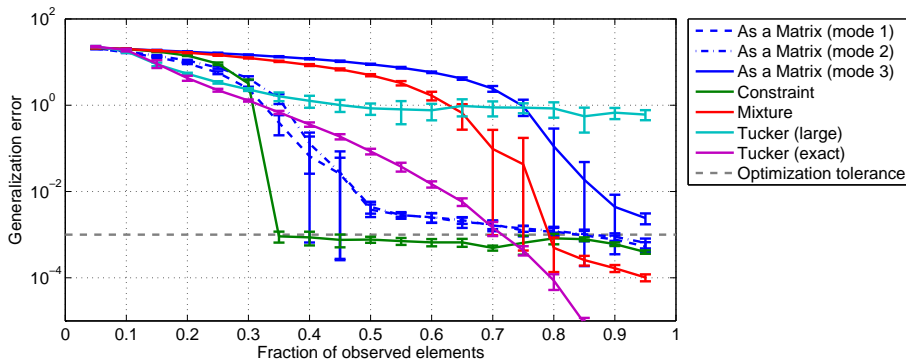
テンソルの穴埋め問題への ADMM の適用

拡張ラグランジアン

$$L_{\eta}(\mathbf{x}, \{\mathbf{Z}_k\}_{k=1}^K, \{\alpha_k\}_{k=1}^K) = \frac{1}{2\lambda} \|\Omega \mathbf{x} - \mathbf{y}\|^2 + \sum_{k=1}^K \gamma_k \|\mathbf{Z}_k\|_{S_1} \\ + \sum_{k=1}^K \left(\alpha_k^{\top} (\mathbf{P}_k \mathbf{x} - \mathbf{z}_k) + \frac{\eta}{2} \|\mathbf{P}_k \mathbf{x} - \mathbf{z}_k\|^2 \right).$$

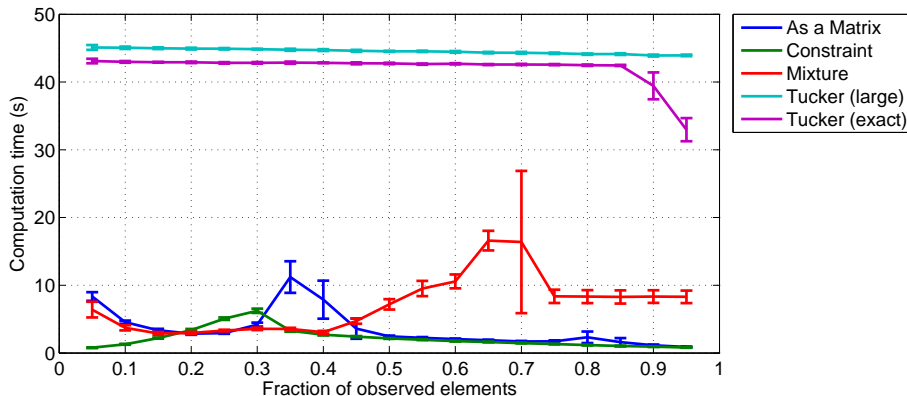
- \mathbf{x} に関する最小化 \mathbf{P}_k が直交行列なので解析的に $O(N)$ で計算可能 .
- \mathbf{Z}_k (\mathbf{z}_k を行列として並べたもの) に関する最小化は Schatten 1-ノルムに関する Prox 作用素 .
- ラグランジュ乗数ベクトルは制約の数 (モードの数) だけ必要 .

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



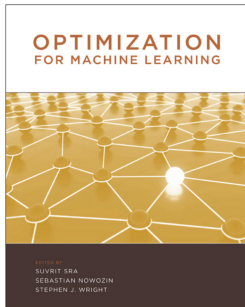
- しかも凸最適化は速い！

- 最適化（凸最適化）：機械学習研究者にとって欠かせないツール
- ブラックボックス最適化から中身を考慮した最適化へ
 - ▶ 単純スパース推定
 - ▶ 構造付きスパース推定
- 理論解析の中でも最適化を含めた話が重要
 - ▶ Stochastic Optimization in Machine Learning (Nathan Srebro, tutorial at ICML 2010)
- 並列化，オンライン化などがホットな話題
 - ▶ LCCC : NIPS 2010 Workshop on Learning on Cores, Clusters and Clouds

ご清聴ありがとうございました！

宣伝

Optimization for Machine Learning (MIT Press, 2011)



謝辞

これらの研究の様々な段階でコメントを頂いた土谷隆先生，小島政和先生，福島雅夫先生に感謝します．この研究は科研費 22700138 および NTT コミュニケーション科学基礎研究所の支援を受けています．

References

Recent surveys

- 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.

IST/FISTA

- Moreau (1965) Proximité et dualité dans un espace Hilbertien. *Bul letin de la S. M. F.*
- 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.

Augmented Lagrangian

- Rockafellar (1976) Augmented Lagrangians and applications of the proximal point algorithm in convex programming. *Math. of Oper. Res.* 1.
- 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. *JMLR* 12.

References

ADMM

- Gabay & Mercier (1976) A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput Math Appl* 2, 17–40.
- Lions & Mercier (1979) Splitting Algorithms for the Sum of Two Nonlinear Operators. *SIAM J Numer Anal* 16, 964–979.
- Eckstein & Bertsekas (1992) On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators.

Matrices/Tensor

- Fazal, Hindi, & Boyd (2001) A Rank Minimization Heuristic with Application to Minimum Order System Approximation. *Proc. of the American Control Conference*.
- 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.
- Mazumder, Hastie, & Tibshirani (2010) Spectral Regularization Algorithms for Learning Large Incomplete Matrices. *JMLR* 11, 2287–2322.
- Tomioka, Hayashi, & Kashima (2011) Estimation of low-rank tensors via convex optimization. *arXiv:1010.0789*.

Total variation

- Rudin, Osher, Fetemi. (1992) Nonlinear total variation based noise removal algorithms. *Physica D: Nonlinear Phenomena*, 60.
- Goldstein & Osher (2009) Split Bregman method for L1 regularization problems. *SIAM J. Imag. Sci.* 2.