

A Fast Augmented Lagrangian Algorithm for Learning Low- Rank Matrices

Ryota Tomioka, Taiji Suzuki, Masashi Sugiyama,
Hisashi Kashima

2010/8/18
ICML2010読む会

お話したいこと

- * 低ランク行列学習はいろいろな応用がある
- * 提案法 (M-DAL) は速くて理論的な性能保証のある実用的なアルゴリズム

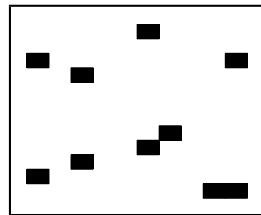
ソフトウェア:

[http://www.ibis.t.u-tokyo.ac.jp/
RyotaTomioaka/Softwares/
SupportPageICML10](http://www.ibis.t.u-tokyo.ac.jp/RyotaTomioaka/Softwares/SupportPageICML10)

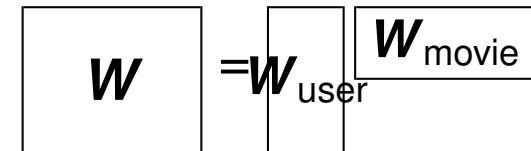
低ランク学習

- 1 Matrix completion [Srebro et al. 05; Abernethy et al. 09] (collaborative filtering, link prediction)

$$Y = W$$



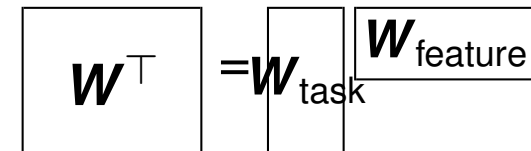
Part of Y is observed



- 2 Multi-task learning [Argyriou et al., 07]

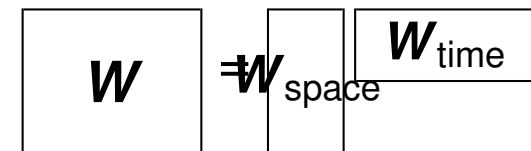
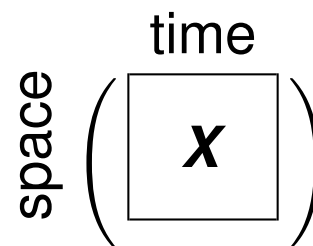
$$y = W^T x + b$$

$$W^T = \begin{bmatrix} w_{\text{task 1}}^T \\ w_{\text{task 2}}^T \\ \vdots \\ w_{\text{task R}}^T \end{bmatrix}$$



- 3 Predicting over matrices (classification/regression) [Tomioka & Aihara, 07]

$$y = \langle W, X \rangle + b$$



低ランク化と特徴抽出

$$f(X) = \langle W, X \rangle + b$$

$$= \text{Tr}(W^\top X) + b$$

$$= \text{Tr}\left(\sum_{j=1}^r \sigma_j \mathbf{v}_j \mathbf{u}_j^\top X\right) + b$$

$$= \sum_{j=1}^r \sigma_j \text{Tr}(\mathbf{v}_j \mathbf{u}_j^\top X) + b$$

$$= \sum_{j=1}^r \sigma_j \text{Tr}(\mathbf{u}_j^\top X \mathbf{v}_j) + b$$

$$\left(W = \sum_{j=1}^r \sigma_j \mathbf{u}_j \mathbf{v}_j^\top \right)$$

特徴抽出
空間フィルタ
時間フィルタ

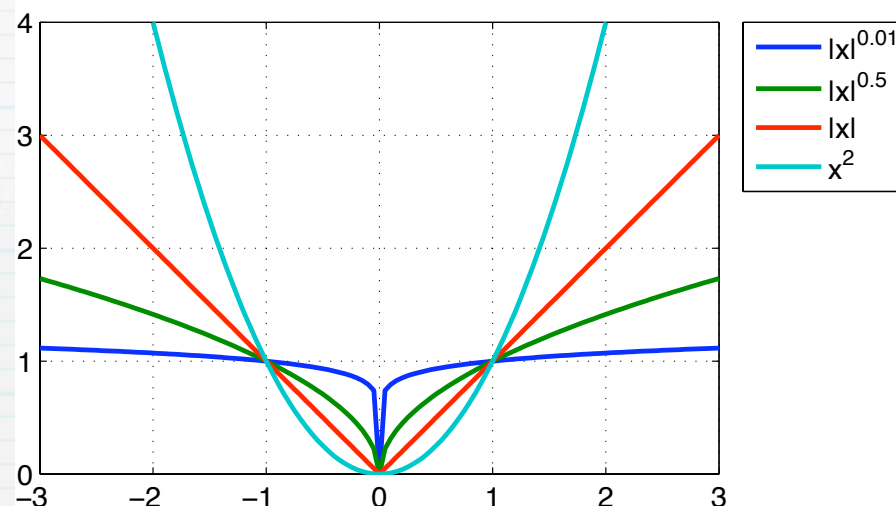
行列のランクの凸近似： トレースノルム

- * L1-ノルム：ベクトルの非ゼロ成分の個数の凸近似

$$\|\mathbf{w}\|_1 = \sum_{j=1}^n |w_j|$$

- * トレースノルム：行列のランク (=非ゼロ特異値の数)の凸近似

$$\|\mathbf{W}\|_* = \sum_{j=1}^r \sigma_j(\mathbf{W})$$



$\sigma_j(\mathbf{W})$ は \mathbf{W} の j 番目に大きい特異値

最適化問題(従来法)

$$\underset{\mathbf{W} \in \mathbb{R}^{R \times C}}{\text{minimize}} \quad L(\mathbf{W}) + \lambda \|\mathbf{W}\|_*$$

損失項(データ) 正則化項

- * Proximal (accelerated) gradient [Ji & Ye, 09]
 - * 途中の解も低ランク
 - * データのスケーリングが悪いと遅い
 - * 収束性は $O(1/k^2)$: First-order black-box の方法としては最適
- * 内点法 [Tomioka & Aihara, 07]
 - * スケーリングにロバスト
 - * 収束性は $O(\exp(-k))$

スケーリングにロバストで途中の解が低ランク
なアルゴリズムが欲しい! → M-DAL(提案法)

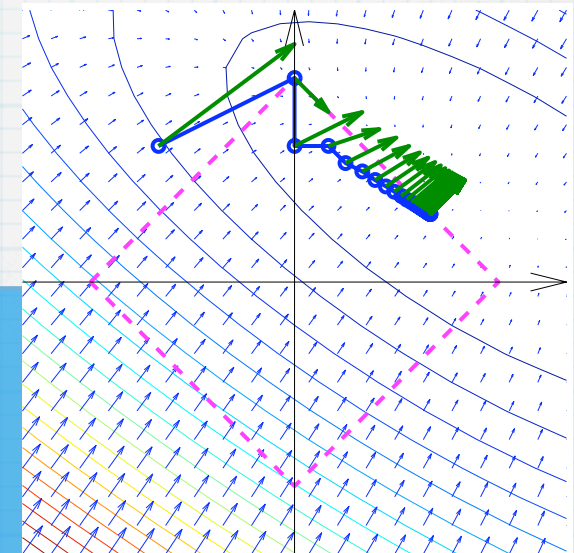
従来法: Proximal Gradient

1. W^0 を適当に初期化. 適当に η_1, η_2, \dots を選ぶ.

2. 停止基準が満たされるまで

$$W^{t+1} := \operatorname{argmin}_W \left(\underbrace{\langle W, \nabla L(W^t) \rangle}_{\text{ロス項の線形近似}} + \underbrace{\lambda \|W\|_*}_{\text{正則化項}} + \underbrace{\frac{1}{2\eta_t} \|W - W^t\|_{\text{fro}}^2}_{\text{Proximity 項}} \right)$$

$$= \underbrace{\text{ST}_{\lambda\eta_t}}_{\text{縮小}} \left(\underbrace{W^t - \eta_t \nabla L(W^t)}_{\text{勾配ステップ}} \right)$$



Spectral Soft-threshold:

$$\text{ST}_\lambda(W) = U \max(S - \lambda I, 0) V^\top$$

ただし、特異値分解

$$W = USV^\top$$

良い点:

- シンプル

悪い点:

- スケールにロバストでない

最適化問題(提案法)

$$\underset{W \in \mathbb{R}^{R \times C}}{\text{minimize}} \quad f_{\ell}(A \text{vec}(W)) + \lambda \|W\|_*$$

損失項 正則化項

データ ↗

- * f_{ℓ} は損失関数(ロジスティックなど)でそれなりに性質がよいとする
- * A はデザイン行列(サンプル数 \times 未知変数の数)で性質が悪いかも
- * f_{ℓ} と A を区別する $\rightarrow A$ (データ) の性質によらない収束性の評価

例: 多入力多出力回帰

$$L(\mathbf{W}) = \sum_{i=1}^m \|\mathbf{y}_i - \mathbf{W} \mathbf{x}_i\|^2$$

$$= \|\mathbf{Y} - \mathbf{W} \mathbf{X}\|_{\text{fro}}^2$$

$$= \left\| \text{vec}(\mathbf{Y}) - \underbrace{(\mathbf{X}^{\top} \otimes \mathbf{I}_R)}_{\text{デザイン行列 } A} \text{vec}(\mathbf{W}) \right\|^2$$

デザイン行列 A

ロス関数

$$\left(\text{vec}(AXB) = (B^{\top} \otimes A) \text{vec}(X) \right)$$

クロネッカー積

* 他のロスでも同様

提案法: Dual Augmented-Lagrangian

1. W^0 を適当に初期化. $\eta_1 < \eta_2 < \dots$ を選ぶ

2. 停止基準が満たされるまで

$$W^{t+1} := \operatorname{argmin}_W \left(\underbrace{f_\ell(\mathbf{A} \operatorname{vec}(\mathbf{W}))}_{\text{ロス項 (線形近似なし)}} + \underbrace{\lambda \|\mathbf{W}\|_*}_{\text{正則化項}} + \underbrace{\frac{1}{2\eta_t} \|\mathbf{W} - \mathbf{W}^t\|^2}_{\text{Proximity 項}} \right)$$
$$= \operatorname{ST}_{\lambda\eta_t} (\mathbf{W}^t + \eta_t \mathbf{A}^\top \boldsymbol{\alpha}^{t+1})$$

ただし $\boldsymbol{\alpha}^{t+1} = \operatorname{argmin} \varphi_t(\boldsymbol{\alpha})$

* $\varphi_t(\boldsymbol{\alpha})$ の最小化は易しい(微分可能)

* ステップサイズ η_t は増加列

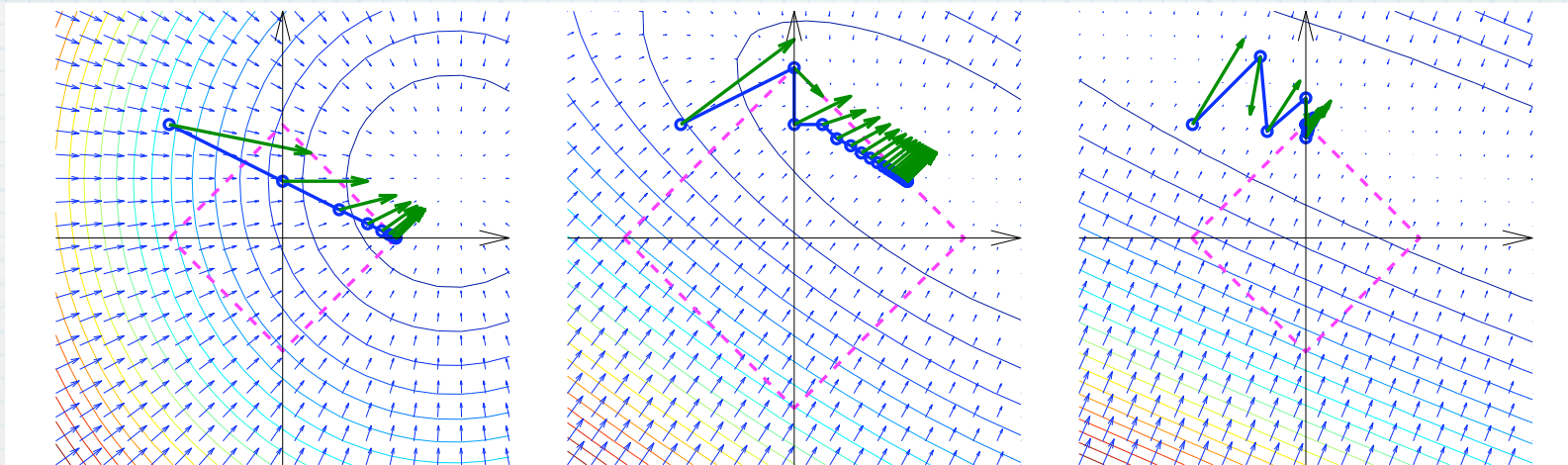
* 導出は Fenchel 双対を使う

* 名前の由来: 双対側での Augmented Lagrangian 法なので [Rockafellar 76]

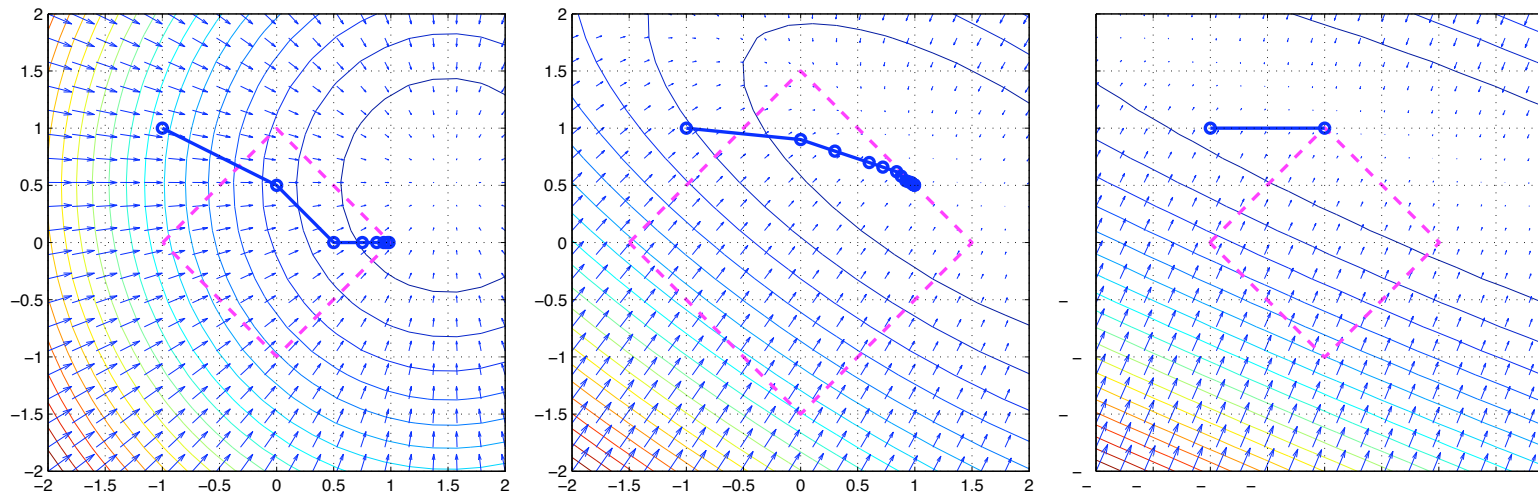
数値例

スケールが悪いほどDALが有利

Proximal Gradient



DAL



定理 (収束レート)

- * W^* : 目的関数 $f(W)$ (ロスと正則化項の和) を最小化する元
- * W^t : DAL アルゴリズムで生成された解. ただし, $\phi_t(\alpha)$ の最小化は以下の条件で止める:

$$\|\nabla\varphi_t(\alpha^{t+1})\| \leq \sqrt{\frac{\gamma}{\eta_t}} \|\mathbf{W}^{t+1} - \mathbf{W}^t\|_{\text{fro}}$$

- * $1/\gamma$ は損失関数の微分 ∇f_{ell} のリプシッツ定数
- * 仮定:

$$f(\mathbf{W}^{t+1}) - f(\mathbf{W}^*) \geq \sigma \|\mathbf{W}^{t+1} - \mathbf{W}^*\|_{\text{fro}}^2 \quad (t = 0, 1, 2, \dots)$$

- * 結果: 超1次収束

$$\|\mathbf{W}^{t+1} - \mathbf{W}^*\|_{\text{fro}} \leq \frac{1}{\sqrt{1 + 2\sigma\eta_t}} \|\mathbf{W}^t - \mathbf{W}^*\|_{\text{fro}}$$

定理の味あいどころ

- * ロス関数の性質 → 定数 γ
 - * 停止基準に影響
- * デザイン行列の性質 → 定数 σ
 - * 停止基準に無関係
 - * 正であれば超1次収束が成立 (小さくてもよい)
- * ステップサイズ η_t
 - * 大きいほど収束は速い
 - * 大きいほど内部最小化は大変

正則化項の一般化

* マルチ行列学習

$$\phi_\lambda(\mathbf{W}) = \lambda \sum_{k=1}^K \|\mathbf{W}^{(k)}\|_* = \lambda \left\| \begin{pmatrix} \mathbf{W}^{(1)} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{W}^{(K)} \end{pmatrix} \right\|_*$$

* 大きい行列を作る必要なし

* 複数の情報源からの情報統合に使える

* 一般のスペクトル正則化

$$\phi_\lambda(\mathbf{W}) = \sum_{j=1}^r g_\lambda(\sigma_j(\mathbf{W}))$$

* g_λ は原点でゼロの凸関数で,以下の操作が計算できる:

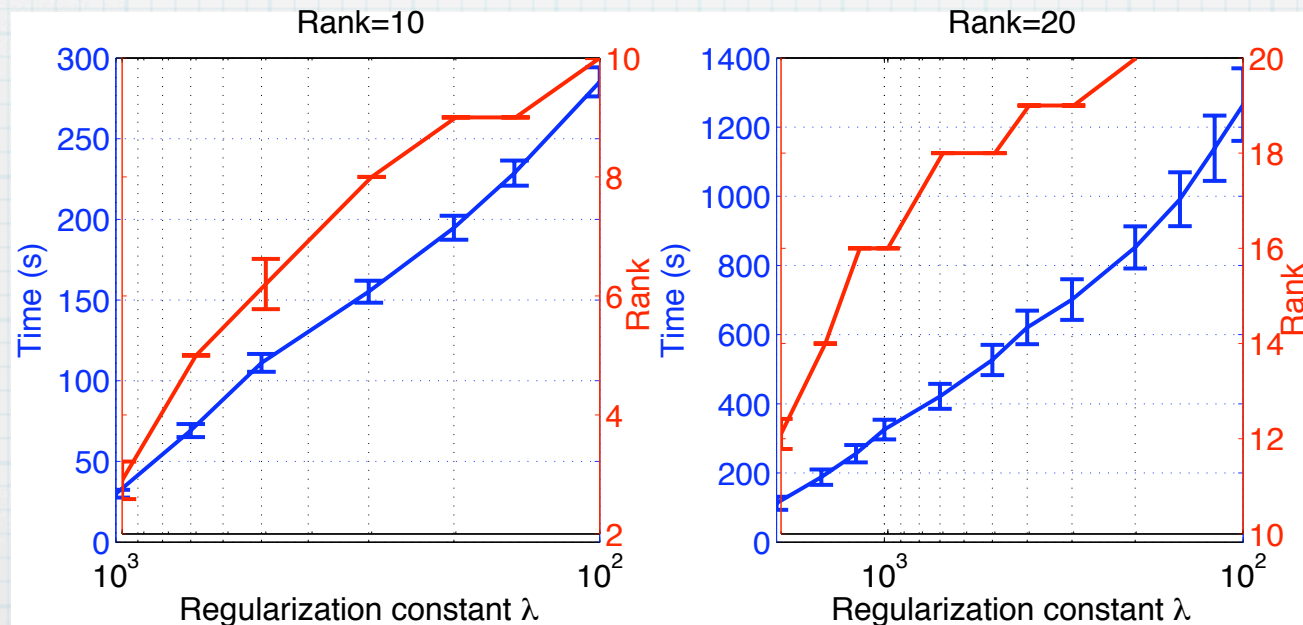
$$\text{ST}_{g_\lambda}(\sigma_j) = \operatorname{argmin}_{x \in \mathbb{R}} \left(g_\lambda(x) + \frac{1}{2}(x - \sigma_j)^2 \right)$$

実験1: 行列穴埋め

- * 特徴: 大規模&構造のあるデザイン行列A
 - * 真の行列 W^* : $10,000 \times 10,000$ (要素数1億) 低ランク
 - * 観測: ランダムに m 個を選ぶ
- ➡ 内部目的関数 $\phi_t(\alpha)$ の最小化に**擬似ニュートン法**を使う
- ➡ $10,000 \times 10,000$ 行列をメモリに保持する必要なし

$m=1,200,000$

$m=2,400,000$



実験1: ランク=10 (詳細)

Rank=10, #observations $m=1,200,000$

λ	time (s)	#outer	#inner	rank	S-RMSE
1000	33.1 (± 2.0)	5 (± 0)	8 (± 0)	2.8 (± 0.4)	0.0158 (± 0.0024)
700	77.1 (± 5.6)	11 (± 0)	18 (± 0)	5 (± 0)	0.0133 (± 0.0008)
500	124 (± 7.2)	17 (± 0)	28 (± 0)	6.4 (± 0.5)	0.0113 (± 0.0015)
300	174 (± 8.0)	23 (± 0)	38.4 (± 0.84)	8 (± 0)	0.00852 (± 0.00039)
200	220 (± 9.9)	29 (± 0)	48.4 (± 0.84)	9 (± 0)	0.00767 (± 0.00031)
150	257 (± 9.9)	35 (± 0)	58.4 (± 0.84)	9 (± 0)	0.00498 (± 0.00026)
100	319 (± 11)	41 (± 0)	70 (± 0.82)	10 (± 0)	0.00743 (± 0.00013)

* 内部反復数はたかだか外部反復数の2倍程度

➡ 内部反復の停止基準は厳密過ぎない

実験1: ランク=20 (詳細)

Rank=20, #observations $m=2,400,000$

λ	time (s)	#outer	#inner	rank	S-RMSE
2000	112 (± 19)	6 (± 0)	15.1 (± 1.0)	12.1 (± 0.3)	0.011 (± 0.002)
1500	188 (± 22)	11 (± 0)	24.1 (± 1.0)	14 (± 0)	0.0094 (± 0.001)
1200	256 (± 25)	15 (± 0)	31.1 (± 1.0)	16 (± 0)	0.0090 (± 0.0008)
1000	326 (± 29)	19 (± 0)	38.1 (± 1.0)	16 (± 0)	0.0073 (± 0.0007)
700	421 (± 36)	24 (± 0)	48.1 (± 1.0)	18 (± 0)	0.0065 (± 0.0004)
500	527 (± 44)	29 (± 0)	57.1 (± 1.0)	18 (± 0)	0.0042 (± 0.0003)
400	621 (± 48)	34 (± 0)	66.1 (± 1.0)	19 (± 0)	0.0044 (± 0.0002)
300	702 (± 59)	38.5 (± 0.5)	74.1 (± 1.5)	19 (± 0)	0.0030 (± 0.0003)
200	852 (± 61)	43.6 (± 0.5)	83.9 (± 2.3)	20 (± 0)	0.0039 (± 0.0001)
150	992 (± 78)	48.4 (± 0.7)	92.5 (± 1.5)	20 (± 0)	0.0024 (± 0.0002)
120	1139 (± 94)	53.4 (± 0.7)	102 (± 1.5)	20 (± 0)	0.0016 ($\pm 6 \times 10^{-5}$)
100	1265 (± 105)	57.7 (± 0.8)	109 (± 2.4)	20 (± 0)	0.0013 ($\pm 8 \times 10^{-5}$)

実験2: 行列の上の教師付き判別

* 真の行列 W^* : 64×64 , $\text{rank}=16$

* サンプル数 $m=1,000$

* ロジスティック損失:

$$f_{\ell}(z) = \sum_{i=1}^m \log(1 + \exp(-y_i z_i))$$

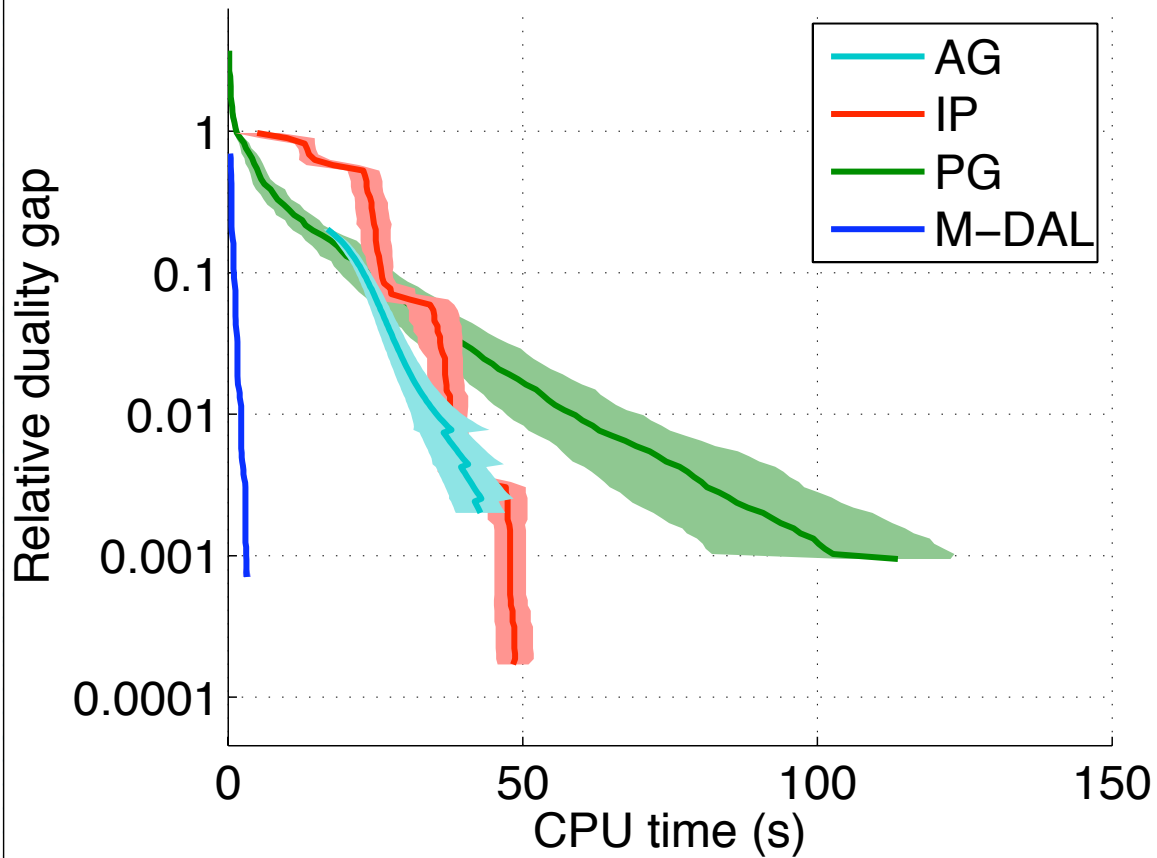
* デザイン行列 A : 1000×64^2 (dense)

* 各訓練データはウィッシュャート分布からサンプル

* 中規模でdense \rightarrow 内部目的関数 $\phi_t(\alpha)$ の最小化に**ニュートン法**を使う

実験2: 結果

双対ギャップ

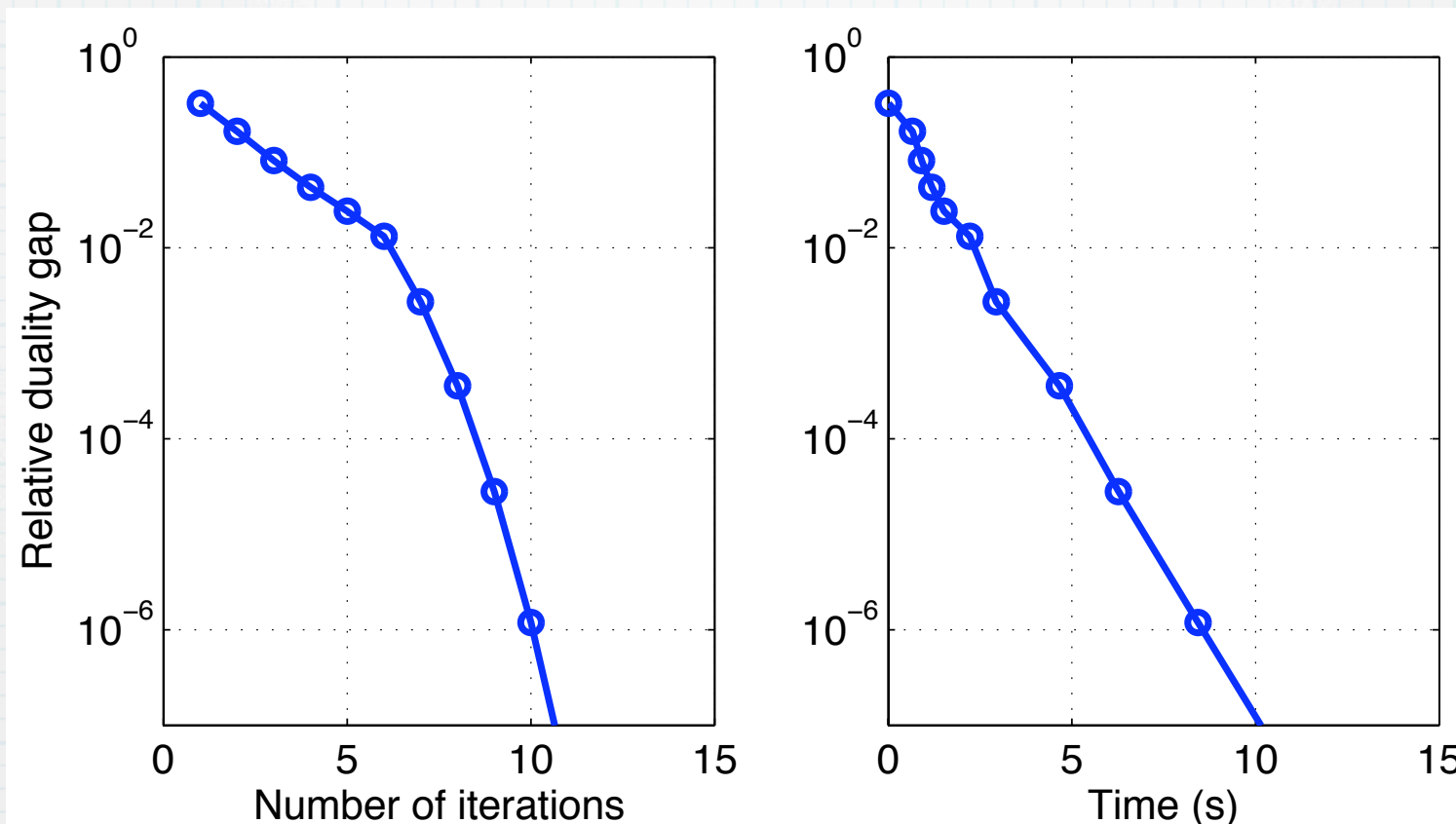


計算時間

- * M-DAL (提案手法)
- * IP (内点法)
- * PG (Projected Gradient)
- * AG (Accelerated Gradient) [Ji & Ye, 09]: Proximal Gradient を速くしたもの

超一次収束

双対ギャップ



反復数

計算時間

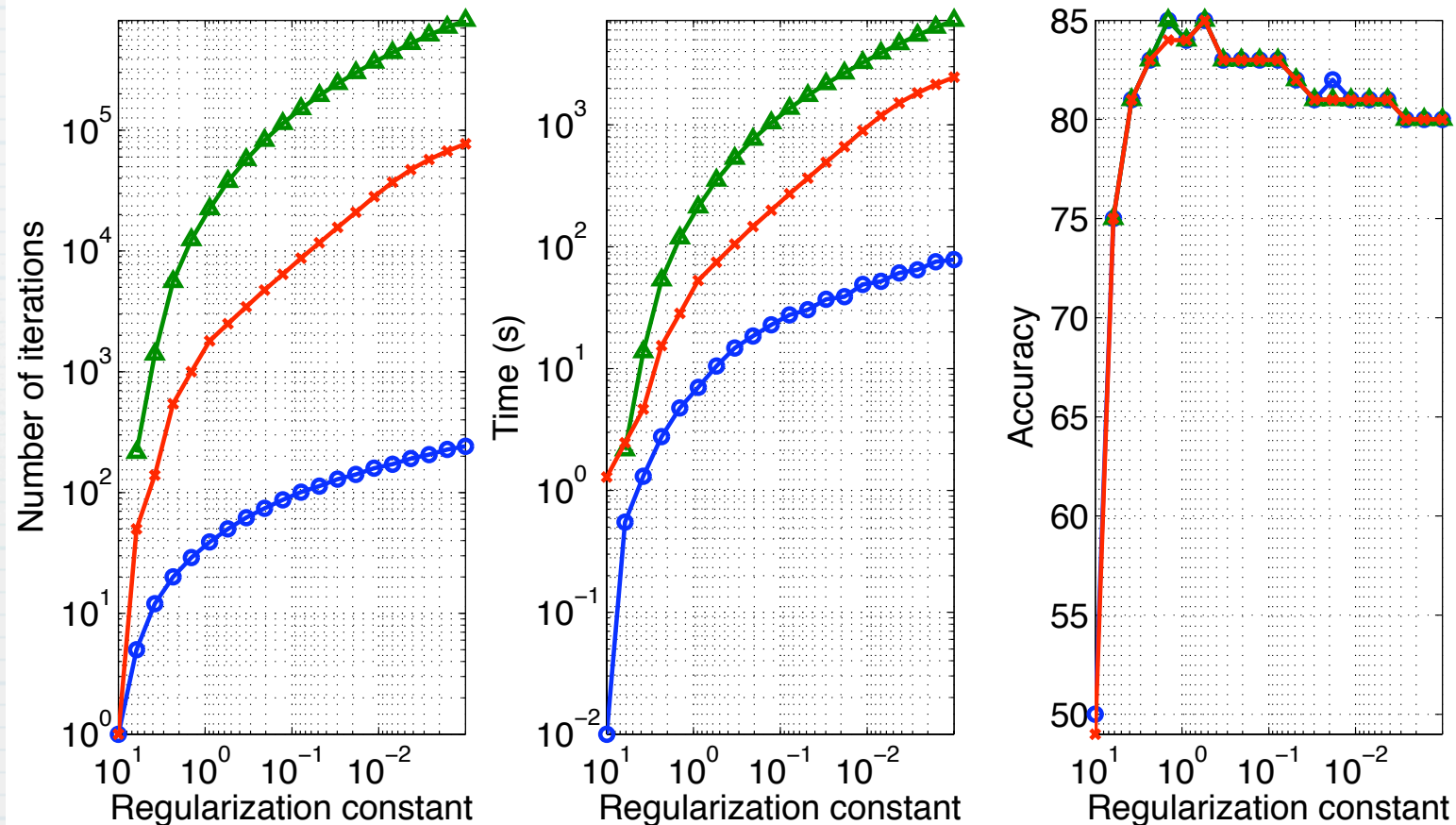
実験3: BCI

- * BCI competition 2003 dataset IV
- * タスク: 次に右 / 左のどちらの手の指でキーを叩くか予測する
- * データ: 多次元脳波データ: 28チャンネルx50時間点 (長さ500ms)
- * それぞれの訓練データを3つの行列に加工
 - * 1次成分 ($<20\text{Hz}$): 28×50 行列
 - * 2次成分 (alpha 7-15Hz): 28×28 行列(共分散)
 - * 2次成分 (beta 15-30Hz): 28×28 行列(共分散)
- * 訓練データ数 316. テストデータ数100.



BCI: 結果

-o- M-DAL (提案法), -△- PG, -x- AG [Ji&Ye 09]

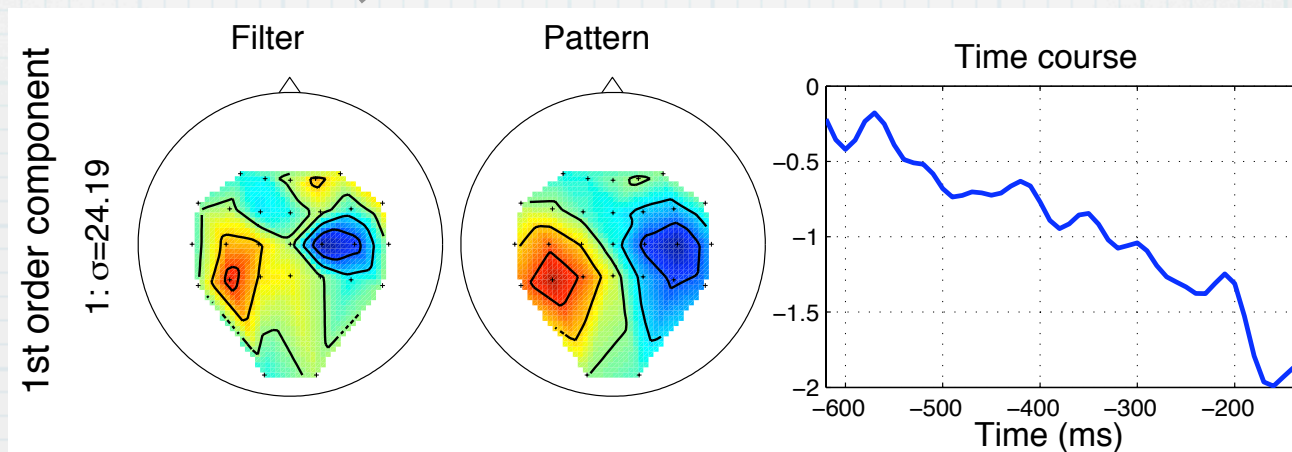
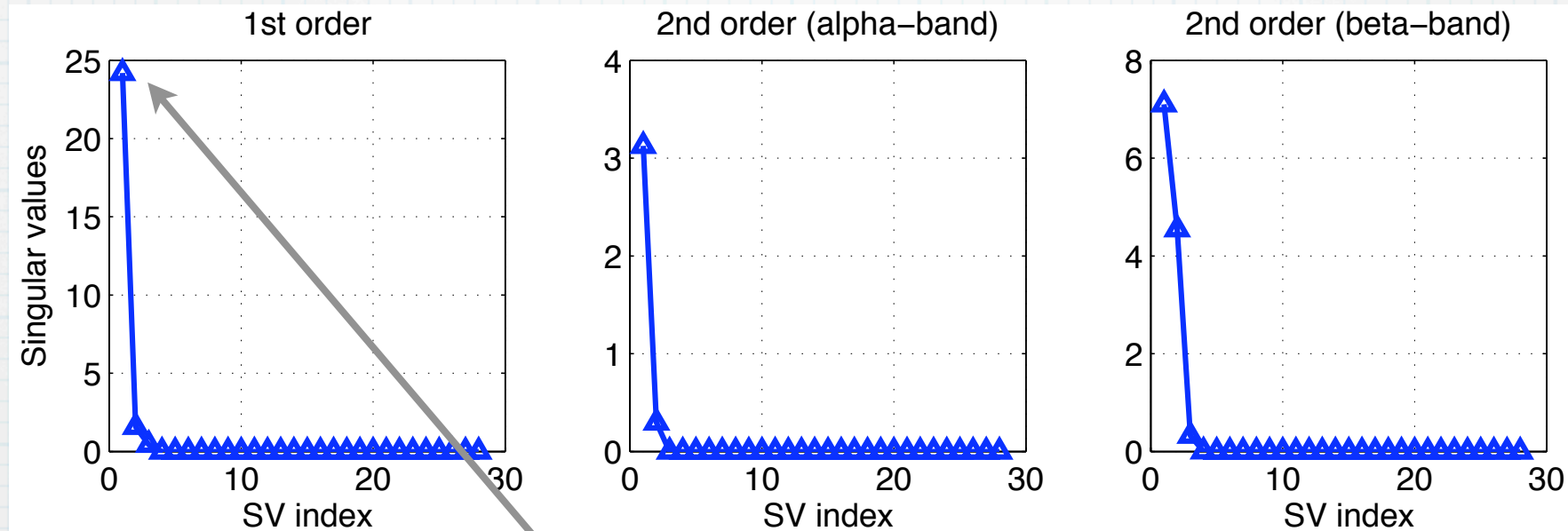


* 外部停止基準: 双対ギャップ $\leq 1e-3$

* 注: 上の計算時間は正則化パスの計算のコスト

BCI: 結果

特異値



まとめ

- * トレースノルム(とその拡張)は低ランク行列の学習を凸最適化で行うための鍵
- * トレースノルムは行列におけるL1正則化に対応
- * DALは途中の解を低ランクで保持することが可能
- * DAL は非漸近的に超1次収束する
 - ➡ 反復数が非常に小さい
- * 特異値に関して分解可能な正則化項+微分可能な損失関数
 - ➡ 1反復あたりの計算量が小さい
- * 低ランク行列の推定/学習だけでなく複数情報源の統合も可能

ICMLのフィードバック

- * Matrix Factorization は流行り気味?
- * セッション1a: Topic Models and Matrix Factorization
- * セッション2a: Matrix Factorization and Recommendation
- * スパース/MKL/最適化はぽつぽつ