

正則化学習法における最適化手法



鈴木 大慈

東京大学

情報理工学系研究科

数理情報学専攻

2013/2/18@九州大学伊都キャンパス

文部科学省委託事業数学協働プログラム

「最適化ワークショップ: 広がっていく最適化」

高次元データ スパース正則化学習法

最適化手法

proximal point algorithm

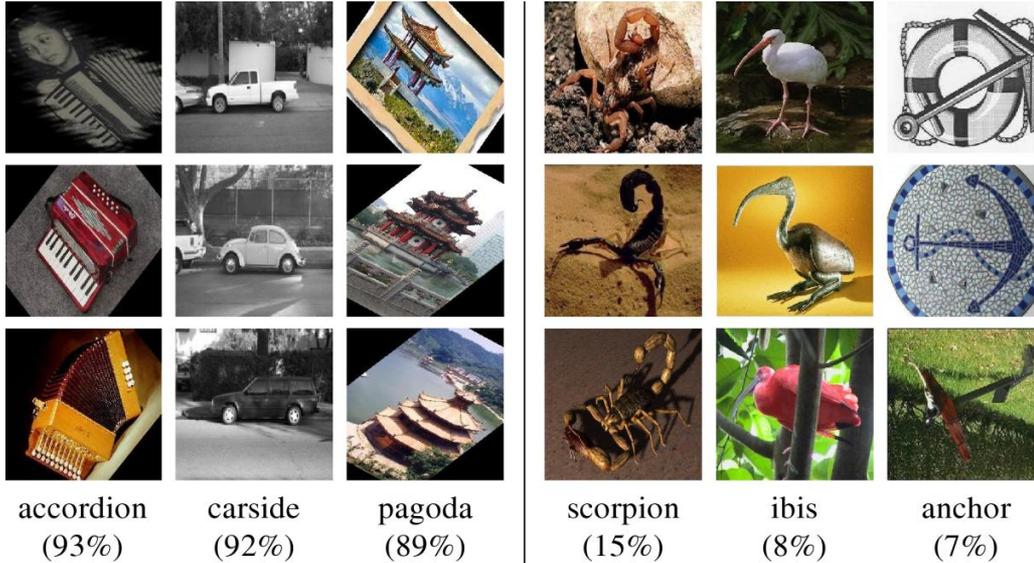
確率最適化手法

問題設定

スパース正則化学習

高次元線形判別

物体認識



音声認識



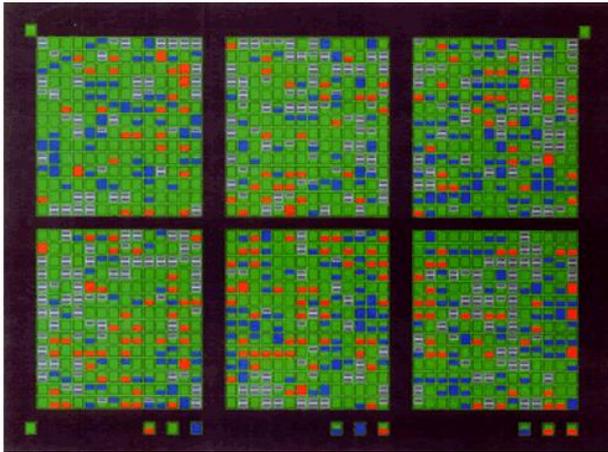
自然言語処理



バイオインフォマティクス

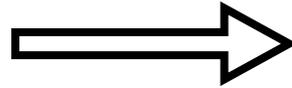


DNAデータ



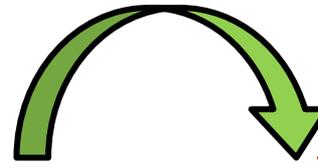
特徴ベクトル

ベクトル化



高次元

判別



癌になりやすい
癌になりにくい

特徴ベクトル

テキストデータ

Syria crisis: 28,000 disappeared, say rights groups

Human rights groups working in Syria say at least 28,000 people have disappeared after being abducted by soldiers or militia.

They say they have the names of 18,000 people missing since anti-government protests began 18 months ago and know of another 10,000 cases.

Online activist group Avaaz says "nobody is safe" from a deliberate government campaign of terror.

It intends to give the UN Human Rights Council a dossier for investigation.

Avaaz has gathered testimony from Syrians who says husbands, sons and daughters have been forcibly abducted by pro-government forces.

Alice Jay, campaign director at Avaaz, said: "Syrians are being plucked off the street by security forces and paramilitaries and being 'disappeared' into torture cells.

"Whether it is women buying groceries or farmers going for fuel, nobody is safe."

She said it was a deliberate strategy to "terrorise families and communities".

"The panic of not knowing whether your husband or child is alive breeds such fear that it silences dissent," she said.

"The fate of each and every one of these people must be investigated and the perpetrators punished."

Fadel Abdulghani, of the Syrian Network for Human Rights, estimates that 28,000 people have disappeared since unrest against the government of President Bashar al-Assad began last year.

Muhannad al-Hasani, of the Syrian human rights organisation Sawasya, put the figure even higher.



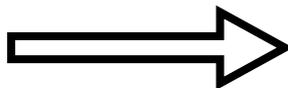
Fighting between government forces and rebels is continuing in the city of Aleppo

Syria conflict

- Turkish town scarred by conflict
- No-man's land
- Turkey-Syria tensions
- Assad heartland

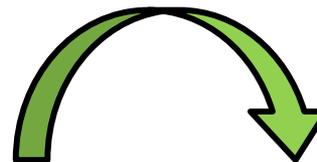
“
The regime is doing this for two reasons - 10

ベクトル化



高次元

判別



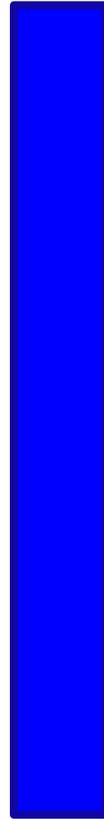
何の話題か？

特徴ベクトル

画像データ

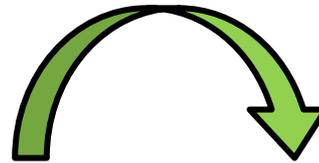


ベクトル化
→



高次元

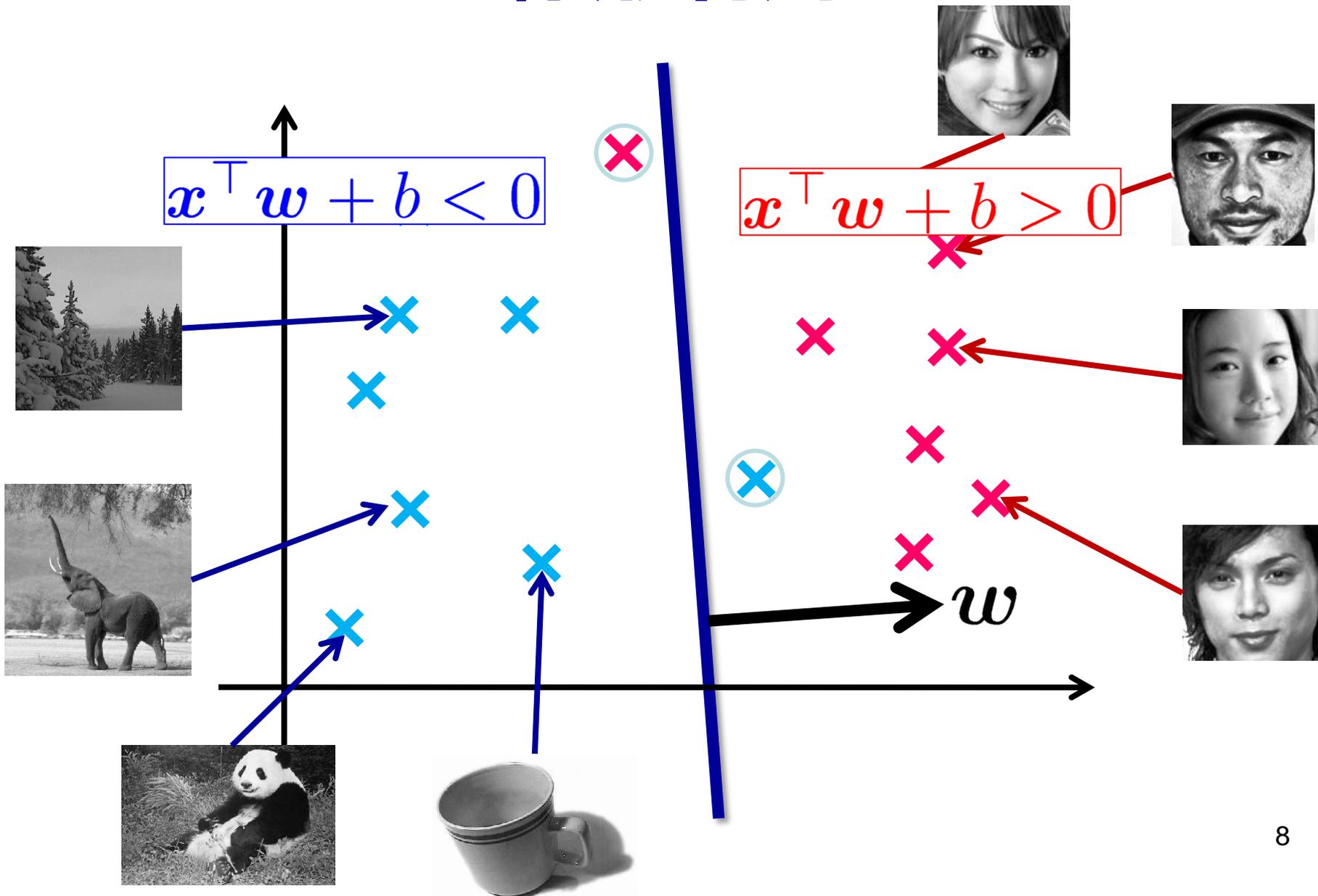
判別

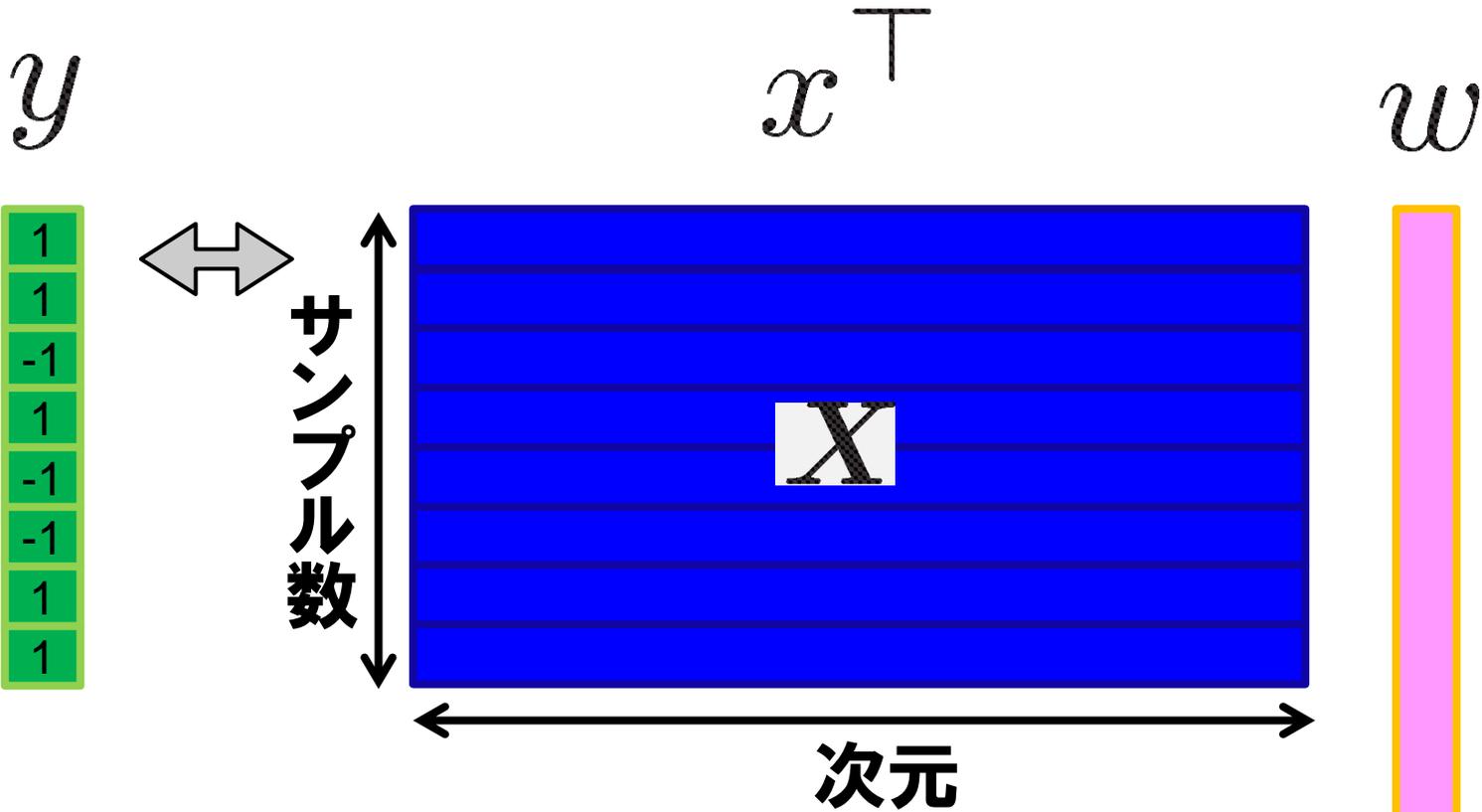


顔か？
顔でないか？



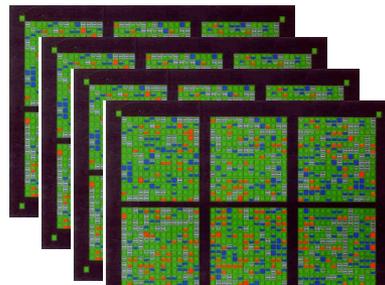
線形判別



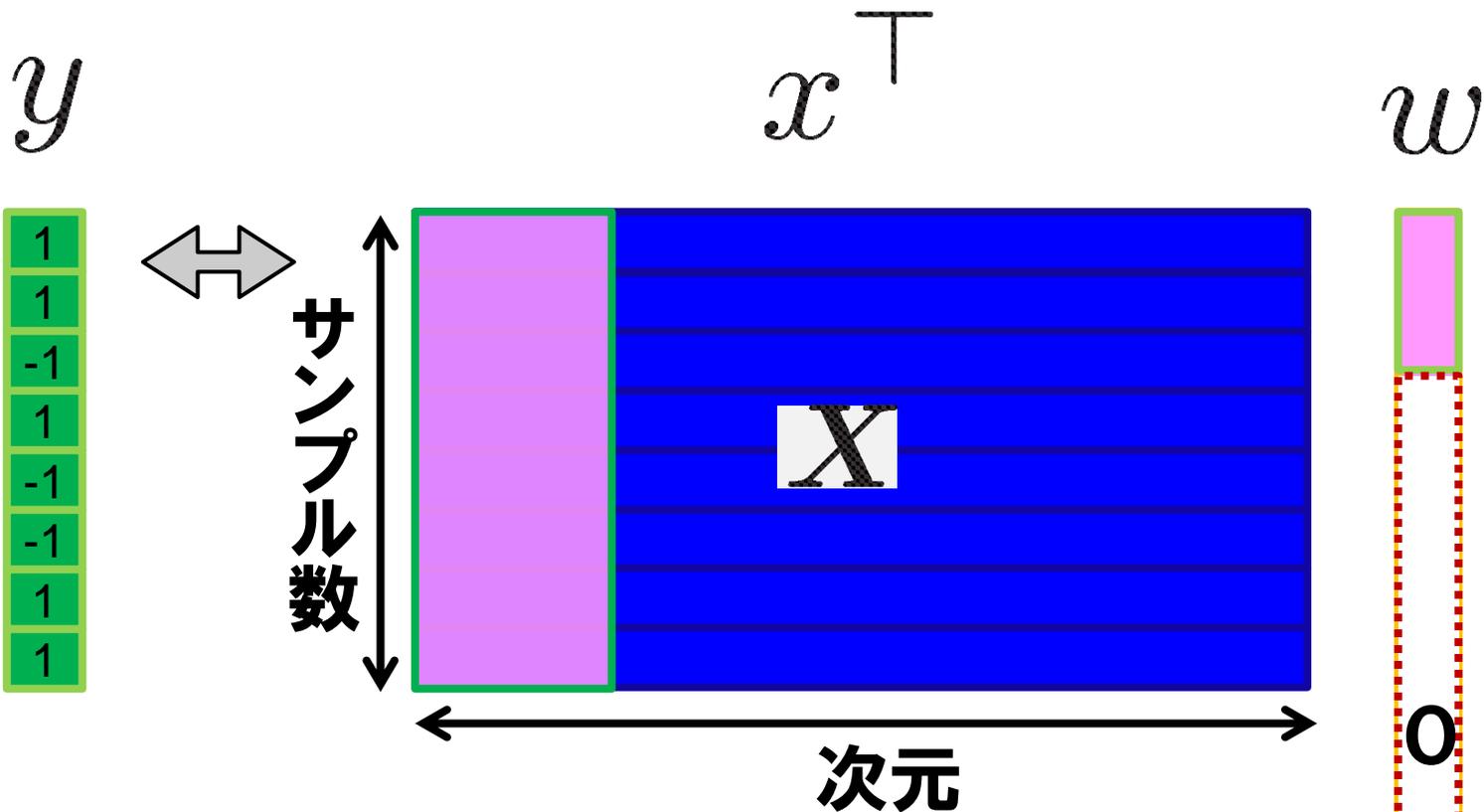


$\{(x_i, y_i)\}_{i=1}^n$: サンプル

$y_i \in \{+1, -1\}$

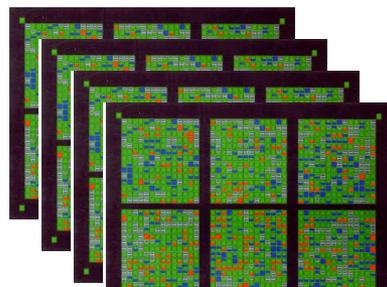


次元 > サンプル数 → 余分な情報を落としたい



$\{(x_i, y_i)\}_{i=1}^n$: サンプル

$y_i \in \{+1, -1\}$



次元 > サンプル数 → 余分な情報を落としたい

スパース推定:L1正則化

$$\min_{\mathbf{w}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i)}_{\text{間違いへのペナルティ}} + \underbrace{C \|\mathbf{w}\|_1}_{\text{平面の「複雑さ」}} \quad \text{L1ノルム} \rightarrow \text{スパース}$$

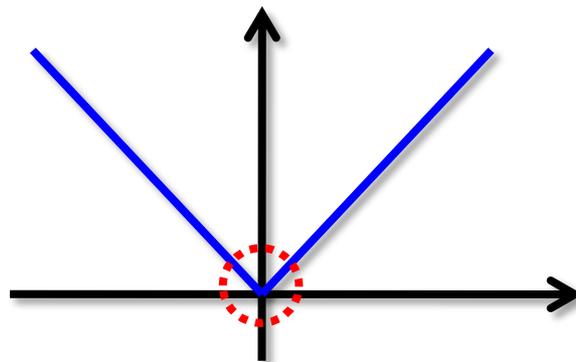
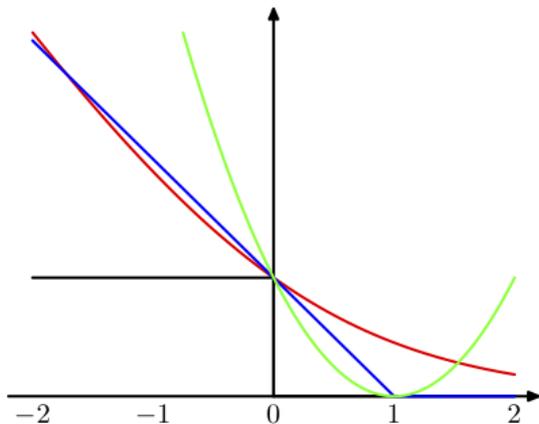
間違いへのペナルティ

平面の「複雑さ」

L1ノルム→スパース

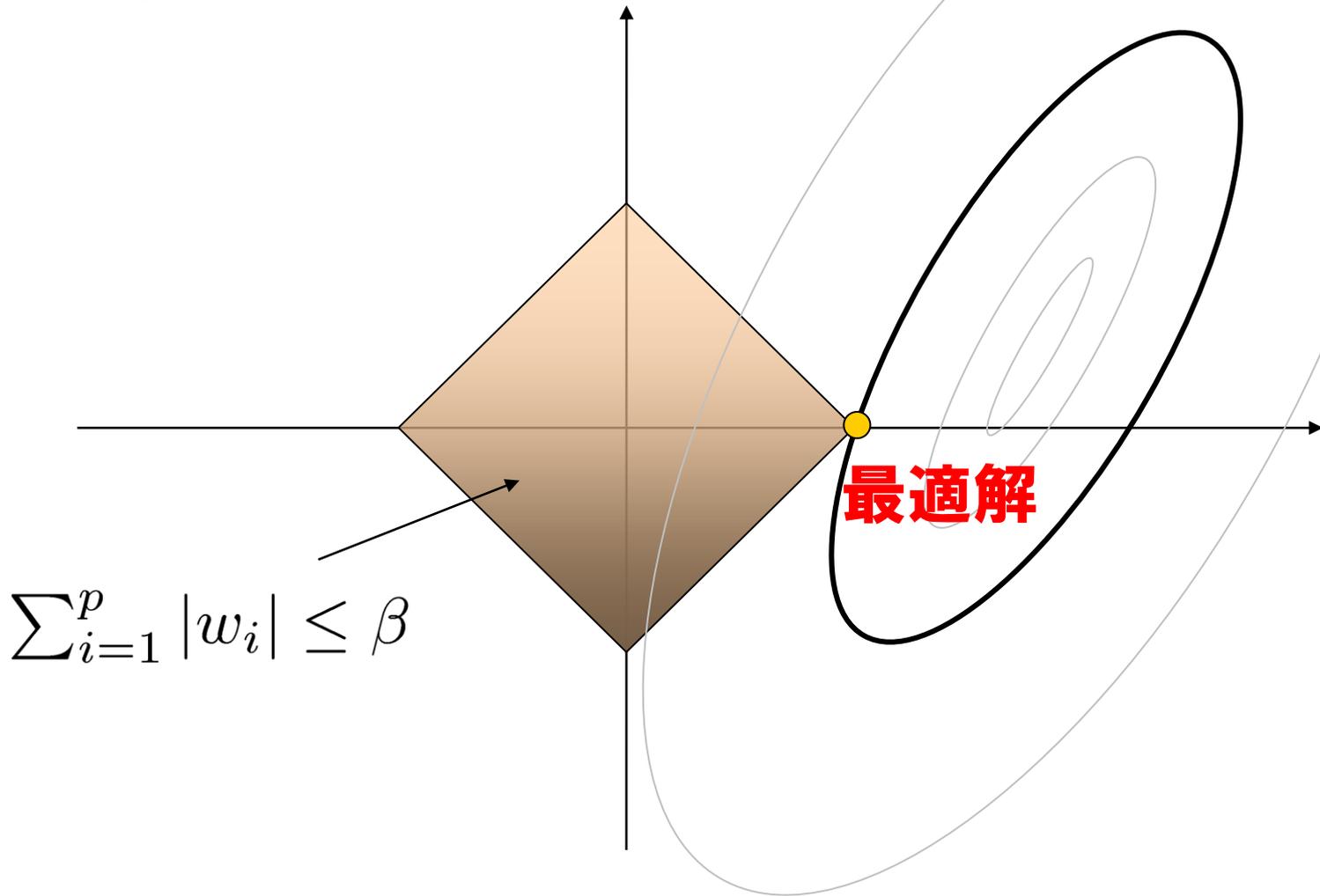
ℓ : ロス関数=間違いの度合いが大きいほど大きな値

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



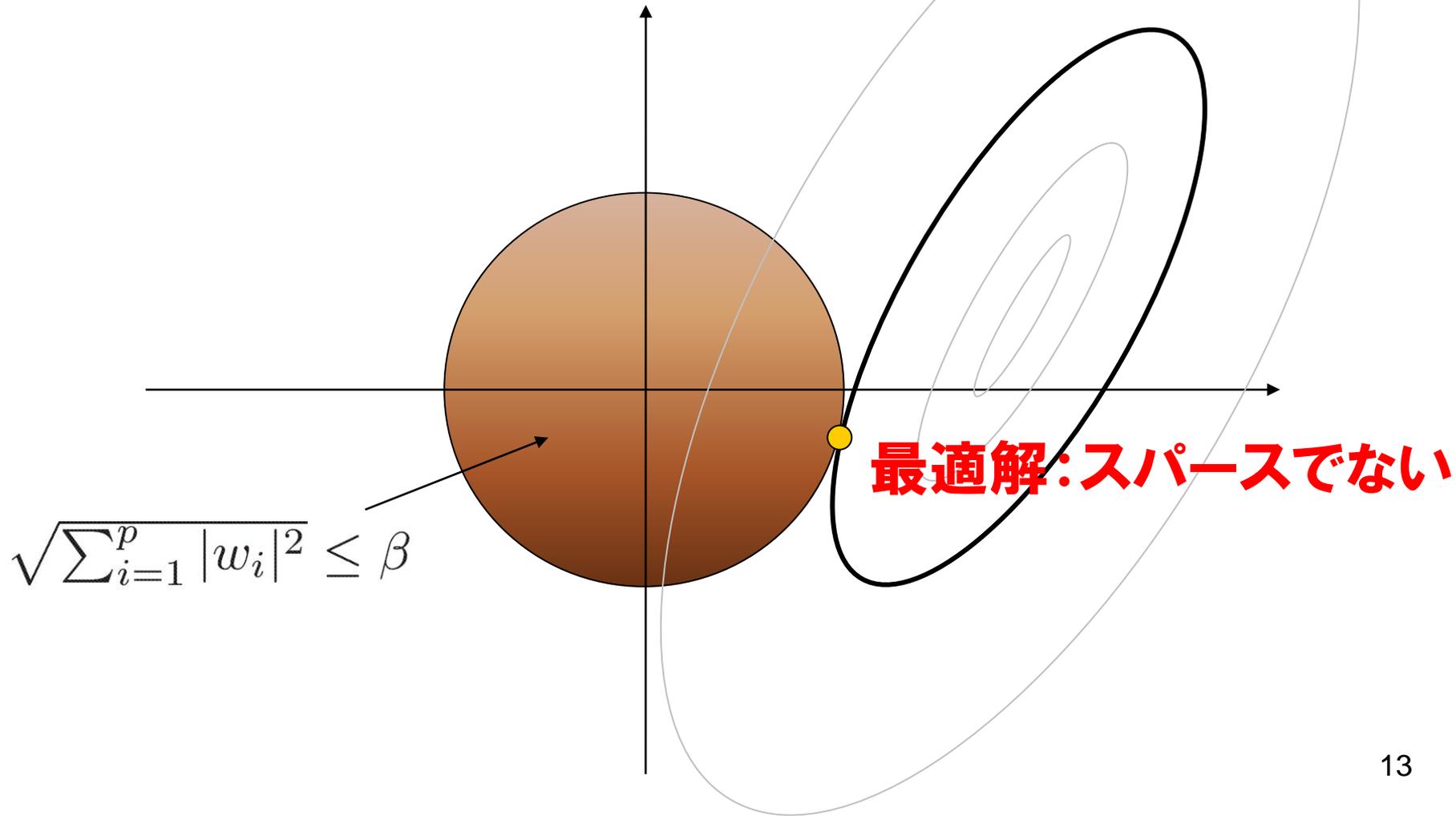
イメージ

$$\min_w \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i) \quad \text{s.t.} \quad \|\mathbf{w}\|_1 \leq \beta$$



L2正則化の場合

$$\min_{\mathbf{w}} \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i) \quad \text{s.t.} \quad \|\mathbf{w}\|_2 \leq \beta$$

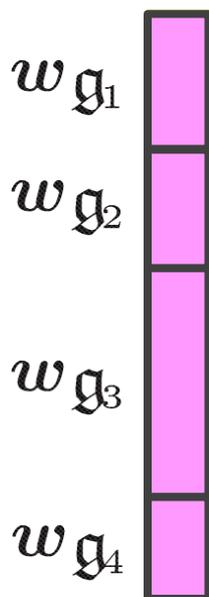


スパース正則化の例

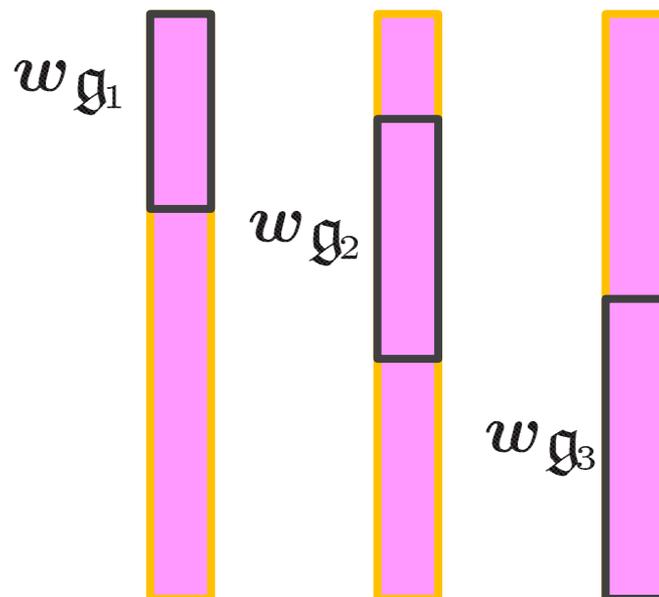
例1: Group Lasso

$$\min_w \frac{1}{n} \sum_{i=1}^n \ell(y_i, w^\top x_i) + C \sum_{g \in \mathcal{G}} \|w_g\|_2$$

グループ構造



重複あり



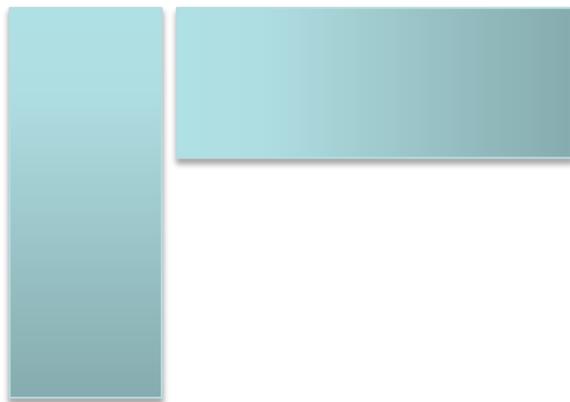
例2:低ランク行列推定

$$\min_w \frac{1}{n} \sum_{i=1}^n \ell(y_i, \text{Tr}(\mathbf{X}_i \mathbf{W})) + C \|\mathbf{W}\|_{\text{tr}}$$

$$\|\mathbf{W}\|_{\text{tr}} = \text{Tr}(\sqrt{\mathbf{W}^\top \mathbf{W}})$$

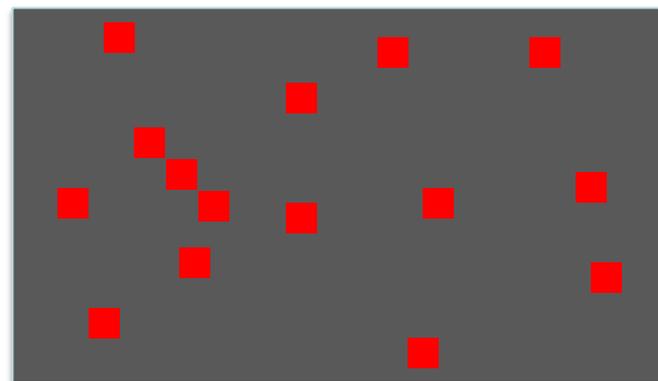
低ランク

$\mathbf{W} =$



ユーザの趣向推定

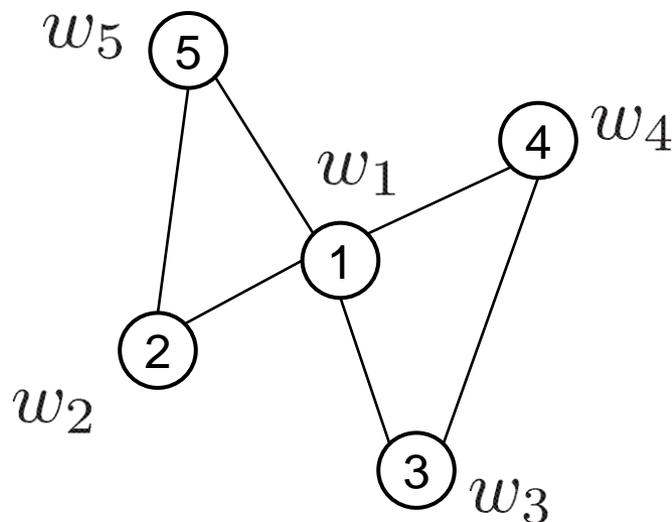
顧客



DVD

例3: グラフ型正則化

$$\min_w \frac{1}{n} \sum_{i=1}^n \ell(y_i, \mathbf{w}^\top \mathbf{x}_i) + C \sum_{(i,j) \in \mathcal{E}} |w_i - w_j|$$



スパース正則化学習の最適化

$$\min_w L(Xw) + \psi(w)$$

proximal point algorithm [Rockafellar, 1976]

$$w_{t+1} \leftarrow \min_w L(Xw) + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$$

$$\min_w \max_{\rho, y} \langle Xw, \rho \rangle - L^*(\rho) + \langle w, y \rangle - \psi^*(y) + \frac{1}{2\eta_t} \|w - w_t\|^2$$

双対

乗数法 [Hestenes, 1969; Powel 1969]

$$\min_{\rho, y} L^*(\rho) + \psi^*(y) - w_t^\top (X^\top \rho - y) + \frac{\eta_t}{2} \|X^\top \rho - y\|^2$$

$$w_{t+1} \leftarrow w_t - \eta_t (X^\top \rho^* - y^*)$$

proximal point alg.: $w_{t+1} \leftarrow \min_w L(Xw) + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$

乗数法: $\min_{\rho, y} L^*(\rho) + \psi^*(y) - w_t^\top (X^\top \rho - y) + \frac{\eta_t}{2} \|X^\top \rho - y\|^2$
 $w_{t+1} \leftarrow w_t - \eta_t (X^\top \rho^* - y^*)$

主問題

prox. point alg.

FOBOS [Duchi&Singer,2009]

FISTA [Beck&Teboulle,2009]

(近似解法・近接勾配法)

乗数法

ADMM

(Alternating Direction
Multiplier Method)

[Glowinski&Marrocco,75;Boyd et.al.,10]

(近似解法)

乗数法

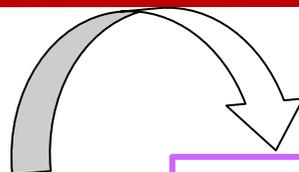
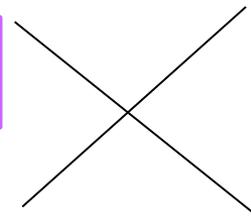
DAL

(Dual Augmented
Lagrangian)

[Tomioka&Sugiyama,2009;
Tomioka,Suzuki&Sugiyama,2011]

prox. point alg.

SpicyMKL



双対問題

FOBOS (Forward Backward Splitting)

Nesterov (2007), Duchi&Singer (2009), FISTA:Beck&Teboulle (2009)

prox. point alg. $w_{t+1} \leftarrow \min_w \underline{L(Xw)} + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$

線形近似

$$g_t \in \nabla_w L(Xw)|_{w=w_t}$$

$$w_{t+1} \leftarrow \min_w \langle g_t, w \rangle + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$$

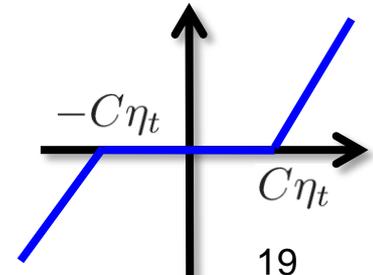
η_t : ステップ幅パラメータ

L1正則化での更新式 ($\psi(w) = C\|w\|_1$)

$$w_{t+1}^{(i)} = \text{sign}(w_t^{(i)} - \eta_t g_t^{(i)}) \max[|w_t^{(i)} - \eta_t g_t^{(i)}| - C\eta_t, 0]$$

最適化の途中でもスパースな解

Soft threshold



Proximal Operation

FOBOSは以下のproximal operationで更新

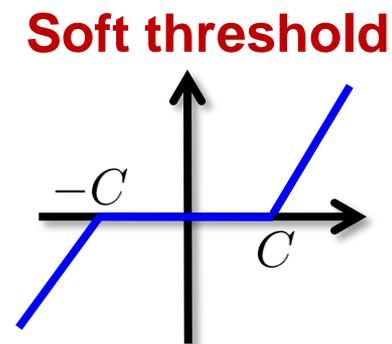
$$\text{prox}(\mathbf{q}|\psi) := \arg \min_{\mathbf{w}} \left\{ \psi(\mathbf{w}) + \frac{1}{2} \|\mathbf{w} - \mathbf{q}\|^2 \right\} \quad (\text{射影の一般化})$$

$$\text{FOBOS: } \mathbf{w}_{t+1} = \text{prox}(\mathbf{w}_t - \eta_t g_t | \eta_t \psi)$$

例: L1ノルムでのproximal operation ($\psi(\mathbf{w}) = C\|\mathbf{w}\|_1$)

$$\text{prox}(\mathbf{q} | C\|\cdot\|_1) = (\text{sign}(q_j) \max(|q_j| - C, 0))_j$$

- 各変数ごとの最適化に分離. 変数間の絡みがない.
- 解析解の存在.



収束レート

- 一般の凸ロス関数

$$L(X\mathbf{w}_T) + \psi(\mathbf{w}_T) - [L(X\mathbf{w}^*) + \psi(\mathbf{w}^*)] \leq C \frac{1}{\sqrt{T}}$$

- 滑らかな凸ロス関数

$$L(X\mathbf{w}_T) + \psi(\mathbf{w}_T) - [L(X\mathbf{w}^*) + \psi(\mathbf{w}^*)] \leq C \frac{1}{T}$$

- $X^\top X$ が正則で強凸かつ滑らかなロス関数

$$L(X\mathbf{w}_T) + \psi(\mathbf{w}_T) - [L(X\mathbf{w}^*) + \psi(\mathbf{w}^*)] \leq C \exp(-Tc)$$

proximal point alg.: $w_{t+1} \leftarrow \min_w L(Xw) + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$

乗数法: $\min_{\rho, y} L^*(\rho) + \psi^*(y) - w_t^\top (X^\top \rho - y) + \frac{\eta_t}{2} \|X^\top \rho - y\|^2$
 $w_{t+1} \leftarrow w_t - \eta_t (X^\top \rho^* - y^*)$

主問題

prox. point alg.

FOBOS [Duchi&Singer,2009]

FISTA [Beck&Teboulle,2009]

(近似解法・近接勾配法)

乗数法

ADMM

(Alternating Direction
Multiplier Method)

[Glowinski&Marrocco,75;Boyd et.al.,10]

(近似解法)

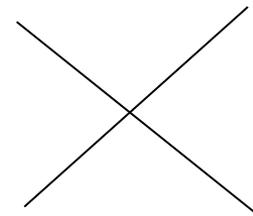
乗数法

DAL

(Dual Augmented
Lagrangian)

[Tomioka&Sugiyama,2009;
Tomioka, Suzuki&Sugiyama,2011]

prox. point alg.



双対問題

Dual Augmented Lagrangian (DAL)

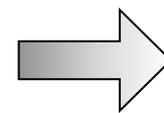
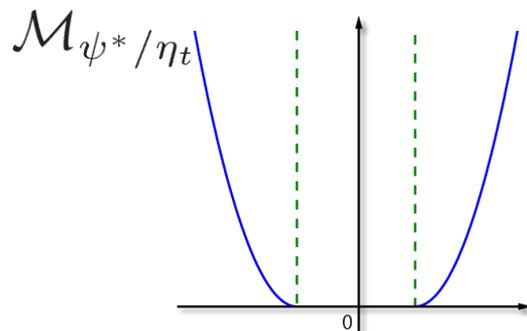
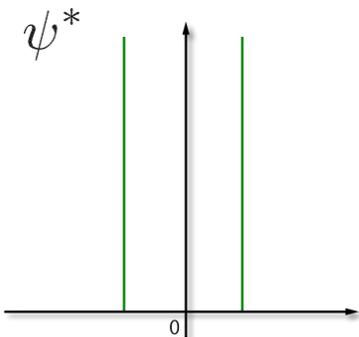
[Tomioka&Sugiyama,2009;
Tomioka,Suzuki&Sugiyama,2011]

$$(\boldsymbol{\rho}_{t+1}, \mathbf{y}_{t+1}) \leftarrow \arg \min_{\boldsymbol{\rho}, \mathbf{y}} L^*(\boldsymbol{\rho}) + \psi^*(\mathbf{y}) - \mathbf{w}_t^\top (\mathbf{X}^\top \boldsymbol{\rho} - \mathbf{y}) + \frac{\eta_t}{2} \|\mathbf{X}^\top \boldsymbol{\rho} - \mathbf{y}\|^2$$

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta_t (\mathbf{X}^\top \boldsymbol{\rho}_{t+1} - \mathbf{y}_{t+1})$$

$$\min_{\boldsymbol{\rho}} L^*(\boldsymbol{\rho}) + \underbrace{\eta_t \min_{\mathbf{y}} \left\{ \frac{\psi^*(\mathbf{y})}{\eta_t} + \frac{1}{2} \left\| \mathbf{y} - \mathbf{X}^\top \boldsymbol{\rho} + \frac{\mathbf{w}_t}{\eta_t} \right\|^2 \right\}}_{\text{Moreau's envelope}}$$

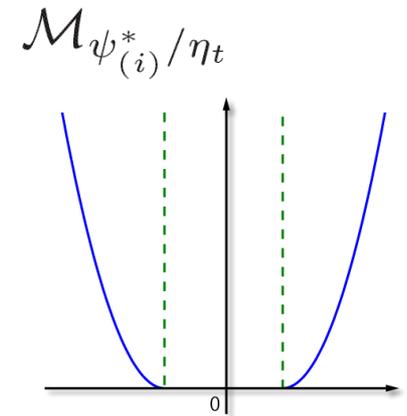
$\mathcal{M}_{\psi^*/\eta_t}(\mathbf{X}^\top \boldsymbol{\rho} - \mathbf{w}_t/\eta_t)$: Moreau's envelope



L^* が滑らかな場合,
 $\boldsymbol{\rho}$ に関するNewton法

DALの性質

$$\min_{\rho} L^*(\rho) + \eta_t \mathcal{M}_{\frac{\psi^*}{\eta_t}} \left(\mathbf{X}^\top \rho - \frac{\mathbf{w}_t}{\eta_t} \right)$$
$$\underbrace{\hspace{10em}}_{\sum_{i=1}^p \mathcal{M}_{\frac{\psi_{(i)}^*}{\eta_t}} \left(\mathbf{X}_i^\top \rho - \frac{w_t^{(i)}}{\eta_t} \right)}$$



- $\mathcal{M}_{\psi_{(i)}^*} / \eta_t$ の微分をスキップできる
→ **スパース性を利用した高速化**
- (超)一次収束

$$f(\mathbf{w}) - f(\mathbf{W}^*) \geq \sigma \|\mathbf{w} - \mathbf{W}^*\|^2 \text{ なら}$$

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

最適解が一意である必要はない

SpicyMKL [Suzuki&Tomioka, 2011]

DALのMKLへの拡張

Multiple Kernel Learning (MKL) [Lanckriet et al. 2004]

グループ正則化の各グループを無限次元の再生核ヒルベルト空間とした方法.

沢山のカーネル関数とそれに付随した再生核ヒルベルト空間 $\{\mathcal{H}_m\}_{m=1}^M$

$$\min_{f_m \in \mathcal{H}_m} \frac{1}{n} \sum_{i=1}^n \ell \left(y_i, \sum_{m=1}^M f_m(\mathbf{x}_i) \right) + C \sum_{m=1}^M \|f_m\|_{\mathcal{H}_m}$$

Sparse learning

Lasso

グループ化

Group Lasso

カーネル化

Multiple Kernel Learning (MKL)

ソフトウェアの公開

SpicyMKL

SpicyMKL is an optimization method of MKL (multiple kernel learning) that scales well against the number of kernels. The software affords hinge, logistic and square losses, and L1-norm and elasticnet regularization.

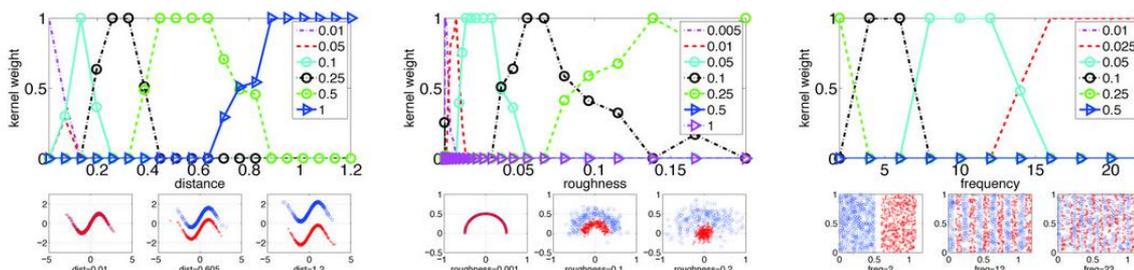
Matlab Implementation: [SpicyMKL.zip](#) (updated 1st Nov. 2010)

- SpicyMKL.m is the main function.
 - SpicyMKL.m calls normKj.mex and HessAugMexbias.mex. Please link blas library when compiling these files like

```
mex -lmwblas normKj.c  
mex -lmwblas HessAugMexbias.c
```

- (startup.m contains this code).
- demo.m is a demo script.
 - The demo code requires [SimpleMKL Toolbox](#) in the path.

Examples



<http://www.simplex.t.u-tokyo.ac.jp/~s-taiji/software/SpicyMKL/>

DAL



What is DAL?

DAL is an efficient and flexible MATLAB toolbox for solving the sparsity-regularized minimization problems, which arises often in machine learning, of the following form:

$$\underset{w \in \mathbb{R}^n}{\text{minimize}} \quad f(Aw) + \lambda g(w)$$

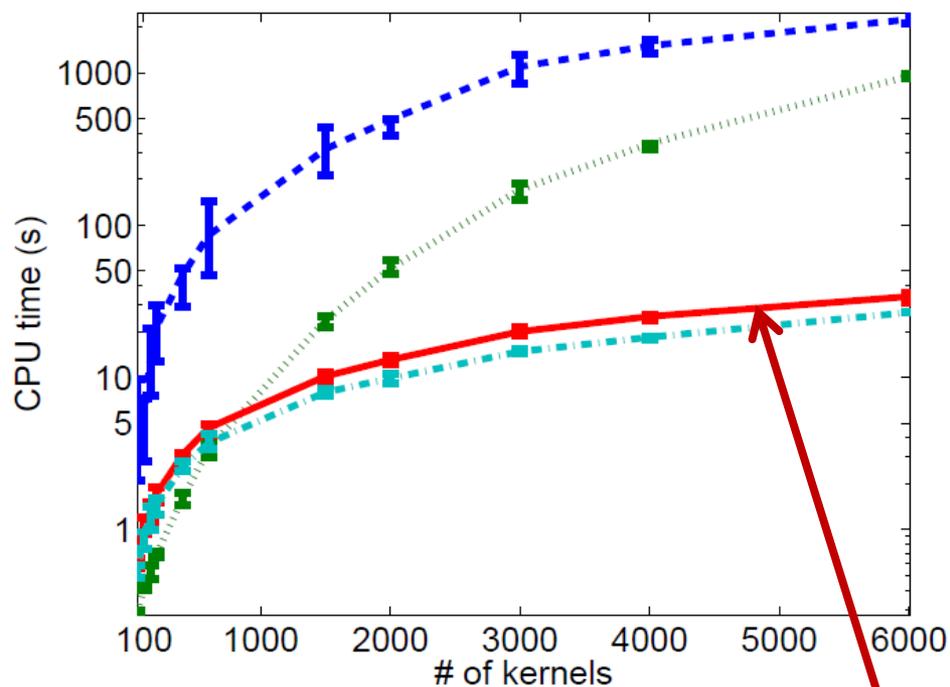
<http://www.ibis.t.u-tokyo.ac.jp/ryotat/dal/>

- DAL is efficient when $m \leq n$ (m : #samples, n : #unknowns) or the matrix A is poorly conditioned.

Matlabコード

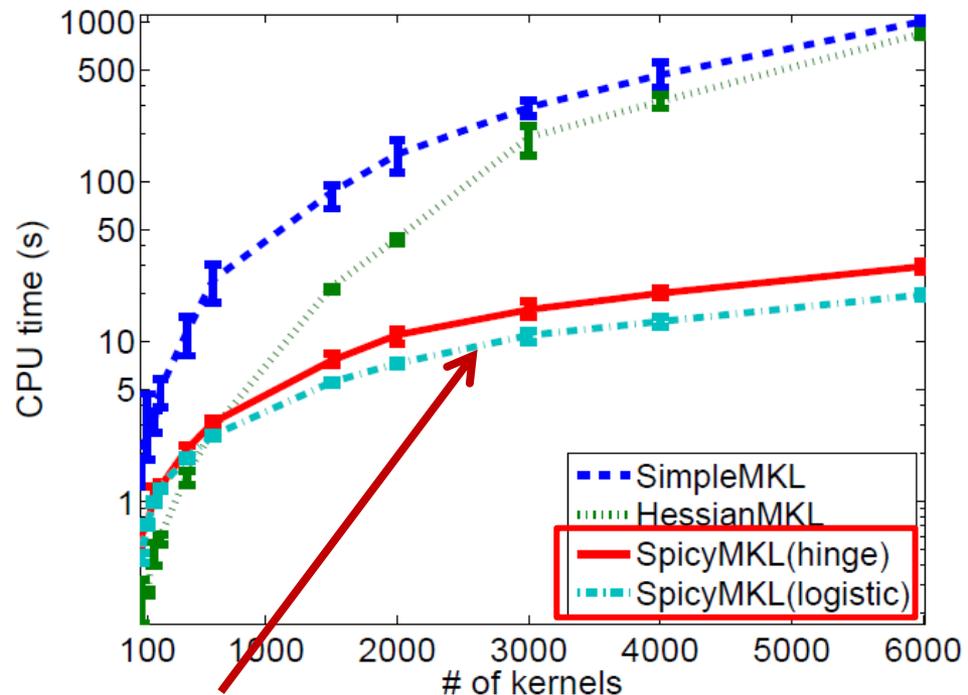
数値実験

CPU time v.s. # of kernels



Splice

SpicyMKL

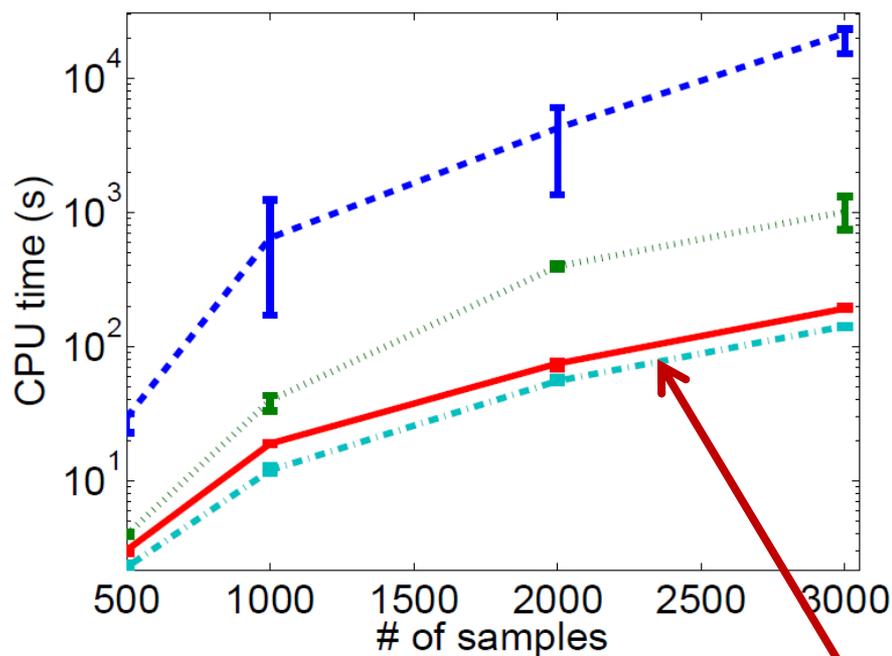


Ringnorm

•カーネルの数に対し良くスケールする

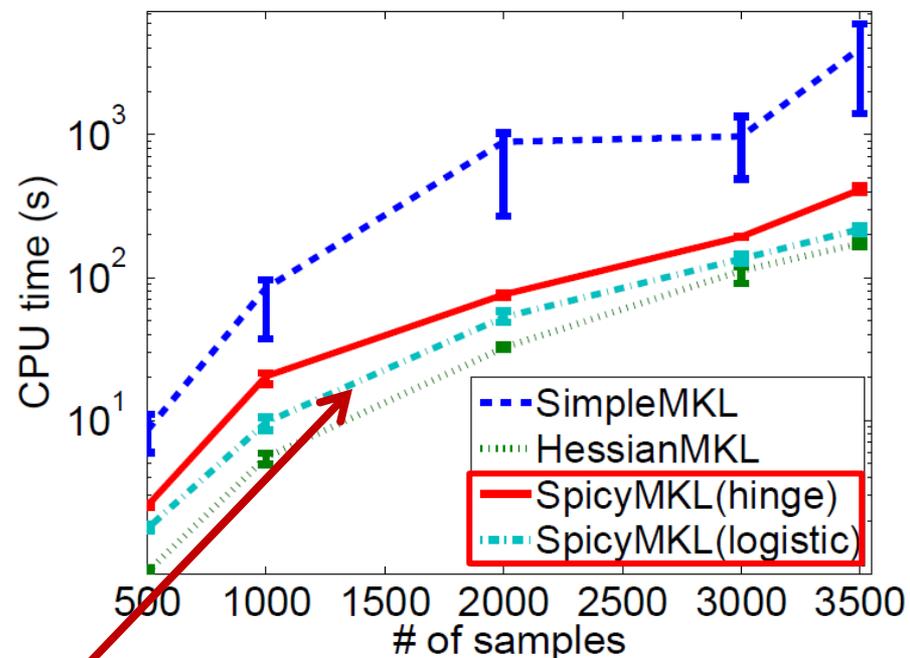
IDA data set, L1 regularization.

CPU time v.s. # of samples



Splice

SpicyMKL



Ringnorm

- サンプル数に対しては既存手法とほぼ同スケール

proximal point alg.: $w_{t+1} \leftarrow \min_w L(Xw) + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$

乗数法: $\min_{\rho, y} L^*(\rho) + \psi^*(y) - w_t^\top (X^\top \rho - y) + \frac{\eta_t}{2} \|X^\top \rho - y\|^2$
 $w_{t+1} \leftarrow w_t - \eta_t (X^\top \rho^* - y^*)$

主問題

prox. point alg.

FOBOS [Duchi&Singer,2009]

FISTA [Beck&Teboulle,2009]

(近似解法・近接勾配法)

乗数法

ADMM

(Alternating Direction
Multiplier Method)

[Glowinski&Marrocco,75;Boyd et al.,10]

(近似解法)

双対問題

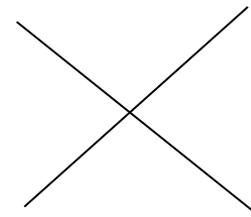
乗数法

DAL

(Dual Augmented
Lagrangian)

[Tomioka&Sugiyama,2009;
Tomioka,Suzuki&Sugiyama,2011]

prox. point alg.



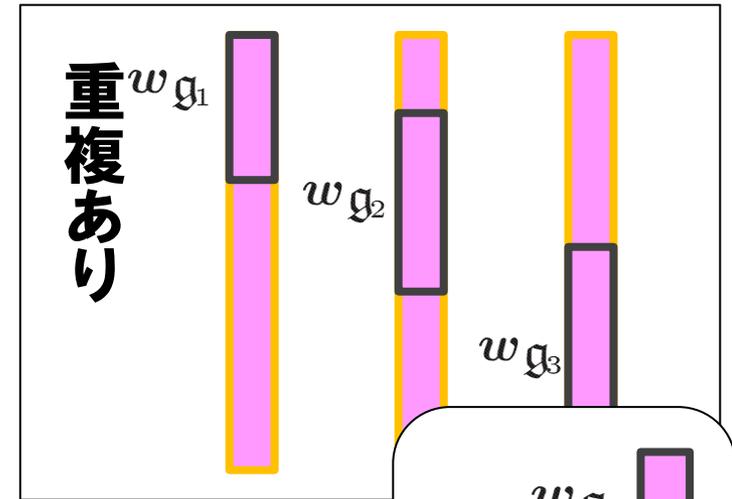
prox. op.が計算しにくい例

$$\text{prox}(q|\psi) := \arg \min_w \left\{ \psi(w) + \frac{1}{2} \|w - q\|^2 \right\}$$

- 重複ありグループ正則化

$$\psi(w) = C \sum_{g \in \mathcal{G}} \|w_g\|_2$$

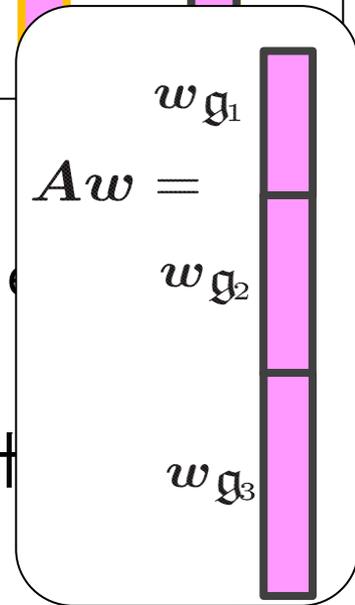
変数間に絡み



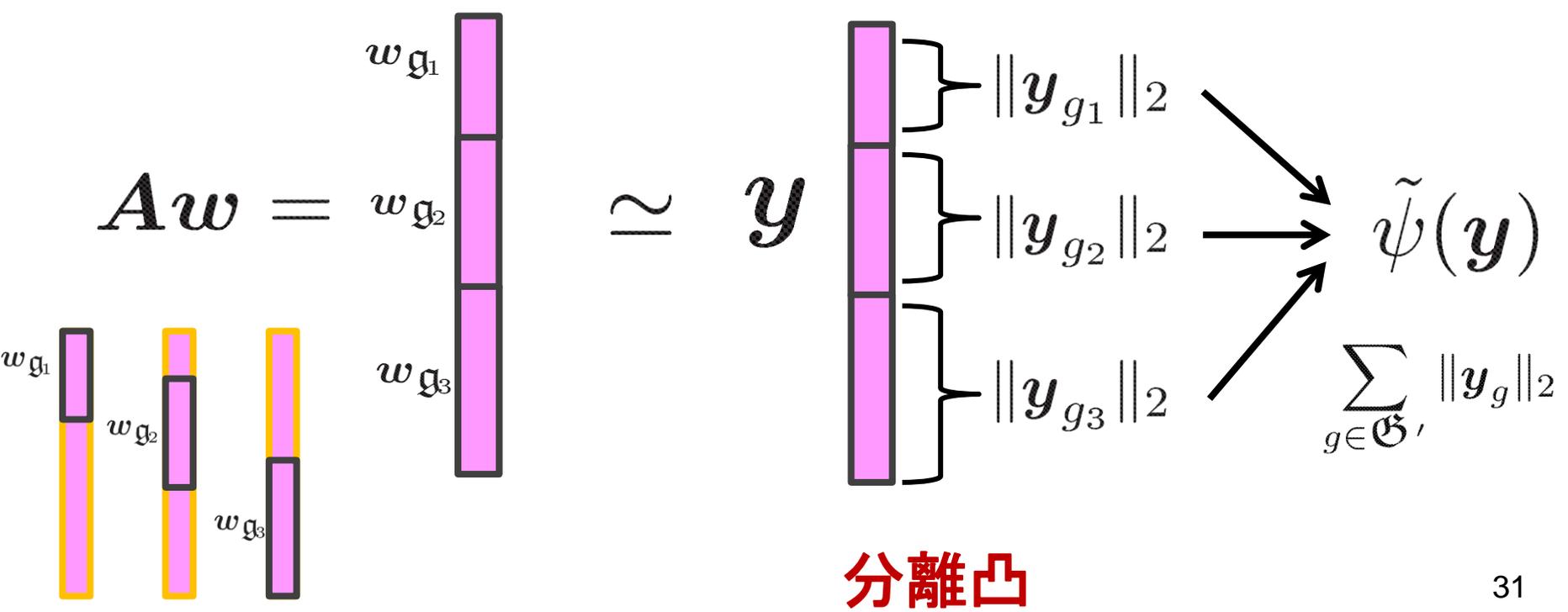
- 解決策

- 各正則化関数に応じた賢い方法で解く [Yuan et al.]
- 変数を増やして問題を簡単にする (汎用的)

idea: $\tilde{\psi}(Aw) = \psi(w)$ を満たし $\text{prox}(\cdot|\tilde{\psi})$ が計算可能な $\tilde{\psi}$ を利用する.



$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}} \quad & L(\mathbf{X}\mathbf{w}) + \tilde{\psi}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{w} = \mathbf{y} \end{aligned}$$



ADMM

(Alternating Direction Multiplier Method)

- 乗数法: w と y を同時最適化

$$(w_{t+1}, y_{t+1}) \leftarrow \min_{w, y} L(Xw) + \tilde{\psi}(y) - \lambda_t^\top (Aw - y) + \frac{\eta_t}{2} \|Aw - y\|^2$$

$$\lambda_{t+1} \leftarrow \lambda_t - \eta_t (Aw_{t+1} - y_{t+1})$$

- Splitting Technique

w と y の最適化を分離

ADMM

$$w_{t+1} \leftarrow \arg \min_w L(Xw) - \lambda_t^\top (Aw - y_t) + \frac{\rho}{2} \|Aw - y_t\|^2$$

$$y_{t+1} \leftarrow \arg \min_y \tilde{\psi}(y) + \lambda_t^\top y + \frac{\rho}{2} \|Aw_{t+1} - y\|^2$$

(prox($Aw_{t+1} - \lambda_t/\rho | \tilde{\psi}/\rho$))

$$\lambda_{t+1} \leftarrow \lambda_t - \rho (Aw_{t+1} - y_{t+1})$$

ADMMの収束レート

- リプシッツ連続凸ロス関数

$$L(X\mathbf{w}_T) + \psi(\mathbf{w}_T) - [L(X\mathbf{w}^*) + \psi(\mathbf{w}^*)] \leq C \frac{1}{T}$$

- $X^\top X$ が正則で強凸かつ滑らかなロス関数

$$L(X\mathbf{w}_T) + \psi(\mathbf{w}_T) - [L(X\mathbf{w}^*) + \psi(\mathbf{w}^*)] \leq C \exp(-Tc)$$

確率的最適化 (オンライン学習)

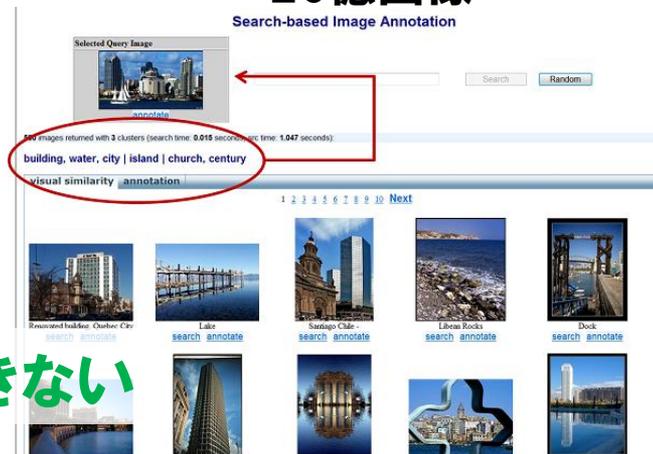
近年のデータ

Flickr
100万枚/日



ARISTA
20億画像

Search-based Image Annotation



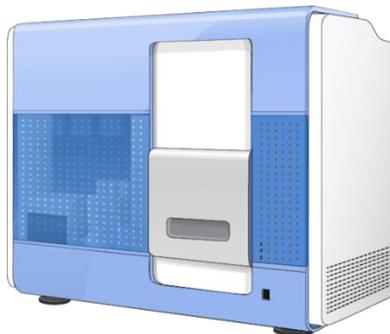
twitter

2億ツイート/日 (2011/6/30)

トルストイ「戦争と平和」8163冊分



次世代シーケンサ
60億本 x100塩基



多量

全データをメモリに保持できない

確率的最適化

prox. point alg. $w_{t+1} \leftarrow \min_w \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i^\top w) + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$

全データを保持するのではなく、
一つサンプルを見たらwを更新してサンプルを捨てる方法

- FOBOS (Forward Backward Splitting) [Duchi&Singer,2009]

$$g_t \in \nabla \ell_t(w_t)$$

$$\ell_t(w) = \ell(y_t, x_t^\top w)$$

$$w_{t+1} \leftarrow \min_w \langle g_t, w \rangle + \psi(w) + \frac{1}{2\eta_t} \|w - w_t\|^2$$

一つのサンプル・線形近似

- RDA (Regularized Dual Averaging) [Xiao,09; Nesterov,09]

$$\bar{g}_t = \frac{1}{t} \sum_{\tau=1}^t g_\tau$$

$$w_{t+1} \leftarrow \min_w \langle \bar{g}_t, w \rangle + \psi(w) + \frac{1}{2\eta_t} \|w\|^2$$

双対変数(勾配)の平均で過去の情報を保持

確率的ADMM

[Suzuki, ICML2013]

ADMM + 確率的最適化

$$\begin{aligned} \min_{\mathbf{w}, \mathbf{y}} \quad & L(\mathbf{X}\mathbf{w}) + \tilde{\psi}(\mathbf{y}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{w} = \mathbf{y} \end{aligned}$$

- FOBOS型ADMM

$$\mathbf{w}_{t+1} \leftarrow \Pi_{\mathcal{W}} \left[-\frac{\eta_t}{\gamma} \left\{ g_t - \mathbf{A}^\top (\boldsymbol{\lambda}_t - \rho(\mathbf{A}\mathbf{w}_t - \mathbf{y}_t)) \right\} + \mathbf{w}_t \right]$$

- RDA型ADMM

$$\mathbf{w}_{t+1} \leftarrow \Pi_{\mathcal{W}} \left[-\frac{\eta_t}{\gamma} \left\{ \bar{g}_t - \mathbf{A}^\top (\bar{\boldsymbol{\lambda}}_t - \rho(\mathbf{A}\bar{\mathbf{w}}_t - \bar{\mathbf{y}}_t)) \right\} \right]$$

$$\mathbf{y}_{t+1} \leftarrow \text{prox}(\mathbf{A}\mathbf{w}_{t+1} - \boldsymbol{\lambda}_t / \rho | \tilde{\psi} / \rho)$$

$$\boldsymbol{\lambda}_{t+1} \leftarrow \boldsymbol{\lambda}_t - \rho(\mathbf{A}\mathbf{w}_{t+1} - \mathbf{y}_{t+1})$$

実装が簡単！

収束レート

データ: $D_{1:T-1} = (\mathbf{x}_t, y_t)_{t=1}^{T-1}$

- 一般の凸ロス関数

$$\mathbb{E}_{D_{1:T-1}}[\mathbb{E}_{(\mathbf{x}, y)}[\ell(y, \mathbf{x}^\top \bar{\mathbf{w}}_T) + \psi(\bar{\mathbf{w}}_T) - \ell(y, \mathbf{x}^\top \mathbf{w}^*) - \psi(\mathbf{w}^*)]] \leq \frac{C}{\sqrt{T}}$$

- 強凸正則化関数

$$\mathbb{E}_{D_{1:T-1}}[\mathbb{E}_{(\mathbf{x}, y)}[\ell(y, \mathbf{x}^\top \bar{\mathbf{w}}_T) + \psi(\bar{\mathbf{w}}_T) - \ell(y, \mathbf{x}^\top \mathbf{w}^*) - \psi(\mathbf{w}^*)]] \leq C \frac{\log(T)}{T}$$

条件

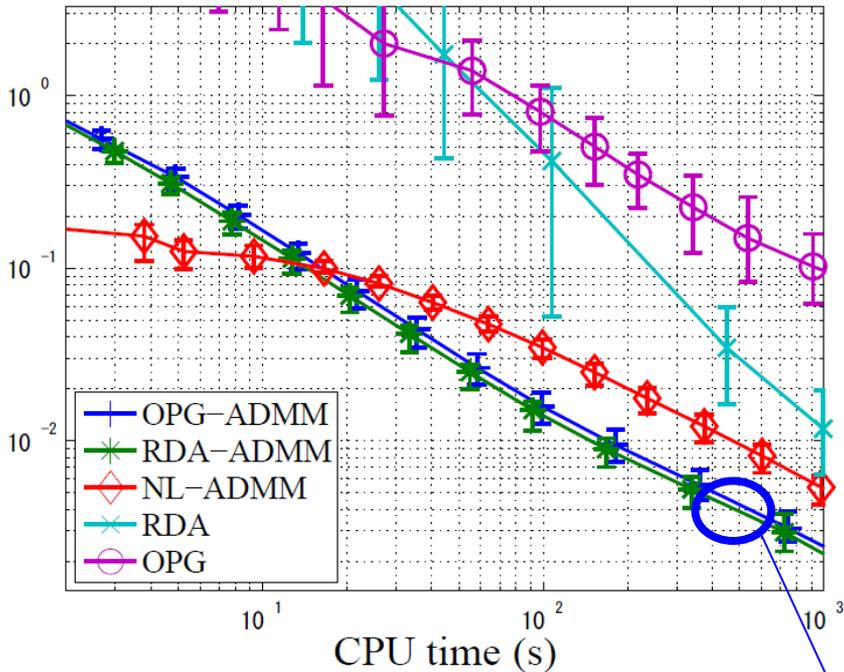
- データはi.i.d.系列
- ロスと正則化項はLipschitz連続
- \mathbf{w} のドメインは有界

数値実験: 確率的ADMM

人工データ

実データ (Adult, a9a
@LIVSVM data sets)

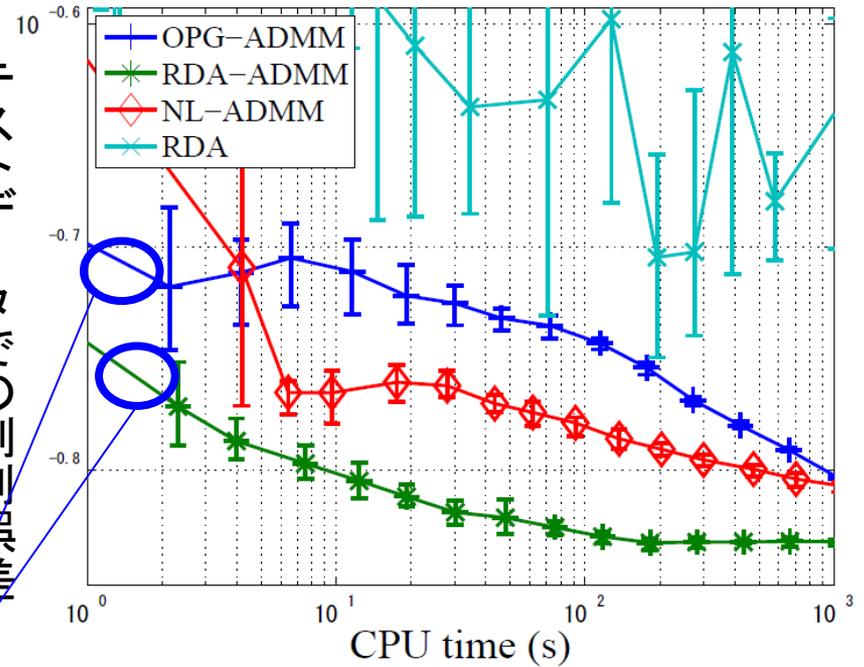
最適値との差



1,024次元
512サンプル
重複ありグループ正則化

提案手法

テストデータでの判別誤差



15,252次元
32,561サンプル
重複ありグループ正則化+ L1正則化

バッチデータに対する確率的最適化

Stochastic Dual Coordinate Ascent

[Shalev-Shwartz&Zhang,2012]

$$\min_{\rho} \frac{1}{n} \sum_{i=1}^n \ell_i^*(\rho_i) + \psi^* \left(-\frac{1}{n} X \rho \right)$$

1. i をランダムに選択($1 \leq i \leq n$)

2. 次元 i 方向に最適化

$$\min_{\Delta \rho_i} \frac{1}{n} \ell_i^*(\Delta \rho_i + \rho_i) + \psi^* \left(-\frac{1}{n} X \rho - X_i \frac{\Delta \rho_i}{n} \right)$$

3. 上の1,2を繰り返す.

ℓ_i^* が強凸で ψ^* が滑らかな時,

$$\text{双対ギャップの期待値} \leq C \left(1 - \frac{s}{n}\right)^T$$

まとめ

- 正則化学習の最適化法

proximal point algorithm \Leftrightarrow 乗数法
双対の関係

prox. point alg. FOBOS (近似解法・近接勾配法)	乗数法 ADMM (近似解法)
乗数法 DAL	prox. point alg.

- 確率的最適化法

- FOBOS・RDA

- 確率的ADMM

←最近のトレンド