

# オンライン学習

定式化

評価法: Regretなど

パーセプトロン

Passive Aggressive Algorithm

(アルゴリズムと損失の限界の評価)

Confidence Weighted Algorithm

Pegasos

Coordinate Descent

バッチ、オンライン、ストリームの比較

ビッグデータへの対応

# オンライン(あるいは逐次)学習とは

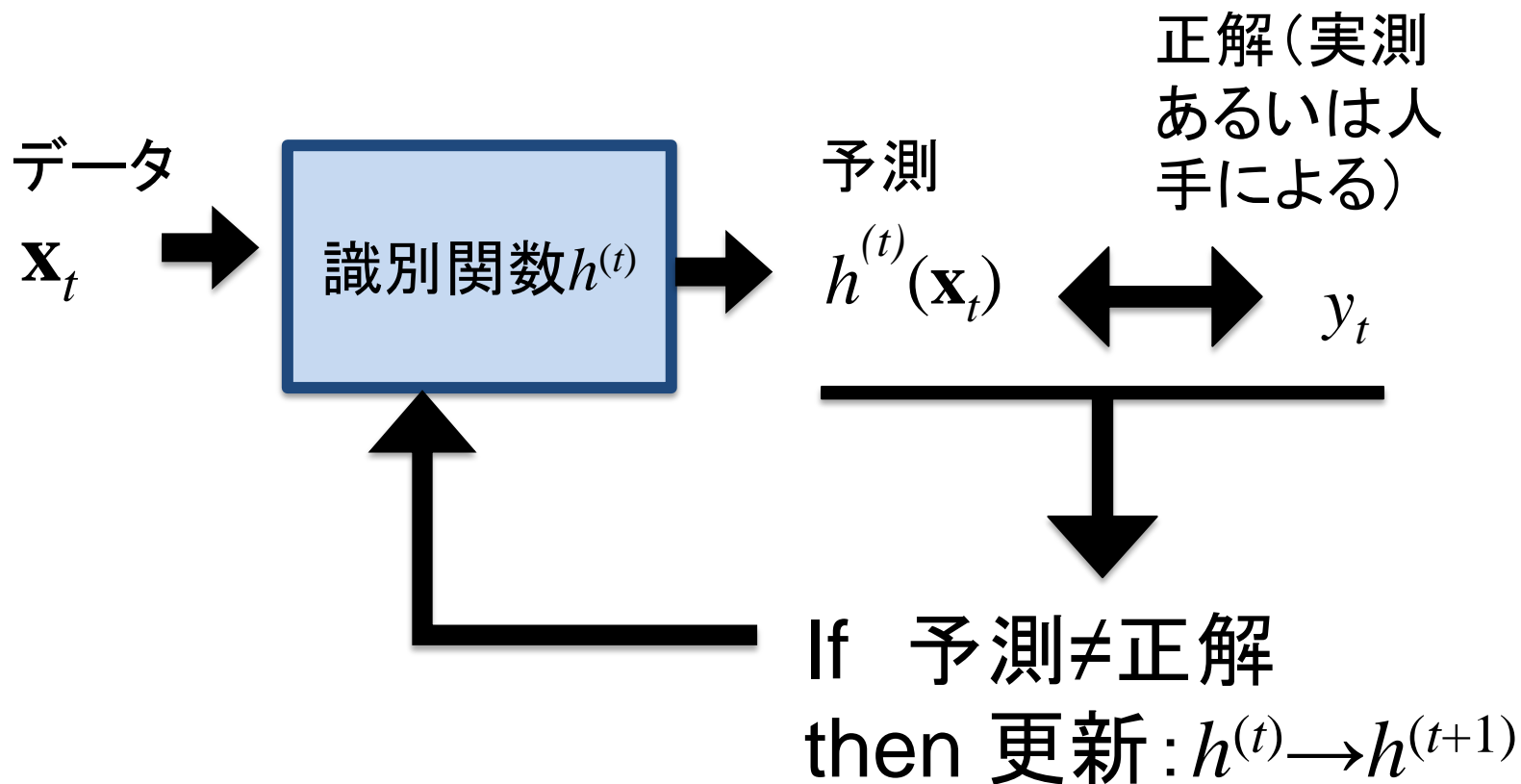
- データを1つずつ読み込んで、それまでの学習結果を更新する。
- 2つの利用局面
  1. データ全体は保持しているが、学習を1データ毎に行う
  2. データが1こずつ時系列としてやってくる
  - この場合はストリームという。
- データ全体をメモリの乗せなくてよいのでマシンに必要なメモリ少、あるいはメモリに乗りきらないほど大きなデータから学習可能
- 1個のデータからの学習(これを1roundという)だけなら高速

# オンライン学習の概観

以下1,2,3を時刻  $t=1,2,\dots,T$ で繰り返す

1. 時刻 $t$ において、仮説 $h_t$ 、入力データ $\mathbf{x}_t$ 、正しい結果データ  $y_t \in \mathcal{Y}$ が与えられる。
2. 仮説 $h_t$  による結果 $h^{(t)}(\mathbf{x}_t)$ を計算し、その後で $y_t$ との比較を損失関数 $\ell$ によって行う。つまり $\ell(h^{(t)}, (\mathbf{x}_t, y_t))$ を計算
  - $\ell$ としては2乗損失やヒンジ損失など
3. 損失関数 $\ell$ の値に応じて $h^{(t)}$  を更新し、新しい仮説  $h^{(t+1)}$ を求める
  - 最終的な目的の累積損失  $\sum_{t=1}^T \ell(h^{(t)}, (\mathbf{x}_t, y_t))$ などを最小化すること
  - 簡単(線形)な仮説として重みベクトル $\mathbf{w}$ と $\mathbf{x}$ の内積 $\langle \mathbf{w}, \mathbf{x} \rangle$ を使う場合は $h$ を $\mathbf{w}$ と書き、 $f_t(\mathbf{w}) = \ell(\mathbf{w}, (\mathbf{x}_t, y_t))$  と定義

# オンライン学習のイメージ



# オンライン学習の評価法

- 仮説 $h$ のなす空間を $\mathcal{H}$ ,  $t$ ラウンドの予測値を $h^{(t)}(\mathbf{x}_t)$
- 累積損失  $\sum_{t=1}^T \ell(h^{(t)}(\mathbf{x}_t), (\mathbf{x}_t, y_t))$  (最小化したい):

$$\text{Regret}_T(h^*) = \sum_{t=1}^T \ell(h^{(t)}(\mathbf{x}_t), (x_t, y_t)) - \sum_{t=1}^T \ell(h^*, (x_t, y_t)): h^* \in \mathcal{H}$$

$$\begin{aligned} \text{Regret}_T(\mathcal{H}) &= \max_{h^* \in \mathcal{H}} \text{Regret}_T(h^*) \\ &= \sum_{t=1}^T \ell(h^{(t)}(\mathbf{x}_t), (x_t, y_t)) - \min_{h^* \in \mathcal{H}} \sum_{t=1}^T \ell(h^*, (x_t, y_t)) \end{aligned}$$

- Mistake(失敗回数)のupper bound
  - 以後は識別に失敗しなくなるまでの学習回数＝学習データ数

# オンライン学習をオンライン凸最適化の観点から定式化

By Shai Shalev-Shwartz

- 以下では  $\mathcal{L}(\mathbf{w}, (x_i, y_i))$  を  $f_i(\mathbf{w})$  と略記することに留意。
- 最も簡単なオンライン学習は、過去の全roundの損失を最小化するような  $\mathbf{w}$  を選ぶ方法: Follow-The-Leader(FTL)

Follow – The – Leader (FTL)

$$\forall t \quad \mathbf{w}_t = \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} f_i(\mathbf{w}) \quad S \text{ は } \mathbf{w} \text{ の取り得る範囲で凸}$$

Lemma 10  $\mathbf{w}_1, \mathbf{w}_2 \dots$  を *FTL* で生成された重みベクトル

$$\forall \mathbf{u} \in S$$

$$\text{Regret}_T(\mathbf{u}) = \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1}))$$

Proof

$\sum_t f_t(\mathbf{w}_t)$ をLemma 10の不等式の両辺から引き移項すると

$$\sum_{t=1}^T f_t(\mathbf{w}_{t+1}) \leq \sum_{t=1}^T f_t(\mathbf{u}) \quad \text{そのこの不等式を帰納法で導く。}$$

base case :  $T = 1$ の場合は $\mathbf{w}_{t+1}$ の定義 $\mathbf{w}_2 = \arg \min_{\mathbf{w}} f_1(\mathbf{w})$ より、式の成立する $\mathbf{u}$ を選べる。

induction step :  $t = T - 1$ で不等式が成立すると仮定する。つまり

$$\forall \mathbf{u} \in S \quad \sum_{t=1}^{T-1} f_t(\mathbf{w}_{t+1}) \leq \sum_{t=1}^{T-1} f_t(\mathbf{u}) \quad \text{両辺に} f_T(\mathbf{w}_{T+1}) \text{を加えると}$$

$$\sum_{t=1}^T f_t(\mathbf{w}_{t+1}) \leq f_T(\mathbf{w}_{T+1}) + \sum_{t=1}^{T-1} f_t(\mathbf{u})$$

この式は $\forall \mathbf{u} \in S$ で成立し、 $\mathbf{u} = \mathbf{w}_{T+1}$ でも成り立つ

$$\Rightarrow \sum_{t=1}^T f_t(\mathbf{w}_{t+1}) \leq \sum_{t=1}^T f_t(\mathbf{w}_{T+1}) = \min_{\mathbf{u} \in S} \sum_{t=1}^T f_t(\mathbf{u})$$

最後の等式は $\mathbf{w}_{T+1}$ の定義 $\mathbf{w}_{T+1} = \arg \min_{\mathbf{w} \in S} \sum_{t=1}^T f_t(\mathbf{w})$ より。  $\square$

# Follow-The-Regularized-Leader (FoReL)

- FTLでは $\mathbf{w}$ に制約がないので、過学習が危ぶまれる。  
そこで、正則化項(Regularizer)を加えたものを最適化(FoReL)

FoReL

$$\forall t \quad \mathbf{w}_t = \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} f(\mathbf{w})_i + R(\mathbf{w}) \quad S \text{は} \mathbf{w} \text{の取り得る範囲で凸}$$

Lemma 20  $\mathbf{w}_1, \mathbf{w}_2 \dots$ をFoReLで生成された重みベクトル

$$\forall \mathbf{u} \in S \quad \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1}))$$

Proof Lemma 10で $t = 0..T$ とし、 $f_0 = R$

とおけばよい

□



# Example of FoReL: $R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2$

$f_t(\mathbf{w}) = \langle \mathbf{w}, \mathbf{z}_t \rangle$ かつ  $S = R^d$ で正則化項  $R(\mathbf{w}) = \frac{1}{2\eta} \|\mathbf{w}\|_2^2$  where  $\eta > 0$

この場合はFoReLは  $\mathbf{w}_{t+1} = \arg \min \sum_{i=1}^t \langle \mathbf{w}, \mathbf{z}_i \rangle + \frac{1}{2\eta} \|\mathbf{w}\|_2^2$

より、 $0 = \frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^t \langle \mathbf{w}, \mathbf{z}_i \rangle + \frac{1}{2\eta} \|\mathbf{w}\|_2^2 = \sum_{i=1}^t \mathbf{z}_i + \frac{\mathbf{w}}{\eta}$ を使えば

$$\Rightarrow \mathbf{w}_{t+1} = -\eta \sum_{i=1}^t \mathbf{z}_i = \mathbf{w}_t - \eta \mathbf{z}_t$$

つまり  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t) \quad (30)$

Online Gradient  
Descent: OGD

# FoReLのRegretのUpper Bound

## ➤ Theorem 30

$$f_t(\mathbf{w}) = \langle \mathbf{w}, \mathbf{z}_t \rangle, \quad S = R^d, \quad R(\mathbf{w}) = \frac{1}{2\eta} \|\mathbf{w}\|_2^2 \quad \text{とする}$$

$$\forall \mathbf{u} \quad \text{Regret}_T(\mathbf{u}) \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_2^2$$

$$\text{if } U = \{\mathbf{u} : \|\mathbf{u}\| \leq B\} \text{ and } \frac{1}{T} \sum_{t=1}^T \|\mathbf{z}_t\|_2^2 \leq L^2$$

$$\text{then } \eta = \frac{B}{L\sqrt{2T}} \quad \text{とおくと} \quad \text{Regret}_T(U) = \inf_{\mathbf{u} \in U} \text{Regret}_T(\mathbf{u}) \leq BL\sqrt{2T}$$

Proof Lemma 20と式(30)より

$$\begin{aligned} \text{Regret}_T(\mathbf{u}) &\leq R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1})) \\ &= \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{w}_{t+1}, \mathbf{z}_t \rangle = \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_2^2 \end{aligned}$$

Regretの上限が  
 $\sqrt{T}$ に比例してい  
ることに注目！

# 損失 $f$ が連続でない場合 Sub-gradient(劣勾配)のFoReL

## ➤ $f$ の凸性が重要

Lemma 50  $S$ が凸集合、 $f$ が $S$ 上の凸関数

$$\text{iff } \forall \mathbf{w} \in S, \exists \mathbf{z} \forall \mathbf{u} \in S, f(\mathbf{u}) \geq f(\mathbf{w}) + \langle \mathbf{u} - \mathbf{w}, \mathbf{z} \rangle \quad (50)$$

(50)を満たす $\mathbf{z}$ の集合を $f$ のsub - gradient と呼び  $\partial f(\mathbf{w})$  と書く。  
連続なら $\nabla f(\mathbf{w})$ と同じ。

これを使ったOnline Gradient Descentが以下。

$$\eta > 0, \mathbf{w}_1 = 0,$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{z}_t \text{ where } \mathbf{z}_t \in \partial f_t(\mathbf{w}_t)$$

# Sub-gradient の場合のFoReLのRegret Bound

再掲 : Lemma 20 and Theorem30:  $\mathbf{w}_1, \mathbf{w}_2 \dots$  をFoReLで生成された重みベクトル

$$\begin{aligned} \forall \mathbf{u} \in S \quad \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) &= \text{Regret}_T(\mathbf{u}) \leq R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1})) \\ &= \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{w}_{t+1}, \mathbf{z}_t \rangle = \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_2^2 \end{aligned} \quad (60)$$

Lemma 50を少し変形して再掲 :  $S$ が凸集合、 $f$ が $S$ 上の凸関数

$$\text{iff } \forall \mathbf{w}_t \in S, \exists \mathbf{z}_t \forall \mathbf{w}_{t+1} \in S, f(\mathbf{w}_{t+1}) \geq f(\mathbf{w}_t) + \langle \mathbf{w}_{t+1} - \mathbf{w}_t, \mathbf{z}_t \rangle \quad (50)$$

凸だと各round  $t$ に対して  $f_t(\mathbf{w}_{t+1}) \geq f_t(\mathbf{w}_t) + \langle \mathbf{w}_{t+1} - \mathbf{w}_t, \mathbf{z}_t \rangle$  だから

$$f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1}) \leq \langle \mathbf{w}_t - \mathbf{w}_{t+1}, \mathbf{z}_t \rangle \Rightarrow \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1})) \leq \sum_{t=1}^T (\langle \mathbf{w}_t, \mathbf{z}_t \rangle - \langle \mathbf{w}_{t+1}, \mathbf{z}_t \rangle)$$

これをLemma 20 and Theorem30にplug inすると sub - gradient OGDでも

$$\text{Regret}_T(\mathbf{u}) \leq R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1}))$$

$$\leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{w}_{t+1}, \mathbf{z}_t \rangle = \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_2^2$$

問題はこの部分

$$\text{Regret}_T(\mathbf{u}) \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_2^2 \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta TL^2 \quad \text{にしたいが}$$

そのためには  $\|\mathbf{z}_t\|$  が上から押さえられることを示す必要がある  $\Rightarrow$

sub - gradient の定義から

$$f(\mathbf{w}_t) - f(\mathbf{w}_{t+1}) \geq \langle \mathbf{z}_t, \mathbf{w}_t - \mathbf{w}_{t+1} \rangle \quad \text{where } \mathbf{z}_t \in \partial f(\mathbf{w}_t)$$

$$f \text{ が } L\text{-Lipshitz} \text{ だとすると } L\|\mathbf{w}_t - \mathbf{w}_{t+1}\| \geq f(\mathbf{w}_t) - f(\mathbf{w}_{t+1}) \quad \text{where } L < \infty$$

$$\text{上の2式を合わせると } L\|\mathbf{w}_t - \mathbf{w}_{t+1}\| \geq \langle \mathbf{z}_t, \mathbf{w}_t - \mathbf{w}_{t+1} \rangle$$

$$\eta \mathbf{z}_t = \mathbf{w}_t - \mathbf{w}_{t+1} \text{ だったので } L\|\eta \mathbf{z}_t\| \geq \eta \|\mathbf{z}_t\|_2^2 \Rightarrow \infty > L \geq \|\mathbf{z}_t\|$$

$$\Rightarrow \text{Regret}_T(\mathbf{u}) \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \eta TL^2 \leq BL\sqrt{2T} \quad \text{same as Theorem 30}$$

このboundは  $\frac{1}{\sqrt{2}}$  にできる。  
付録参照

# FoReLの上界を厳しくする

- まず、FoReLの別形式を導入する

$$\begin{aligned}\mathbf{w}_{t+1} &= \arg \min_{\mathbf{w}} R(\mathbf{w}) + \sum_{i=1}^t \langle \mathbf{w}, \mathbf{z}_i \rangle \\ &= \arg \max_{\mathbf{w}} \left\langle \mathbf{w}, -\sum_{i=1}^t \mathbf{z}_i \right\rangle - R(\mathbf{w})\end{aligned}$$

ここで  $g(\boldsymbol{\theta}) = \arg \max_{\mathbf{w}} \langle \mathbf{w}, \boldsymbol{\theta} \rangle - R(\mathbf{w})$  とおく とFoReLは次式で書ける

initialize  $\boldsymbol{\theta}_1 = 0$

for  $t = 1, \dots, T$

$\mathbf{w}_t = g(\boldsymbol{\theta}_t);$

$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{z}_t; \quad (\mathbf{z}_t \in \partial f_t(\mathbf{w}_t))$

Online Mirror Descent  
(OMD)という

# 数学的ツールの準備

- Fenchel-Young 不等式

$$f^*(\boldsymbol{\theta}) = \max_{\mathbf{u}} \langle \mathbf{u}, \boldsymbol{\theta} \rangle - f(\mathbf{u}) \text{ とすると}$$

$$\forall \mathbf{u}, \quad f^*(\boldsymbol{\theta}) \geq \langle \mathbf{u}, \boldsymbol{\theta} \rangle - f(\mathbf{u})$$

$$\mathbf{u} \in \partial f^*(\boldsymbol{\theta}) \quad \text{あるいは} \boldsymbol{\theta} \text{ で微分可能なら} \quad \mathbf{u} = \nabla f^*(\boldsymbol{\theta})$$

$$\Rightarrow \quad f^*(\boldsymbol{\theta}) = \langle \mathbf{u}, \boldsymbol{\theta} \rangle - f(\mathbf{u})$$

$$R(\mathbf{w}) = \frac{\|\mathbf{w}\|^2}{2\eta} \Rightarrow R^*(\boldsymbol{\theta}) = \max_{\mathbf{w}} \left( \langle \mathbf{w}, \boldsymbol{\theta} \rangle - \frac{\|\mathbf{w}\|^2}{2\eta} \right) = \frac{\eta \|\boldsymbol{\theta}\|^2}{2}$$

# 数学的ツールの準備

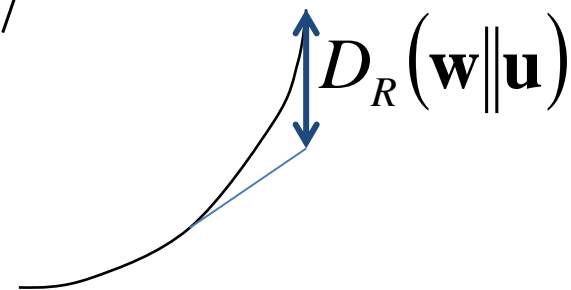
- Bregman Divergence:  $D_R$

$$D_R(\mathbf{w} \parallel \mathbf{u}) = R(\mathbf{w}) - (R(\mathbf{u}) + \langle \nabla R(\mathbf{u}), \mathbf{w} - \mathbf{u} \rangle)$$

$$R^*(\mathbf{z}) = \frac{\eta \|\mathbf{z}\|^2}{2} \quad \text{だと} \quad \nabla R^*(\mathbf{z}) = \eta \|\mathbf{z}\|$$

$$D_{R^*}(\mathbf{z}_1 \parallel \mathbf{z}_2) = \frac{\eta}{2} \|\mathbf{z}_1\|^2 - \frac{\eta}{2} \|\mathbf{z}_2\|^2 - \langle \eta \|\mathbf{z}_2\|, \mathbf{z}_1 - \mathbf{z}_2 \rangle$$

$$= \frac{\eta}{2} \|\mathbf{z}_1 - \mathbf{z}_2\|^2 \quad (DR1)$$





補題 OML1: Online Mirror Descent で  $g = \nabla R^*$  であるなら

$$\sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{u}, \mathbf{z}_t \rangle \leq R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T D_{R^*} \left( -\sum_{i=1}^t \mathbf{z}_i \parallel -\sum_{i=1}^{t-1} \mathbf{z}_i \right)$$

Proof

Fenchel - Young 不等式から

$$R(\mathbf{u}) + \sum_{t=1}^T \langle \mathbf{u}, \mathbf{z}_t \rangle = R(\mathbf{u}) - \left\langle \mathbf{u}, -\sum_{t=1}^T \mathbf{z}_t \right\rangle \geq -R^* \left( -\sum_{t=1}^T \mathbf{z}_t \right) \quad (1)$$

$$\mathbf{w}_t = g(\boldsymbol{\theta}_t) = \nabla R^* \left( -\sum_{i=1}^{t-1} \mathbf{z}_i \right) \quad (2)$$

Bregman Divergenceの定義より

$$\begin{aligned} -R^* \left( -\sum_{t=1}^T \mathbf{z}_t \right) &= -R^*(0) - \sum_{t=1}^T \left( R^* \left( -\sum_{i=1}^t \mathbf{z}_i \right) - R^* \left( -\sum_{i=1}^{t-1} \mathbf{z}_i \right) \right) \\ &= -R^*(0) + \sum_{t=1}^T \left( \langle \mathbf{w}_t, \mathbf{z}_t \rangle - D_{R^*} \left( -\sum_{i=1}^t \mathbf{z}_i \parallel -\sum_{i=1}^{t-1} \mathbf{z}_i \right) \right) \quad (3) \end{aligned}$$

なお、 $R^*(0) = \max_{\mathbf{w}} \langle 0, \mathbf{w} \rangle - R(\mathbf{w}) = -\min_{\mathbf{w}} R(\mathbf{w}) = -R(\mathbf{w}_1)$  (4) これらを合わせると

$$\begin{aligned} -\sum_{t=1}^T \langle \mathbf{u}, \mathbf{z}_t \rangle &\leq R(\mathbf{u}) + R^* \left( -\sum_{t=1}^T \mathbf{z}_t \right) = R(\mathbf{u}) - R^*(0) + \sum_{t=1}^T \left( \langle \mathbf{w}_t, \mathbf{z}_t \rangle - D_{R^*} \left( -\sum_{i=1}^t \mathbf{z}_i \parallel -\sum_{i=1}^{t-1} \mathbf{z}_i \right) \right) \\ &= R(\mathbf{u}) - R(\mathbf{w}_1) + \sum_{t=1}^T \left( -\langle \mathbf{w}_t, \mathbf{z}_t \rangle + D_{R^*} \left( -\sum_{i=1}^t \mathbf{z}_i \parallel -\sum_{i=1}^{t-1} \mathbf{z}_i \right) \right) \end{aligned}$$

## 定理 OMD2

Online Mirror Descentにおいて

$$R(\mathbf{w}) = \frac{\|\mathbf{w}\|^2}{2\eta}, \left( \text{つまり } R(\boldsymbol{\theta})^* = \frac{\eta\|\boldsymbol{\theta}\|^2}{2} \text{である} \right) \text{とする。}$$

$$\Rightarrow \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{u}, \mathbf{z}_t \rangle \leq \frac{\|\mathbf{u}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{z}_t\|^2 \quad (OMD10)$$

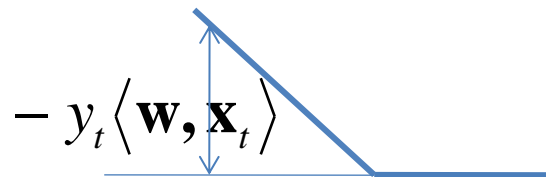
Proof 補題OMD1 と (DR1)より明らか ■

$$L^2 = \frac{1}{T} \sum_{t=1}^T \|\mathbf{z}_t\|^2 \quad \text{かつ} \quad \eta = \frac{B}{L\sqrt{T}} \text{ただし } \|\mathbf{u}\| \leq B \quad \text{とすると}$$

$$\sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{u}, \mathbf{z}_t \rangle \leq BL\sqrt{T} \quad (OMD20)$$

# パーセプトロン(Perceptron)

- FoReLから導出されたOnline Gradient Descentの例としてパーセプトロンを紹介する。
- パーセプトロンはF. Rosenblattが1956年に提案した線形識別の繰り返しアルゴリズム
- 効率がよく、現在でもその価値が高い
- 入力 $x_t$ が目的のクラスに
  - 属する場合に $y_t=1$ , 属さない場合に $y_t=-1$
- $f_t(\mathbf{w}) = [-y_t \langle \mathbf{w}, \mathbf{x}_t \rangle]_+$  右図
- $f_t(0) = 0$  これより失敗側では正



➤  $f_t(w)$ のsub-gradientを計算すると

$$\mathbf{z}_t \in \partial f_t(\mathbf{w}_t) \Rightarrow \text{if } y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle > 0 \text{ then } \nabla f_t(\mathbf{w}_t) = \mathbf{z}_t = 0$$

$$\text{otherwise } \nabla f_t(\mathbf{w}_t) = \mathbf{z}_t = -y_t \mathbf{x}_t \in \partial f_t(\mathbf{w}_t)$$

Online Gradient Descentの形式 :  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{z}_t$  に当てはめると

$$\mathbf{w}_{t+1} = \begin{cases} \mathbf{w}_t & \text{if } y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle > 0 \\ \mathbf{w}_t + \eta y_t \mathbf{x}_t & \text{otherwise} \end{cases}$$

➤  $\eta$ は正

➤ 次にPerceptronのアルゴリズムが得られる。

➤ FoReLの別形式として導入したOnline Mirror Descentとみれば、(OMD20)の上界が使える

# Perceptron アルゴリズム

初期化 :  $\mathbf{w}_1 = 0$

for  $t = 1, 2, \dots, T$

入力 =  $\mathbf{x}_t$

if  $y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle \leq 0$

then  $\mathbf{w}_{t+1} = \mathbf{w}_t + \eta y_t \mathbf{x}_t$

else  $\mathbf{w}_{t+1} = \mathbf{w}_t$

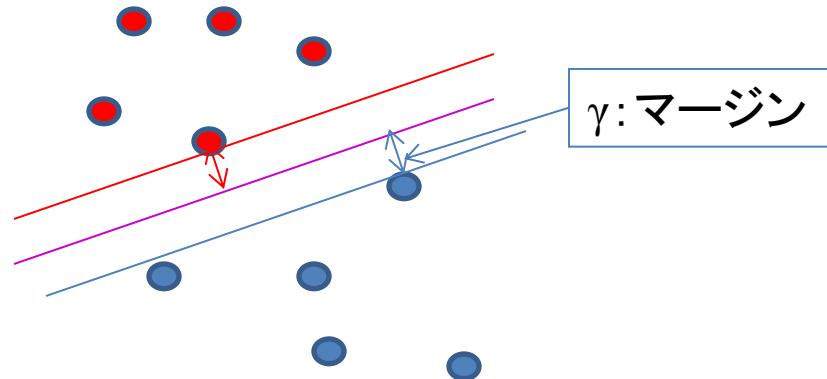
分類に失敗したときだけ

そのデータを分類器Wに足し込むという至って単純な更新

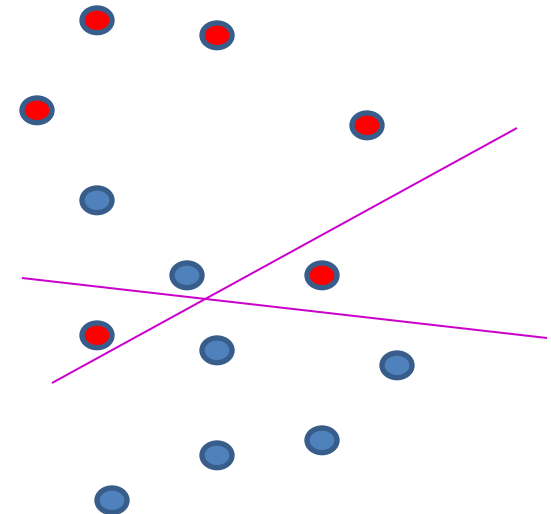
# 線形分離可能性

- 線形分離可能: クラスを識別する超平面が存在する場合
  - そうでない場合を線形分離不可能という。
  - 下図参照

線形分離可能



線形分離不可能



# Perceptronアルゴリズムの分析

FoReLの別形式として導入したOnline Mirror Descentとみれば、(OMD20)の上界が使える

$$\text{Regret}_T(\mathbf{u}) = \sum_{t=1}^T f_t(\mathbf{w}_t) - \sum_{t=1}^T f_t(\mathbf{u}) \leq \frac{1}{2\eta} \|\mathbf{u}\|_2^2 + \frac{\eta}{2} \sum_{t=1}^T \|\mathbf{z}_t\|_2^2 \quad (60)$$

ここで解析の容易さのため  $f_t = [1 - y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle]_+$  で近似する。  
(もとの  $f_t$  より必ず大きいので上界は甘くなる。)

sub - gradientは  $\mathbf{z}_t = -\mathbf{1}_{[y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle \leq 0]} y_t \mathbf{x}_t$

$\mathbf{w}_t, \mathbf{x}_t$  をスケール変換して、最も0に近い正例で  $\langle \mathbf{w}_t, \mathbf{x}_t \rangle = 1$  となったと見なしたと考えてもよい

# Perceptronアルゴリズムの分析

判定の失敗: mistakeを起こした $\mathbf{x}_t$ の集合を $\mathcal{M}$ 、mistake回数 $=|\mathcal{M}|$

$$\Rightarrow f_t \text{の形より明らかに} \sum_{t=1}^T f_t(\mathbf{w}_t) \geq |\mathcal{M}| \quad R = \max_t \|\mathbf{x}_t\| \text{とおくと}$$

$$(60) \text{の右辺の最小値は} \eta = \frac{\|\mathbf{u}\|}{\sqrt{|\mathcal{M}|}R} \text{のときなので}$$

$$|\mathcal{M}| - \sum_{t=1}^T f_t(\mathbf{u}) \leq \sqrt{|\mathcal{M}|}R\|\mathbf{u}\| \leq \frac{1}{2\eta}\|\mathbf{u}\|^2 + \frac{\eta}{2}|\mathcal{M}|R^2 \quad (70)$$

$$\text{if } \exists \mathbf{u} \quad y_t \langle \mathbf{u}, \mathbf{x}_t \rangle \geq 1 \text{ then } \sum_{t=1}^T f_t(\mathbf{u}) = 0: \text{データが線形分離可能}$$

$$\Rightarrow |\mathcal{M}| \leq \sqrt{|\mathcal{M}|}R\|\mathbf{u}\| \Rightarrow |\mathcal{M}| \leq R^2\|\mathbf{u}\|^2 = \frac{R^2}{\gamma^2}$$

$$\Rightarrow \text{線形分離できる場合} \frac{1}{\|\mathbf{u}\|} \text{は境界面に最も近いデータと境界面の距離} \gamma$$

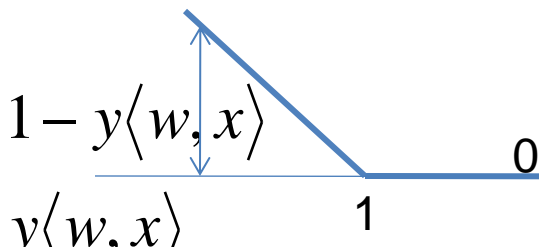


# Passive Aggressive Algorithm

- K.Crammer et.al. Online Passive-Aggressive Algorithms  
Journal of Machine Learning Research 7 (2006) 551–585
- 識別(あるいは分類)の問題設定
  - round  $t$  で  $n$ 次元データ  $x_t \in R^n$  が到着
  - $x_t$  の正負は  $y_t = \{+1: \text{正}, -1: \text{負}\}$  のように与えられる
  - 重みベクトル:  $w \in R^n \Rightarrow \text{sign}\langle w, x \rangle$ : 正負を表すので
  - 正しい(誤った)判定:  $y_t \langle w_t, x_t \rangle > 0$  、 ( $< 0$ )
- $w_t$  はデータが到着するたびに更新されている

# 損失関数による定式化

- 境界面そのもので判定はきわどいのでマージンを持たせる。マージンを1とした場合の損失関数(hinge-loss function)は以下の通り

$$\ell(\mathbf{w}; (\mathbf{x}, y)) = \begin{cases} 0 & y\langle \mathbf{w}, \mathbf{x} \rangle \geq 1 \\ 1 - y\langle \mathbf{w}, \mathbf{x} \rangle & \text{otherwise} \end{cases}$$


以下では  $\ell_t = \ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))$  と書く

- この設定で、round  $t$  の更新は次の条件付き最適化問題となる。

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \quad s.t. \quad \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) = 0 \quad (PA-1)$$

# FoReLとして見ると

FoReL

$$\begin{aligned}\forall t \quad \mathbf{w}_t &= \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} f(\mathbf{w})_i + R(\mathbf{w}) \quad S \text{は} \mathbf{w} \text{の取り得る範囲で凸} \\ &= \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} \ell(\mathbf{w}_i; (\mathbf{x}_i, y_i)) + \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|_2^2\end{aligned}$$

次ページのようにFoReLの定式化では $\eta$ に相当する $\tau$ は  
個別の $y_t, \mathbf{x}_t, \mathbf{w}_t$ に依存するため  
(60)のような簡単な解析ができない。

# 最適化問題(PA-1)を解く

➤ If  $\ell_t=0$  then  $\mathbf{w}_t$  minimizes  $\frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2$

$$\mathbf{w}_{t+1} = \mathbf{w}_t$$

Passive

➤ If  $\ell_t \neq 0$  then Lagrange 未定乗数法で解く。

$$L(\mathbf{w}, \tau) = \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + \tau(1 - y_t \langle \mathbf{w}, \mathbf{x}_t \rangle)$$

$$\frac{\partial L(\mathbf{w}, \tau)}{\partial \mathbf{w}} = \mathbf{w} - \mathbf{w}_t - \tau y_t \mathbf{x}_t = 0$$

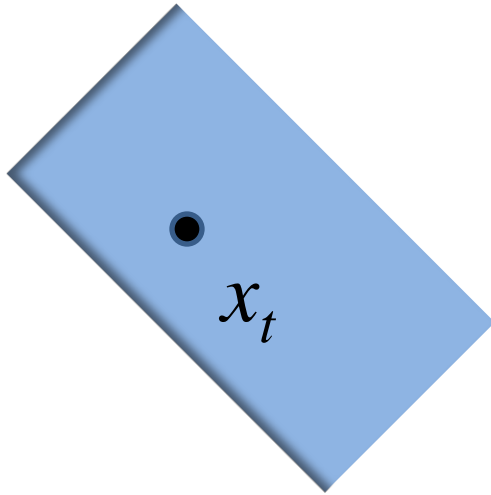
$$\Rightarrow \mathbf{w} = \mathbf{w}_t + \tau y_t \mathbf{x}_t \quad \Rightarrow L(\tau) = -\frac{1}{2} \tau^2 \|\mathbf{x}_t\|^2 + \tau(1 - y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle)$$

$$\Rightarrow \frac{\partial L(\tau)}{\partial \tau} = -\tau \|\mathbf{x}_t\|^2 + (1 - y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle) = 0$$

$$\Rightarrow \tau_t = \frac{1 - y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle}{\|\mathbf{x}_t\|^2} \quad \mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$$

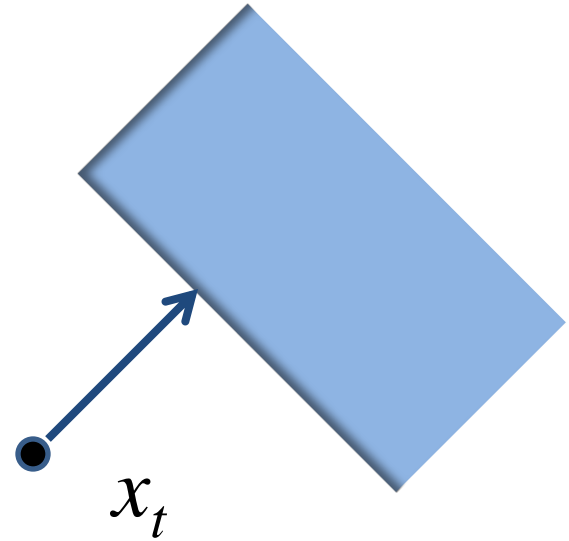
Aggressive

# Passive Aggressive



$$\mathbf{w}_{t+1} = \mathbf{w}_t$$

Passive



$$\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$$

Aggressive

$$\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2}$$

# soft marginの学習法 PA-I, PA-II

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in R^n} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi \quad s.t. \quad \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) \leq \xi, \quad \xi \geq 0 \quad (PA-I)$$

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in R^n} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi^2 \quad s.t. \quad \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) \leq \xi \quad (PA-II)$$

$\Rightarrow$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$$

$$\tau_t = \min \left\{ C, \frac{\ell_t}{\|\mathbf{x}_t\|^2} \right\} \quad (PA-I)$$

$$\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} \quad (PA-II)$$

# Passive Aggressive Algorithm

INPUT: aggressiveness parameter  $C > 0$

INITIALIZE:  $\mathbf{w}_1 = (\mathbf{0}, \dots, \mathbf{0})$

For  $t = 1, 2, \dots$

- receive instance:  $\mathbf{x}_t \in \mathbb{R}^n$
- predict:  $\hat{y}_t = \text{sign}\langle \mathbf{w}_t, \mathbf{x}_t \rangle$
- receive correct label:  $y_t \in \{-1, +1\}$
- suffer loss:  $\ell_t = \max\{0, 1 - y_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle\}$
- update:

1. set:  $\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2} \quad (\text{PA})$

$$\tau_t = \min \left\{ C, \frac{\ell_t}{\|\mathbf{x}_t\|^2} \right\} \quad (\text{PA-I})$$

$$\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} \quad (\text{PA-II})$$

2. update:  $\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t \langle \mathbf{w}_t, \mathbf{x}_t \rangle$

# 付録: PA-Iの導出

$PA-I$ のLagrangianは以下の通り

$$\begin{aligned} L(w, \xi, \tau, \lambda) &= \frac{1}{2} \|w - w_t\|^2 + C\xi + \tau(1 - \xi - y(w \cdot x_t)) - \lambda\xi \\ &= \frac{1}{2} \|w - w_t\|^2 + \xi(C - \tau - \lambda) + \tau(1 - y(w \cdot x_t)) \quad \tau \geq 0, \quad \lambda \geq 0 \quad (pa1-1) \end{aligned}$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = w_t + \tau y_t x_t$$

$\xi(C - \tau - \lambda)$ の最小値は0で  $C - \tau - \lambda = 0$  ( $pa1-2$ )のとき。

$C - \tau - \lambda \geq 0$ である。そうでないとする、 $\xi(C - \tau - \lambda)$ はいくらでも小さくなれるので、 $L$ の最小化ができない。

KKT条件より  $\lambda \geq 0$ なので、( $pa1-2$ )より  $C - \tau \geq 0 \Rightarrow C \geq \tau$  ( $pa1-3$ )

以下では  $C \geq \ell_t / \|x_t\|^2$  (case 1) と  $C < \ell_t / \|x_t\|^2$  (case 2) に分けて考える。

case 1

( $pa1-2$ )を( $pa1-1$ )に代入し

$$L(w, \xi, \tau, \lambda) = \frac{1}{2} \|w - w_t\|^2 + \tau(1 - y(w \cdot x_t)) \text{ となるので、これを } w \text{ について最適化すると}$$

$$\text{元々のPAと同じく} \quad \tau_t = \ell_t / \|x_t\|^2 \quad (pa1-4)$$



$$(\text{case 2}) \quad \ell_t / \|x_t\|^2 > C \quad \Rightarrow \quad C\|x_t\|^2 < 1 - y_t \langle w_t, x_t \rangle \quad (pa1-6)$$

元々の optimization

$$w_{t+1} = \arg \min_w \frac{1}{2} \|w - w_t\|^2 + C\xi \quad s.t \quad 1 - y_t \langle w, x_t \rangle \leq \xi \quad (pa1-7) \quad \text{and } \xi \geq 0$$

$$\text{と } w = w_t + \tau y_t x_t \quad \text{により } 1 - y_t \langle w_t, x_t \rangle - \tau \|x_t\|^2 \leq \xi \quad (pa1-8)$$

$$(pa1-8) \text{と} (pa1-6) \text{を組み合わせると} \quad C\|x_t\|^2 - \tau\|x_t\|^2 < \xi$$

$(pa1-3)$ で  $C \geq \tau$  だったから、  $0 < \xi$

$$KKT \text{条件から } \lambda \xi = 0 \text{なので、 } \lambda = 0 \quad \Rightarrow \quad (pa1-2) \text{より} \quad \tau = C$$

$$(\text{case 1}) (\text{case 2}) \text{を合せると} \quad \tau_t = \min \left\{ C, \frac{\ell_t}{\|x_t\|^2} \right\}$$

## 付録: PA-IIの導出

$\ell_t = 0 \Rightarrow \tau_t = 0 \Rightarrow \ell_t > 0$ の場合について考えればよい

$$L(w, \xi, \tau) = \frac{1}{2} \|w - w_t\|^2 + C\xi^2 + \tau(1 - \xi - y(\underline{w} \cdot x_t)) \quad \tau \geq 0 \quad (pa2-1)$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = w_t + \tau y_t x_t$$

$$\frac{\partial L}{\partial \xi} = 2C\xi - \tau = 0 \Rightarrow \xi = \tau/2C \Rightarrow \text{この}\xi\text{と}w\text{を}L\text{に代入すると}$$

$$L(w, \tau) = -\frac{\tau^2}{2} \left( \|x_t\|^2 + \frac{1}{2C} \right) + \tau(1 - y_t(\underline{w}_t \cdot x_t))$$

$$\frac{\partial L}{\partial \tau} = 0 \Rightarrow \tau = \frac{1 - y_t(w_t \cdot x_t)}{\|x_t\|^2 + \frac{1}{2C}} = \frac{\ell_t}{\|x_t\|^2 + \frac{1}{2C}} \quad \blacksquare$$

# 損失の限界の評価

任意の固定された重みベクトル  $u$  に対する損失を  $l_t^*$  とする。

$$l_t = l(w_t; (x_t, y_t)) \quad l_t^* = l(u; (x_t, y_t))$$

## Lemma 1

$(x_1, y_1), \dots, (x_T, y_T)$  はデータ列。ただし  $x_t \in R^n, y_t \in \{+1, -1\}$   
 $\tau_t$  は PA, PA - I, PA - II の Algorithm における更新式 parameter  
上記の  $u \in R^n$  に対して

$$\sum_{t=1}^T \tau_t \left( 2\ell_t - \tau_t \|x_t\|^2 - 2\ell_t^* \right) \leq \|u\|^2$$

# Proof

$$\Delta_t \stackrel{\text{def}}{=} \|w_t - u\|^2 - \|w_{t+1} - u\|^2$$

$$\sum_{t=1}^T \Delta_t = \sum_{t=1}^T \left( \|w_t - u\|^2 - \|w_{t+1} - u\|^2 \right) = \|w_1 - u\|^2 - \|w_{T+1} - u\|^2$$

$$w_1 = 0 \quad \|w_{T+1} - u\|^2 \geq 0 \quad \text{なので}$$

$$\sum_{t=1}^T \Delta_t \leq \|u\|^2$$

minimum marginをviolateしない場合は、 $\ell_t = 0, \tau_t = 0 \quad \Rightarrow \quad \Delta_t = 0$

よって、 $\ell_t > 0$ の場合にだけ注目： $w_{t+1} = w_t + y_t \tau_t x_t$

$$\begin{aligned} \Delta_t &= \|w_t - u\|^2 - \|w_{t+1} - u\|^2 = \|w_t - u\|^2 - \|w_t - u + y_t \tau_t x_t\|^2 \\ &= \|w_t - u\|^2 - \left( \|w_t - u\|^2 + 2\tau_t y_t \langle (w_t - u), x_t \rangle + \tau_t^2 \|x_t\|^2 \right) \\ &= -2\tau_t y_t \langle (w_t - u), x_t \rangle - \tau_t^2 \|x_t\|^2 \end{aligned}$$

$$\ell_t > 0 \text{ の場合 } \quad \ell_t = 1 - y_t \langle w_t, x_t \rangle \quad \Rightarrow \quad y_t \langle w_t, x_t \rangle = 1 - \ell_t$$

$$\text{同じく} \quad y_t \langle u, x_t \rangle = 1 - \ell_t^*$$

$$\begin{aligned} \Rightarrow \quad \Delta_t &= -2\tau_t y_t \langle (w_t - u), x_t \rangle - \tau_t^2 \|x_t\|^2 \\ &\geq 2\tau_t \left( (1 - \ell_t^*) - (1 - \ell_t) \right) - \tau_t^2 \|x_t\|^2 \\ &= \tau_t \left( 2\ell_t - \tau_t \|x_t\|^2 - 2\ell_t^* \right) \end{aligned}$$



## Theorem 2

Lemma 1 と同じ設定。  $\exists u \quad s.t. \quad \forall t \left[ \ell_t^* = 0 \right], \max_t \|x_t\| \leq R$

$$\Rightarrow \sum_{t=1}^T \ell_t^2 \leq \|u\|^2 R^2$$

### Proof

$$\forall t \left[ \ell_t^* = 0 \right] \text{ and Lemma 1} \quad \Rightarrow \quad \sum_{t=1}^T \tau_t \left( 2\ell_t - \tau_t \|x_t\|^2 \right) \leq \|u\|^2$$

$$\text{PA では} \quad \tau_t = \ell_t / \|x_t\|^2 \quad \Rightarrow \quad \sum_{t=1}^T \ell_t^2 / \|x_t\|^2 \leq \|u\|^2$$

$$\forall t \left[ \|x_t\|^2 \leq R^2 \right] \quad \Rightarrow \quad \sum_{t=1}^T \ell_t^2 / R^2 \leq \|u\|^2 \quad \Rightarrow \quad \sum_{t=1}^T \ell_t^2 \leq \|u\|^2 R^2 \quad \blacksquare$$

- Theorem 2では次の制約が厳しい。  $\exists u \quad s.t. \quad \forall t [l_t^* = 0]$
- この制約は、 $u$ で完全な識別ができること。
- この制約を外す定理を考える

Theorem 3

Lemma 1 と同じ設定。

$\forall t \quad \|x_t\|^2 = 1$       このとき  $\ell_t^* = \ell(u; (x_t, y_t))$ であるような  $u \in R^n$ に対して

$$\sum_{t=1}^T \ell_t^2 \leq \left( \|u\| + 2\sqrt{\sum_{t=1}^T (\ell_t^*)^2} \right)^2$$

Proofは次ページ

## Proof

$$\|x_t\|^2 = 1 \quad \Rightarrow \quad \tau_t = l_t$$

$$\therefore \sum_{t=1}^T \tau_t (2l_t - \tau_t \|x_t\|^2 - 2l_t^*) \leq \|u\|^2 \quad \Rightarrow \quad \sum_{t=1}^T l_t^2 \leq 2 \sum_{t=1}^T l_t l_t^* + \|u\|^2$$

Cauchy - Schwartz より  $\sum_{t=1}^T l_t l_t^* \leq \sqrt{\sum_{t=1}^T l_t^2} \cdot \sqrt{\sum_{t=1}^T (l_t^*)^2}$  だから

$$\sum_{t=1}^T l_t^2 \leq 2 \sqrt{\sum_{t=1}^T l_t^2} \cdot \sqrt{\sum_{t=1}^T (l_t^*)^2} + \|u\|^2 \quad \text{ここで } \sqrt{\sum_{t=1}^T l_t^2} = LT \text{ とおくと}$$

$$LT^2 - 2LT \cdot \sqrt{\sum_{t=1}^T (l_t^*)^2} - \|u\|^2 \leq 0$$

$LT$  の最大値  $\max LT$  は上式で  $\leq$  を  $=$  とした 2 次式の大きなほうの解。

$$\max LT = \sqrt{\sum_{t=1}^T (l_t^*)^2} + \sqrt{\sum_{t=1}^T (l_t^*)^2} + \|u\|^2$$

$$\sqrt{a+b} \leq \sqrt{a} + \sqrt{b} \text{ を使 う と } \max LT = \max \sqrt{\sum_{t=1}^T l_t^2} \leq 2 \sqrt{\sum_{t=1}^T (l_t^*)^2} + \|u\|$$

$$\Rightarrow \quad \max \sum_{t=1}^T l_t^2 \leq \left( 2 \sqrt{\sum_{t=1}^T (l_t^*)^2} + \|u\| \right)^2$$





## PA-Iにおける入力データ識別の失敗回数の上限

$$\|x_t\|^2 \leq R^2 \quad \tau_t = \min \left\{ \frac{\ell_t}{\|x_t\|^2}, C \right\}$$
$$\Rightarrow \#mistakes \leq \max \{R^2, 1/C\} \left( \|u\|^2 + 2C \sum_{t=1}^T \ell_t^* \right)$$

Proofは次のページ

$t$ 回目の繰り返しで *mistake* が起きたとすると  $\ell_t \geq 1$

$\|x_t\|^2 \leq R^2$  と  $\tau_t = \min\{\ell_t / \|x_t\|^2, C\}$  より  $\min\{1/R^2, C\} \leq \tau_t \ell_t$

$M$  を繰り返し全体における *mistake* 回数とする。

$$0 \leq \tau_t \ell_t \quad \text{なので} \quad \min\{1/R^2, C\} M \leq \sum_{t=1}^T \tau_t \ell_t \quad (\text{pa1-10})$$

$$\ell_t^* = \ell(u; (x_t, y_t)) \Rightarrow \tau_t \ell_t^* \leq C \ell_t^* \text{ かつ } \tau_t \|x_t\|^2 \leq \ell_t$$

$$\text{これをLemma 1} \quad \sum_{t=1}^T \tau_t (2\ell_t - \tau_t \|x_t\|^2 - 2\ell_t^*) \leq \|u\|^2 \text{ に代入すると}$$

$$\sum_{t=1}^T \tau_t \ell_t \leq \|u\|^2 + 2C \sum_{t=1}^T \ell_t^* \quad (\text{pa1-20})$$

$$(\text{pa1-10})(\text{pa1-20}) \quad \text{より} \quad \min\{1/R^2, C\} M \leq \|u\|^2 + 2C \sum_{t=1}^T \ell_t^* \quad (\text{pa1-30})$$

(pa1-30)の両辺に  $\max\{R^2, 1/C\}$  を掛けると

$$M \leq \max\{R^2, 1/C\} \left( \|u\|^2 + 2C \sum_{t=1}^T \ell_t^* \right) \quad \blacksquare$$

## PA-IIIにおける累積損失の上限

$$\begin{aligned} \|x_t\|^2 &\leq R^2 & \tau_t &= \frac{\ell_t}{\|x_t\|^2 + \frac{1}{2C}} \\ \Rightarrow \sum_{t=1}^T \ell_t^2 &\leq \left(R^2 + \frac{1}{2C}\right) \left(\|u\|^2 + 2C \sum_{t=1}^T \ell_t^*\right) \end{aligned}$$

Proofは次のページ

Lemma 1:  $\|u\|^2 \geq \sum_{t=1}^T \tau_t (2\ell_t - \tau_t \|x_t\|^2 - 2\ell_t^*)$  で右辺の  $\Sigma$  内で  $(\alpha\tau_t - \ell_t^*/\alpha)^2$  を差し引く

ただし  $\alpha = 1/\sqrt{2C}$

$$\|u\|^2 \geq \sum_{t=1}^T \tau_t \left( 2\ell_t - \tau_t \|x_t\|^2 - 2\ell_t^* - (\alpha\tau_t - \ell_t^*/\alpha)^2 \right)$$

$$= \sum_{t=1}^T \left( 2\tau_t \ell_t - \tau_t^2 \|x_t\|^2 - 2\tau_t \ell_t^* - \alpha^2 \tau_t^2 + 2\tau_t \ell_t^* - \ell_t^{*2}/\alpha^2 \right)$$

$$= \sum_{t=1}^T \left( 2\tau_t \ell_t - \tau_t^2 (\|x_t\|^2 + \alpha^2) - \ell_t^{*2}/\alpha^2 \right)$$

$\alpha = 1/\sqrt{2C}$  を代入すると

$$\|u\|^2 \geq \sum_{t=1}^T \left( 2\tau_t \ell_t - \tau_t^2 \left( \|x_t\|^2 + \frac{1}{2C} \right) - 2C \ell_t^{*2} \right)$$

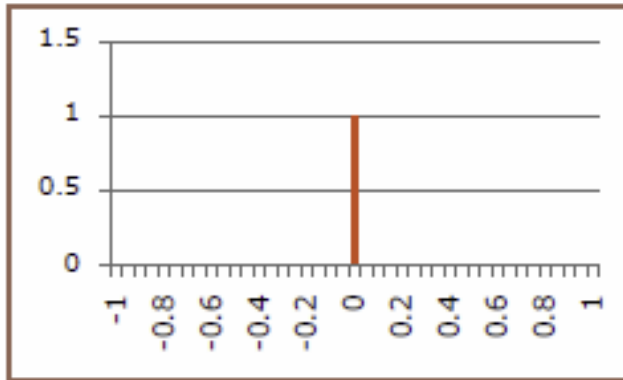
$$\tau_t = \ell_t / \left( \|x_t\|^2 + 1/(2C) \right) \text{ を代入すると } \|u\|^2 \geq \sum_{t=1}^T \left( \frac{\ell_t^2}{\|x_t\|^2 + \frac{1}{2C}} - 2C \ell_t^{*2} \right)$$

$$\|x_t\|^2 \leq R^2 \text{ を使えば、 } \sum_{t=1}^T \ell_t^2 \leq \left( R^2 + \frac{1}{2C} \right) \left( \|u\|^2 + 2C \sum_{t=1}^T \ell_t^* \right) \blacksquare$$

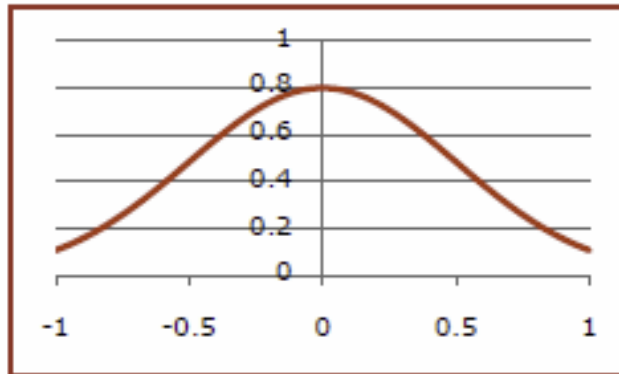
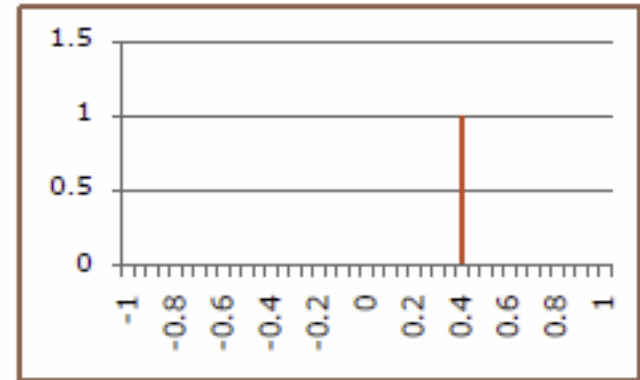
# Confidence Weighted Algorithm

Crammer et al. 2008

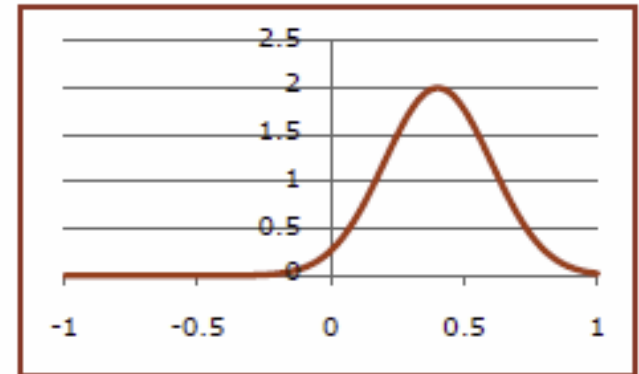
学習する重みベクトル $W$ を点ではなく分布(正規分布)にする  
→  $W$ の期待値と分散を更新する



既存手法



CW



# Pegasos:

## Primal Estimated sub-GrAdient Solver for SVM

### ➤ L2正則化+L1損失のSVM

$$\min_{\mathbf{w}} f(\mathbf{w}; B) \equiv \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{k} \sum_{i \in B} \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \quad A \text{は学習に使うデータ、} k = |A|$$

### ➤ Pegasosの更新は次式による

$$\mathbf{w}_{t+1/2} \leftarrow \mathbf{w}_t - \eta \nabla f(\mathbf{w}_t; A) \quad \nabla f(\mathbf{w}_t; A) \text{は} f \text{の劣微分} \partial f(\mathbf{w}_t; A) \text{の要素}$$

$$\text{where} \quad \nabla f(\mathbf{w}; A) = \lambda \mathbf{w}_t - \frac{1}{|A|} \sum_{i \in A^+} y_i \mathbf{x}_i \quad A^+ \equiv \{i \mid i \in A, 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0\},$$

$$\eta = \frac{1}{\lambda t}, \quad l \text{はベクトル} \mathbf{w}, \mathbf{x}_i \text{の次元、} t \text{は繰り返し回数}$$

### ➤ 更新の後、 $\mathbf{w}$ を半径 $\min(1, \sqrt{1/\lambda})$ の球にproject

$$\mathbf{w} \leftarrow \min(1, \sqrt{1/\lambda} / \|\mathbf{w}\|_2) \mathbf{w}$$

### ➤ 以上を収束するまで繰り返す。データ集合Aごとなので、onlineというよりはmini-batch

# Pegasos: Primal Estimated sub-GrAdient SOLver for SVM のアルゴリズム

初期化:  $\mathbf{w}_1 = 0$

For  $t = 1, 2, \dots, T$

全データ  $D$  から  $A_t$  を選ぶ。  $|A_t| = k$

$$A_t^+ = \{i \in A_t \mid 1 - y_i \langle \mathbf{w}_t, \mathbf{x}_i \rangle > 0\}$$

$$\eta^{(t)} = \frac{1}{\lambda t}$$

$$\mathbf{w}_{t+1/2} = (1 - \eta^{(t)} \lambda) \mathbf{w}_t + \frac{\eta^{(t)}}{|A_t^+|} \sum_{i \in A_t^+} y_i \mathbf{x}_i$$

$$\mathbf{w}_{t+1} = \min \left\{ 1, \frac{1/\sqrt{\lambda}}{\|\mathbf{w}_{t+1/2}\|} \right\} \mathbf{w}_{t+1/2}$$

Output:  $\mathbf{w}_{T+1}$



# $f(\mathbf{w})$ の評価

まず  $\mathbf{w}^* = \arg \min_{f(\mathbf{w})} f(\mathbf{w})$  とする。

また、関数 $f$ は、 $f(\mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2$ が凸のとき、 $\lambda$ -strongly convex という。

Lemma1.  $f_1, \dots, f_T$ を $\lambda$ -strongly convexとし、 $B$ を閉凸集合とする。

$\mathbf{w}_1, \dots, \mathbf{w}_{T+1}$ を以下のようなベクトルの列とする： $\nabla_t$ は $f_t$ の劣微分の要素

$\mathbf{w}_1 \in B, \quad \forall t \geq 1 \quad \mathbf{w}_{t+1} = \arg \min_{\mathbf{w}' \in B} \|\mathbf{w}_t - \eta_t \nabla_t(\mathbf{w}_t) - \mathbf{w}'\|$ ただし $\eta_t = 1/(\lambda t)$

$\|\nabla_t(\mathbf{w}_t)\| \leq G$ のとき、 $\forall \mathbf{u} \in B$ に対して次式が成り立つ

$$\frac{1}{T} \sum_{t=1}^T f_t(\mathbf{w}_t) \leq \frac{1}{T} \sum_{t=1}^T f_t(\mathbf{u}) + \frac{G^2(1 + \ln(T))}{2\lambda T}$$

Proof :  $f_t$ はstrongly - convexで、 $\nabla_t$ は劣微分の要素なので

$$\langle \mathbf{w}_t - \mathbf{u}, \nabla_t \rangle \geq f_t(\mathbf{w}_t) - f_t(\mathbf{u}) + \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \quad (10)$$

$\mathbf{w}'_t = \mathbf{w}_t - \eta_t \nabla_t$  とおき、 $\mathbf{w}_{t+1}$ は $\mathbf{w}'_t$ の $B$ へのprojectionであり

$\mathbf{u} \in B$  なので  $\|\mathbf{w}'_t - \mathbf{u}\|^2 \geq \|\mathbf{w}_{t+1} - \mathbf{u}\|^2$  である。よって

$$\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 \geq \|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}'_t - \mathbf{u}\|^2 = 2\eta_t \langle \mathbf{w}_t - \mathbf{u}, \nabla_t \rangle - \eta_t^2 \|\nabla_t\|^2 \quad (15)$$

$\|\nabla_t\| \leq G$ かつ $\eta_t = 1/(\lambda t)$ なので(15)により

$$\langle \mathbf{w}_t - \mathbf{u}, \nabla_t \rangle \leq \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta_t} + \frac{\eta_t}{2} G^2 \quad (20)$$

$$(10)と(20)を組み合わせると  $f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta_t} + \frac{\eta_t}{2} G^2 - \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \quad (30)$$$

(30)を $t=1, T$ まで総和をとると

$$\sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq \sum_{t=1}^T \left( \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta_t} - \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \right) + \sum_{t=1}^T \frac{\eta_t}{2} G^2 \quad (40)$$

$$\eta_t = 1/(\lambda t) \text{を代入すると} (40) = \sum_{t=1}^T \left( \lambda t \left( \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2} \right) - \frac{\lambda}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \right) + \frac{G^2}{2} \sum_{t=1}^T \frac{1}{\lambda t}$$

$$= \sum_{t=1}^T \lambda \left( (t-1) \|\mathbf{w}_t - \mathbf{u}\|^2 - t \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 \right) + \frac{G^2}{2\lambda} \sum_{t=1}^T \frac{1}{t} = -\frac{\lambda T}{2} \|\mathbf{w}_{T+1} - \mathbf{u}\|^2 + \frac{G^2}{2\lambda} \sum_{t=1}^T \frac{1}{t} \leq \frac{G^2}{2\lambda} (1 + \ln(T))$$

Lemma1を拡張するとさらに強力な次の定理が得られる。  
詳細は: Mathematical Programming 127(1), 2011, pp.3-30  
Pegasos: primal estimated sub-gradient solver for SVM. Shai  
Shalev-Schwartz, et.al.

Theorem 1:  $\forall \mathbf{x}$ (入力データ):  $\|\mathbf{x}\| \leq R$ ,  $\mathbf{w}^* = \arg \min_{f(\mathbf{w})} f(\mathbf{w})$

かつ  $\text{projection}$  したときは  $c = (\sqrt{\lambda} + R)^2$   
 $\text{projection}$  しないときは  $c = 4R^2$  とすると

$T \geq 3$  に対して

$$\frac{1}{T} \sum_{t=1}^T f(\mathbf{w}_t; A_t) \leq \frac{1}{T} \sum_{t=1}^T f(\mathbf{w}^*; A_t) + \frac{c(1 + \ln(T))}{2\lambda T}$$

$A_t$  は全データ  $D$  から選ばれた部分集合。

Proof :

if projectionが起こらない

then  $B$ は半径  $= 1/\sqrt{\lambda}$  の球で、 $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}' \in B} \|\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t; A_t) - \mathbf{w}'\|$

$\Rightarrow$  Theorem1を証明するにはLemma1の前提条件が成立することを証明する。

(1)  $f(\mathbf{w}, A_t)$ は、 $\lambda$  - strongly convex :  $\frac{\lambda}{2} \|\mathbf{w}\|^2$  と  $A_t$  のヒンジ損失の平均値の和なので

$\lambda$  - strongly convex

Proof : cont'd

(2)次に $\|\nabla_t\|$ の上界を求める。

if projection stepが実行された

$$\text{then } \|\mathbf{w}_t\| \leq 1/\sqrt{\lambda} \text{ かつ } \|\mathbf{x}\| \leq R \Rightarrow \nabla f(\mathbf{w}_t; A_t) = \lambda \mathbf{w}^{(t)} - \frac{1}{|A|} \sum_{i \in A^+} y_i \mathbf{x}_i \leq \sqrt{\lambda} + R$$

if projection stepが実行されなかった

$$\text{then } \mathbf{w}_{t+1} = (1 - \eta_t \lambda) \mathbf{w}_t + \frac{\eta_t}{|A_t|} \sum_{i \in A_t^+} y_i \mathbf{x}_i = \left(1 - \frac{1}{t}\right) \mathbf{w}_t + \frac{1}{\lambda t |A_t|} \sum_{i \in A_t^+} y_i \mathbf{x}_i$$

$$\text{where } A_t^+ = \{i \in A_t \mid 1 - y_i (\mathbf{w}^{(t)} \cdot \mathbf{x}_i) > 0\}$$

$$\text{ここで } \mathbf{w}_1 = 0 \Rightarrow \mathbf{w}_2 = \frac{1}{\lambda |A_1|} \sum_{i \in A_1^+} y_i \mathbf{x}_i \Rightarrow \mathbf{w}_3 = \left(\frac{1}{2}\right) \frac{1}{\lambda |A_1|} \sum_{i \in A_1^+} y_i \mathbf{x}_i + \frac{1}{2\lambda |A_2|} \sum_{i \in A_2^+} y_i \mathbf{x}_i$$

$$\Rightarrow \mathbf{w}_4 = \left(\frac{2}{3}\right) \left(\frac{1}{2}\right) \left( \frac{1}{\lambda |A_1|} \sum_{i \in A_1^+} y_i \mathbf{x}_i + \frac{1}{\lambda |A_2|} \sum_{i \in A_2^+} y_i \mathbf{x}_i \right) + \frac{1}{3\lambda |A_3|} \sum_{i \in A_3^+} y_i \mathbf{x}_i$$

$$\Rightarrow \mathbf{w}_t = \frac{1}{\lambda t} \sum_{j=1}^t \frac{1}{|A_j|} \sum_{i \in A_j^+} y_i \mathbf{x}_i$$

$$\Rightarrow \|\mathbf{w}_{t+1}\| \leq \frac{R}{\lambda}$$

この導出は初等的が  
だちょっとした計算

Proof : cont'd

(3)最後に $\mathbf{w}^* \in B$ を示す。

projectionしない場合は、 $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}' \in B} \|\mathbf{w}_t - \eta_t \nabla_t(\mathbf{w}_t) - \mathbf{w}'\|$ により

$\mathbf{w}^* \in B$ と言える。

projectionした場合は、 $\|\mathbf{w}^*\| \leq 1/\sqrt{\lambda}$ を示す。

ここで対象にしているSVMの主問題は以下の形式

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \quad \text{s.t. } \forall i \in [1, m]: \xi_i \geq 0, \xi_i \geq 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i)$$

$C = 1/(\lambda m)$ とおき、以下の双対問題を導く。

ここで双対問題を思い  
つくところがいかにも  
SVM的

$$\min \sum_{i=1}^m \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \right\|^2 \quad \text{s.t. } \forall i \in [1, m]: 0 \leq \alpha_i \leq C \quad (50)$$

主問題の解を $(\mathbf{w}^*, \xi^*)$ 、双対問題の解を $\alpha^*$ とすると

Proof : cont'd

主問題の解を $(\mathbf{w}^*, \xi^*)$  双対問題の解を $\alpha^*$ とすると

$$\mathbf{w}^* = \sum_{i=1}^m \alpha_i^* y_i \mathbf{x}_i \text{であり(50)は次のように書き直せる。} \|\alpha^*\|_1 - \frac{1}{2} \|\mathbf{w}^*\|^2$$

SVMの問題では強双対定理が成り立つので

$$\frac{1}{2} \|\mathbf{w}^*\|^2 + C \|\xi^*\|_1 = \|\alpha^*\|_1 - \frac{1}{2} \|\mathbf{w}^*\|^2$$

強双対定理の実に  
賢い使い方だ

$$\|\alpha^*\|_\infty = \max \{ \alpha_i \in \alpha^* \} \leq C = 1/(\lambda m) \Rightarrow \|\alpha^*\|_1 = \sum_{i=1}^m |\alpha_i| \leq Cm = 1/\lambda$$

$$\Rightarrow \frac{1}{2} \|\mathbf{w}^*\|^2 \leq \frac{1}{2} \|\mathbf{w}^*\|^2 + C \|\xi^*\|_1 = \|\alpha^*\|_1 - \frac{1}{2} \|\mathbf{w}^*\|^2 \leq \frac{1}{\lambda} - \frac{1}{2} \|\mathbf{w}^*\|^2$$

$$\Rightarrow \|\mathbf{w}^*\| \leq 1/\sqrt{\lambda}$$

$\Rightarrow$  以上の結果(1)(2)(3)をLemma1に適用すれば定理が得られる。  $\square$

# Coordinate Descent

C.-J. Hsieh, et.al. ICML2008

- Target:  $L1$  損失-  $L2$  正則化のSVM 双対化して解く。下に定義

$$\min_{\alpha} f^D(\alpha) \equiv \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha$$

subject to  $0 \leq \alpha_i \leq C \quad \forall i = 1, \dots, l$

where  $Q_{ij} \equiv y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$



# Coordinate Descent

- Coordinate Descentは順番に1変数ずつ選び、他の変数は固定して最適化。

$$\min_d f^D(\mathbf{a} + de_i) - f^D(\mathbf{a}) \equiv \frac{1}{2} Q_{ii} d^2 + \nabla_i f^D(\mathbf{a}) d \quad f^D \text{は} f \text{の双対}$$

$$\text{subject to } 0 \leq \alpha_i + d \leq C \quad \text{where } e_i = \left[ \underbrace{0, \dots, 0}_{i-1}, 1, 0, \dots, 0 \right]^T$$

$d$ を計算することによる $\alpha_i$ の更新式は下式

$$\alpha_i \leftarrow \min \left( \max \left( \alpha_i - \frac{\nabla_i f^D(\mathbf{a})}{Q_{ii}}, 0 \right), C \right) \quad (CD10)$$

# Coordinate Descent つづき

➤ (CD10)の $Q_{ii}$ は $\alpha_i$ の最適化の中で1回計算すればよいが

➤  $\nabla_i f^D(\boldsymbol{\alpha}) = (Q\boldsymbol{\alpha})_i - 1 = \sum_{t=1}^l (y_i y_t \langle \mathbf{x}_i, \mathbf{x}_t \rangle) \alpha_t - 1$  は  $\langle \mathbf{x}_i, \mathbf{x}_t \rangle \forall t = 1, \dots, l$

の計算コストが $O(nl)$ でうれしくない。そこで

➤  $\mathbf{u} \equiv \sum_{t=1}^l y_t \alpha_t \mathbf{x}_t$  を保持しておけば

➤  $\nabla_i f^D(\boldsymbol{\alpha}) = (Q\boldsymbol{\alpha})_i - 1 = y_i \langle \mathbf{u}, \mathbf{x}_i \rangle - 1$  (CD20)

となり計算コストは 以下の計算のための $O(n)$ でうれしい。

$$\mathbf{u} \leftarrow \mathbf{u} + y_i (\alpha_i - \bar{\alpha}_i) \mathbf{x}_i \quad (\text{CD30})$$

$\bar{\alpha}_i$  (CD10)の更新前、 $\alpha_i$  (CD10)の更新後

# L1損失-L2正則化のSVMの Coordinate Descent アルゴリズム

$\alpha$ の初期化、および  $\mathbf{u} = \sum_{i=1}^l y_i \alpha_i \mathbf{x}_i$

$Q_{ii} \ \forall i = 1, \dots, l$  の計算

while  $\alpha$  is not optimal

For  $i = 1, \dots, l$

(CD20)により  $G = y_i \langle \mathbf{u}, \mathbf{x}_i \rangle - 1$  の計算

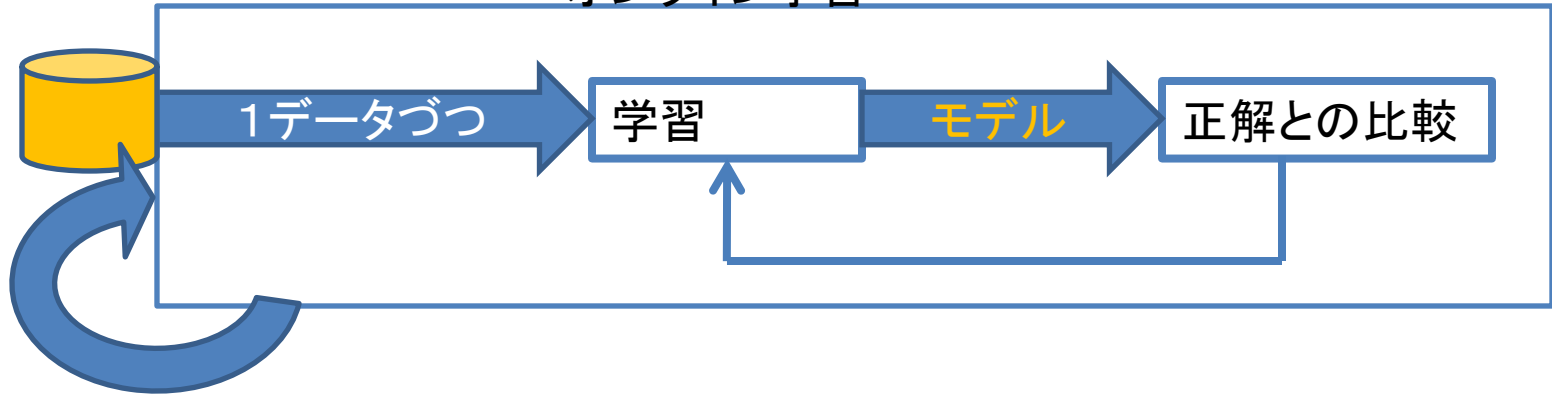
$\bar{\alpha}_i \leftarrow \alpha_i$

$\alpha_i \leftarrow \min \left( \max \left( \alpha_i - \frac{G}{Q_{ii}}, 0 \right), C \right)$

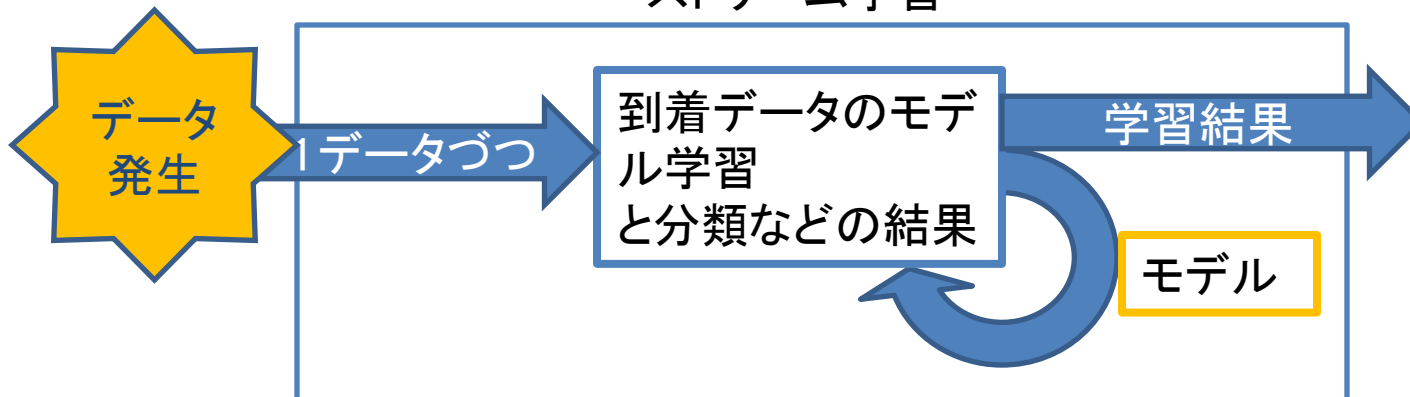
$\mathbf{u} \leftarrow \mathbf{u} + y_i (\alpha_i - \bar{\alpha}_i) \mathbf{x}_i$

# オンライン学習とストリーム学習

## オンライン学習



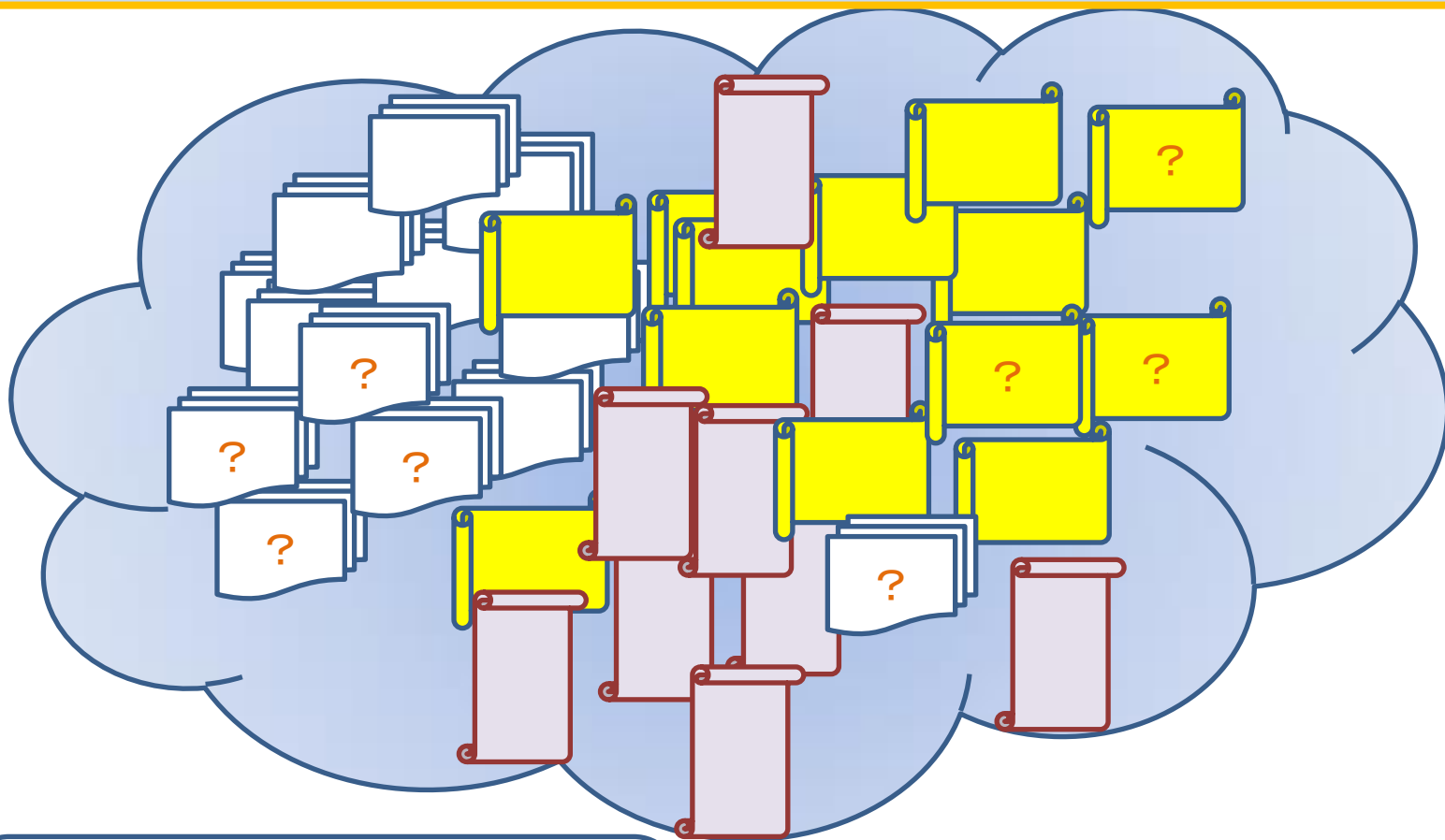
## ストリーム学習



# バッチ、オンライン、ストリームの比較

	バッチ学習	オンライン学習	ストリーム学習
メモリに乗せるデータ	同時に全部	同時には1データ	同時には1データ
メモリ量	大	小でも可能	小
データの到来	全データが揃ってから処理	全データが揃ってから処理	1データ到着ごとに処理
データの消去	消去せず	消去せず	データは処理後に消去
同一データの処理回数	収束するまで繰り返し	収束するまで繰り返し	1回
メモリに保持するモノ	全データと途中で内部状態	内部状態のみでも可能	内部状態のみ
性能	精度高	バッチより劣る。 ただし、最近 はバッチに肉迫	劣る
可能な処理	何でもあり	やや制限あり	限定的

捕捉：世の中、ビッグデータがホットだと言うけれど



異なる分類のデータ



分類されていない生のデータ

# パーセプトロンの別のアルゴリズム

データ  $x_i$  は  $N$  個ある。

$w(0) = 0; k = 0; R = \max_{1 \leq i \leq N} \|x_i\|;$

*repeat*

*for*       $i = 1$       *to*       $N$

{ *if*  $y_i \langle w(k), x_i \rangle \leq 0$  *then*

$\{ w(k+1) = w(k) + \eta y_i x_i; \quad k = k + 1 \ ; \}$

}

*until* {      *for* ループ内で失敗しない  
          ように  $w(k)$  を最適化の結果とする。  
          (すなわち  $y_i \langle w(k), x_i \rangle \leq 0$  の場合なし)

この部分が線形識別

$y_i \langle w(k), x_i \rangle \leq 0$  という識別に失敗したデータに、その値を  
重み (学習率と呼ぶ)  $\eta$  で  $w$  に足しこんで是正を図るアルゴリズム

# パーセプトロンは有限回で収束

→ mistakeのupper bound

## Novikoffの定理(バイアスのない場合)

$$R = \max_{1 \leq i \leq N} \|x_i\| \quad (0)$$

$$y_i \langle w_{opt}, x_i \rangle \geq \gamma : \text{マージン} \quad (1)$$

である $w_{opt}$ が存在するなら、パーセプトロン  
アルゴリズムが失敗する回数はたかだか

$$\left( \frac{R}{\gamma} \right)^2 \|w_{opt}\|^2 \text{回である}$$



# 証明

$t$ 回目の失敗 に先立つ重みを $w_{t-1}$

更新は、 $y_i \langle w_{t-1}, x_i \rangle < 0$  のとき起こる。このとき  $w_t = w_{t-1} + \eta y_i x_i$  (2)

$$(1)より \quad \langle w_t, w_{opt} \rangle = \langle w_{t-1}, w_{opt} \rangle + \eta y_i \langle x_i, w_{opt} \rangle \geq \langle w_{t-1}, w_{opt} \rangle + \eta \gamma \quad (3)$$

$$w_0 = 0とすれば(3)を繰り返し用いて  $\langle w_t, w_{opt} \rangle \geq t \eta \gamma \quad (4)$$$

$$(2)より \quad \|w_t\|^2 = \|w_{t-1}\|^2 + 2\eta y_i \langle w_{t-1}, x_i \rangle + \eta^2 \|x_i\|^2$$

←  $x_i$ は負例なので第2項は負

$$\leq \|w_{t-1}\|^2 + \eta^2 \|x_i\|^2 \leq \|\hat{w}_{t-1}\|^2 + \eta^2 R^2$$

$$\Rightarrow \quad \|w_t\|^2 \leq t \eta^2 R^2 \quad \Rightarrow \quad \|w_t\| \leq \sqrt{t} \eta R \quad (5)$$

$$(4)(5)より \quad \|w_{opt}\| \sqrt{t} \eta R \geq \|w_{opt}\| \|w_t\| \geq \langle w_t, w_{opt} \rangle \geq t \eta \gamma$$

$$\Rightarrow \quad t \leq \left( \frac{R}{\gamma} \right)^2 \|w_{opt}\|^2 \quad \blacksquare$$

# メモリ容量より大きなデータのSVM

Hsiang-Fu Yu et.al KDD2010

主問題の場合はPegasos, 双対問題の場合は  
Coordinate Descent (CD) をブロックごとに適用

1. 全データインデクス  $\{1, \dots, l\}$  を  $m$  ブロック  $B_1, \dots, B_m$  に分割
2. 主問題なら  $w$ , 双対問題なら  $\alpha$  を初期化
3. For  $k=1, 2, \dots$  (外側の繰り返し)
4.   For  $j=1, \dots, m$  (内側の繰り返し)
5.     read  $x_r \ \forall r \in B_j$  from **Disk**
6.      $\{x_r \mid r \in B_j\}$  に関して最適化(PegasosかCD)を行う
7.      $w$ あるいは $\alpha$ を更新

ここが重要

# 主問題をPegasosで解く場合の6.の部分

Pegasosでは次の最適化をする

$$\min_{\mathbf{w}} \frac{1}{2lC} \|\mathbf{w}\|^2 + \frac{1}{l} \sum_{i=1}^l \max(1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle, 0)$$

$\mathbf{B}_j$ を $B_j^1, \dots, B_j^{\bar{r}}$ に分割

For  $r = 1, \dots, \bar{r}$

$$\bar{\mathbf{w}} = \mathbf{w} - \eta^t \nabla^t$$

$$\text{where } \eta^t = \frac{lC}{t} \quad \nabla^t = \frac{1}{lC} \mathbf{w} + \frac{1}{|B|} \sum_{i \in B^+}^l y_i \mathbf{x}_i, \quad B^+ = \{i \in B \mid 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle > 0\}$$

$$\mathbf{w} \leftarrow \min \left( 1, \frac{\sqrt{lC}}{\|\bar{\mathbf{w}}\|} \right) \bar{\mathbf{w}}$$

$$t \leftarrow t + 1$$

end For

# 双対問題でCoordinate Descent(CD)を使う場合

$$\min_{\mathbf{d}_{\mathbf{B}j}} f^D(\boldsymbol{\alpha} + \mathbf{d}_{\mathbf{B}j}) \quad f^D \text{は主問題} f \text{の双対問題} \quad (\text{CD10})$$

$$\text{subject to } \mathbf{d}_{\bar{\mathbf{B}}j} = 0 \text{ and } 0 \leq \alpha_i + d_i \leq C \quad \forall i \in B_j$$

$$\text{where } Q_{ij} \equiv y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad \bar{\mathbf{B}}j = \{1, \dots, l\} \setminus \mathbf{B}j$$

ここで、 $Q$ のうちメモリにいるブロック $\mathbf{B}j$ に関する部分だけを使い、下の最適化を行う ( $i=1, \dots, l$ に注意)

$$\min_{\mathbf{d}_{\mathbf{B}j}} f^D(\boldsymbol{\alpha} + \mathbf{d}_{\mathbf{B}j}) \equiv \frac{1}{2} \mathbf{d}_{\mathbf{B}j}^T Q_{\mathbf{B}j\mathbf{B}j} \mathbf{d}_{\mathbf{B}j} + \boldsymbol{\alpha}^T Q_{\mathbf{B}j,i} \mathbf{d}_{\mathbf{B}j} - \mathbf{e}_{\mathbf{B}j}^T \mathbf{d}_{\mathbf{B}j} - f^D(\boldsymbol{\alpha}) \quad (\text{CD20})$$

$$\Rightarrow 6.\text{の}\boldsymbol{\alpha}\text{更新部分} : \boldsymbol{\alpha}_{\mathbf{B}j} \leftarrow \boldsymbol{\alpha}_{\mathbf{B}j} + \arg \min_{\mathbf{d}_{\mathbf{B}j}} f^D(\boldsymbol{\alpha} + \mathbf{d}_{\mathbf{B}j})$$

さて、 $\mathbf{w} = \sum_{i=1}^l \alpha_i y_i \mathbf{x}_i$  をメモリ中に保持しておけば

$$\boldsymbol{\alpha}^T Q_{r,i} \mathbf{d}_r = y_r \langle \mathbf{w}, \mathbf{x}_r \rangle, \forall r \in \mathbf{B}j \quad \Rightarrow \text{最適化に必要なのはブロック}\mathbf{B}j\text{だけ。}$$

$$\text{なお、} \boxed{\mathbf{w} \leftarrow \mathbf{w} + \sum_{r \in \mathbf{B}j} d_r y_r \mathbf{x}_r} \quad (\text{CD30})$$

という更新で6.の更新部分で $\mathbf{w}$ は更新すればよい。

# 双対化の御利益： 教師データアクセスの観点から

- 主問題と双対問題は最適化するパラメーター数が違う。
  - 主問題パラメーター数  $\gg$  双対問題パラメーター数 なら双対問題を解くほうが楽 → 教科書的
- SVMの場合：
  - 主問題のパラメーターは重みベクトル:  $w$
  - 双対問題にパラメーターは個別データ:  $x_i$
  - → 必ずしも教科書的なお得感ではない。

# 双対化の御利益

## ➤ SVMの場合：

➤ 主問題のパラメターは重みベクトル:  $\mathbf{w}$

➤ 下の定式化なので、全教師データ  $\{t_n, \mathbf{x}_n\}$  が同時に必要

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{subject to} \quad t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) \geq 1 \quad n = 1, \dots, N \quad \dots (SVM 30)$$

➤ データ量が大きくメモリにロード仕切れない場合に困ったことになる。

➤ データ量は最近、増加傾向

# 双対化の御利益

- →必ずしも教科書的なお得感ではない。
- 一方、双対問題では入力データ $\mathbf{x}_i, t_i$ のと最適化する $a_i$ が対応する形で最適化式に現れるので、どのデータを学習で使うか制御しやすい。(下の式参照)
  - 例えば、 $a_i (i \neq j)$ を固定して、 $a_j$ を最適化する操作を $j$ を動かして繰り返すなど。そのときには  $k(\mathbf{x}_i, \mathbf{x}_j) \quad j = 1, \dots, N$  だけしか使わない。

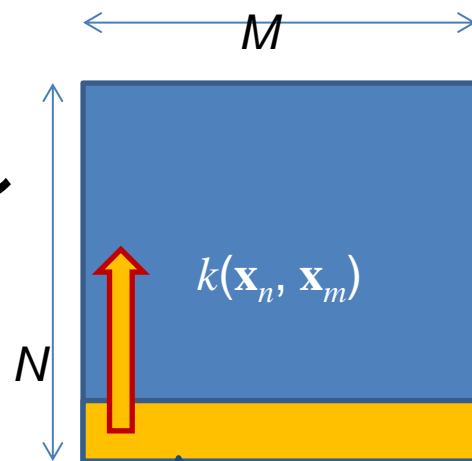
$$\max \tilde{L}(\mathbf{a}) = \max \left[ \sum_{n=1}^N a_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m) \right] \quad \cdots (SVM 70)$$

$$\text{subject to} \quad a_n \geq 0 \quad n = 1, \dots, N \quad \cdots (SVM 80)$$

$$\sum_{n=1}^N a_n t_n = 0 \quad \cdots (SVM 90) \quad \text{where} \quad k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$$

# 双対化の御利益

- 入力データ、あるいはカーネル行列全体がメモリに乗り切らないビッグデータを扱うために、入力（すなわちカーネル  $k(\mathbf{x}_n, \mathbf{x}_m)$ ）の一部を取捨選択してメモリにロードして使う方法が、この双対化で可能になっている。



- →ビッグデータ時代における御利益

- cf. 台湾大学のLIBSVM（SVMの実装）
- 全データからどのようにメモリにロードする部分を切り出すかについてはここで紹介した通り。

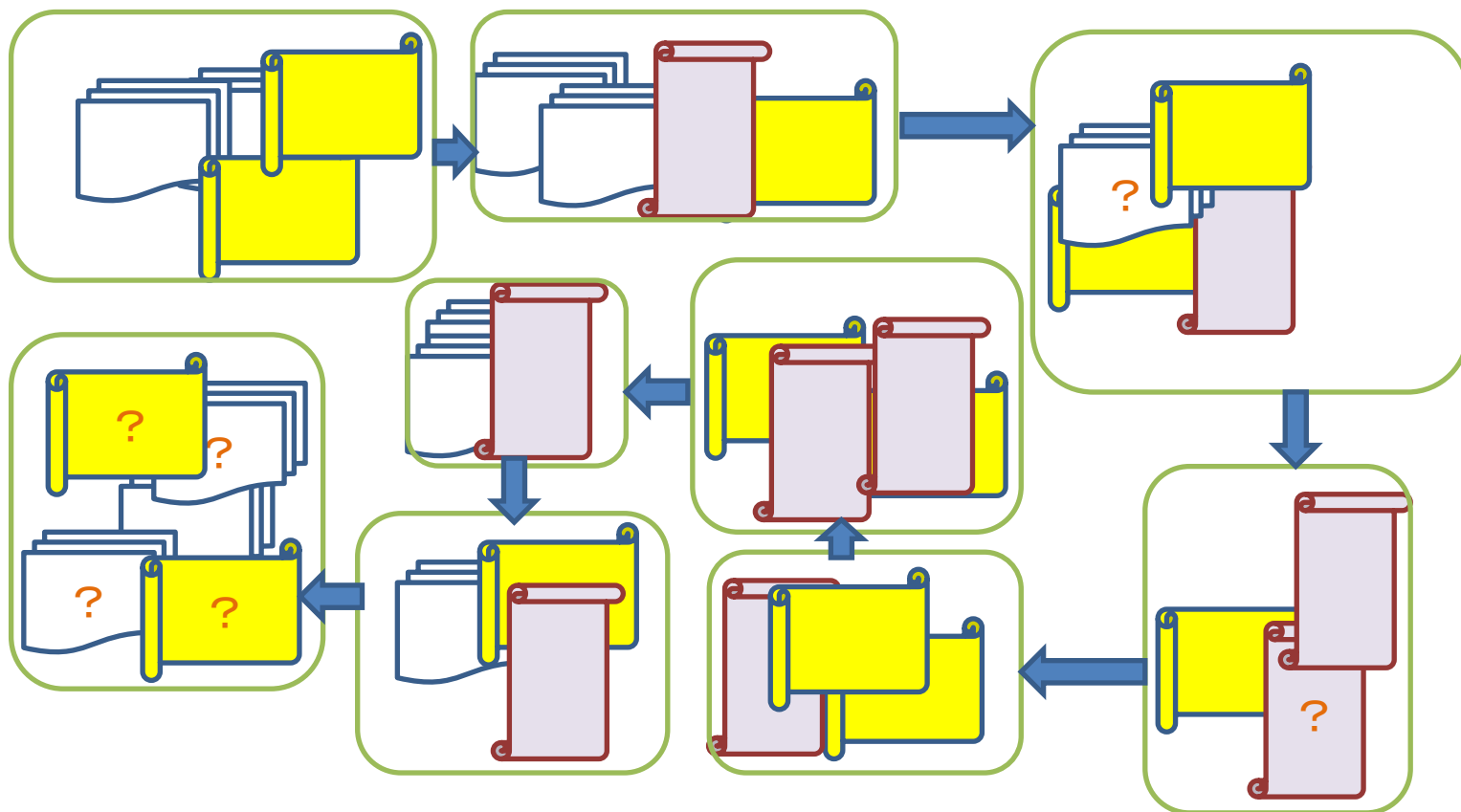
この部分だけ  
使って最適化：  
次に使う部分  
ロードし直して最  
適化：繰り返す



# 内外のバランス など

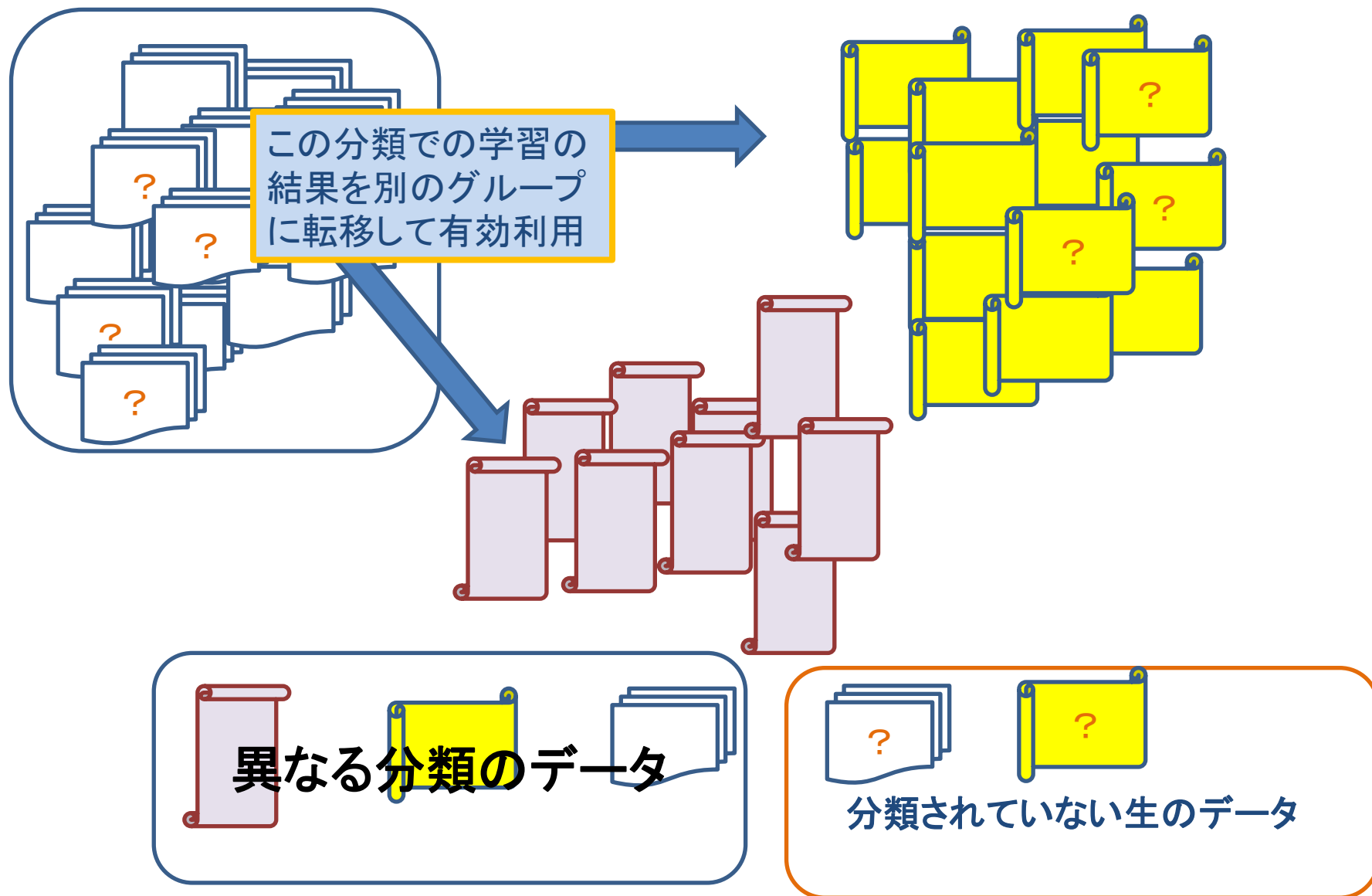
- 内側の繰り返しで最適化でCDにおいて $\alpha$ の更新を1回にし、looseな解を求めると、外側の繰り返しが多数回必要
- 内側の繰り返しで精密な最適化を行えば、外側の繰り返しは少なくてよい。
- $B_j\{j=1,\dots,m\}$ 内の要素の最適化処理における順番は外側の繰り返し毎にランダムに変更した方が収束が早い

# 小さなブロックに分けてデータマイニング、機械学習



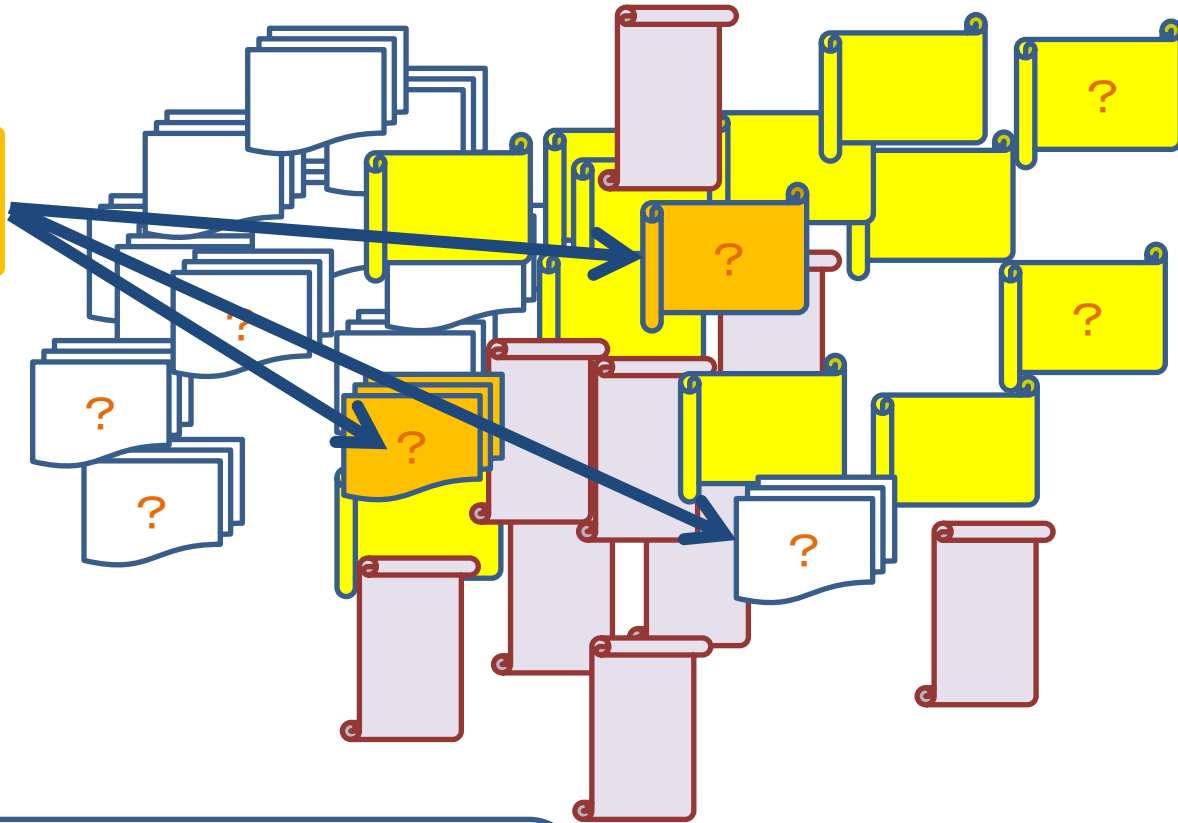
ブロックをメモリに順次ロードして学習し、その結果を活用して次のブロックへと繰り返す：  
例えば Stream SVM

# 転移学習

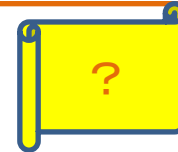


# 人間に正解をつけてもらうデータを絞り込む: Active学習

分類しにくい部分  
のデータ



異なる分類のデータ



分類されていない生のデータ

# 付録： DualityによるFoReLの定式化

Fenchel 双対 :  $f^*(\boldsymbol{\theta}) = \max_{\mathbf{u}} \langle \mathbf{u}, \boldsymbol{\theta} \rangle - f(\mathbf{u})$

明らかに次式が成立 : Fenchel - Young Equality :

$$\forall \mathbf{u}, \quad f^*(\boldsymbol{\theta}) \geq \langle \mathbf{u}, \boldsymbol{\theta} \rangle - f(\mathbf{u})$$

$\mathbf{w}_t = \arg \max_{\mathbf{u}} \langle \mathbf{u}, \mathbf{z}_t \rangle - f(\mathbf{u})$ だとすると

$$\langle \mathbf{w}_t, \mathbf{z}_t \rangle - f(\mathbf{w}_t) \geq \langle \mathbf{u}, \mathbf{z}_t \rangle - f(\mathbf{u})$$

$$\Rightarrow f(\mathbf{u}) - f(\mathbf{w}_t) \geq \langle \mathbf{u} - \mathbf{w}_t, \mathbf{z}_t \rangle$$

$\mathbf{z}_t$ は $f$ の $\mathbf{w}_t$ における sub - gradient

$f$ が微分可能なら、等式が成立するときは $f^*$ の定義より  $\mathbf{z}_t = \nabla f(\mathbf{w}_t)$

# Online Mirror Descent: OMD

FoReLで  $f_t(\mathbf{w}) = \langle \mathbf{w}, \mathbf{z}_t \rangle + R(\mathbf{w})$        $R(\mathbf{w})$ は正則化関数とする  
なお、 $\mathbf{w} \notin S$ だと  $R(\mathbf{w}) = \infty$  とする

$\mathbf{z}_{1:t} = \sum_{i=1}^t \mathbf{z}_i$  と略記するとFoReLは

$$\begin{aligned}\mathbf{w}_{t+1} &= \arg \min_{\mathbf{w}} R(\mathbf{w}) + \sum_{i=1}^t \langle \mathbf{w}, \mathbf{z}_i \rangle = \arg \min_{\mathbf{w}} R(\mathbf{w}) + \langle \mathbf{w}, \mathbf{z}_{1:t} \rangle \\ &= \arg \max_{\mathbf{w}} \langle \mathbf{w}, -\mathbf{z}_{1:t} \rangle - R(\mathbf{w})\end{aligned}$$

$g(\boldsymbol{\theta}) = \arg \max_{\mathbf{w}} \langle \mathbf{w}, \boldsymbol{\theta} \rangle - R(\mathbf{w})$  とおく とFoReLは

1.  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{z}_t$
2.  $\mathbf{w}_{t+1} = g(\boldsymbol{\theta}_{t+1})$

## Online Mirror Descent (OMD)

$$g: R^d \rightarrow S$$

$$\boldsymbol{\theta}_1 = \mathbf{0}$$

for  $t = 1, 2, \dots$

$$\mathbf{w}_t = g(\boldsymbol{\theta}_t) (= \arg \max_{\mathbf{w}} \langle \mathbf{w}, \boldsymbol{\theta} \rangle - R(\mathbf{w}))$$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \mathbf{z}_t \quad \text{where } \mathbf{z}_t \in \partial f_t(\mathbf{w}_t)$$

ここで以下が言える

$$\text{Regret}_T(\mathbf{u}) = \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq R(\mathbf{u}) - \min_{\mathbf{v} \in S} R(\mathbf{v}) + \eta \sum_{t=1}^T \|\mathbf{z}_t\|_*^2 \quad (\text{OMD100})$$

$$R(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \Rightarrow R^* = \frac{1}{2} \|\boldsymbol{\theta}\|^2$$

# パーセプトロンの別のアルゴリズム

データ  $x_i$  は  $N$  個ある。

$w(0) = 0; b(0) = 0; k = 0; R = \max_{1 \leq i \leq N} \|x_i\|;$

repeat

for  $i = 1$  to  $N$

{ if  $y_i (w(k)^T x_i + b(k)) \leq 0$  then

$\{ w(k+1) = w(k) + \eta y_i x_i; \quad b(k+1) = b(k) + \eta y_i R^2; \quad k = k+1 \quad ; \}$

}

until { forループ内で失敗しない  
(すなわち  $y_i (w(k)^T x_i + b(k)) \leq 0$  の場合なし)  
ように  $w(k), b(k)$  を最適化の結果とする。

この部分が線形識別

$y_i (w(k)^T x_i + b(k)) \leq 0$  という識別に失敗したデータに、その値を重み(学習率と呼ぶ)  $\eta$  で  $w$  に足しこんで是正を図るアルゴリズム



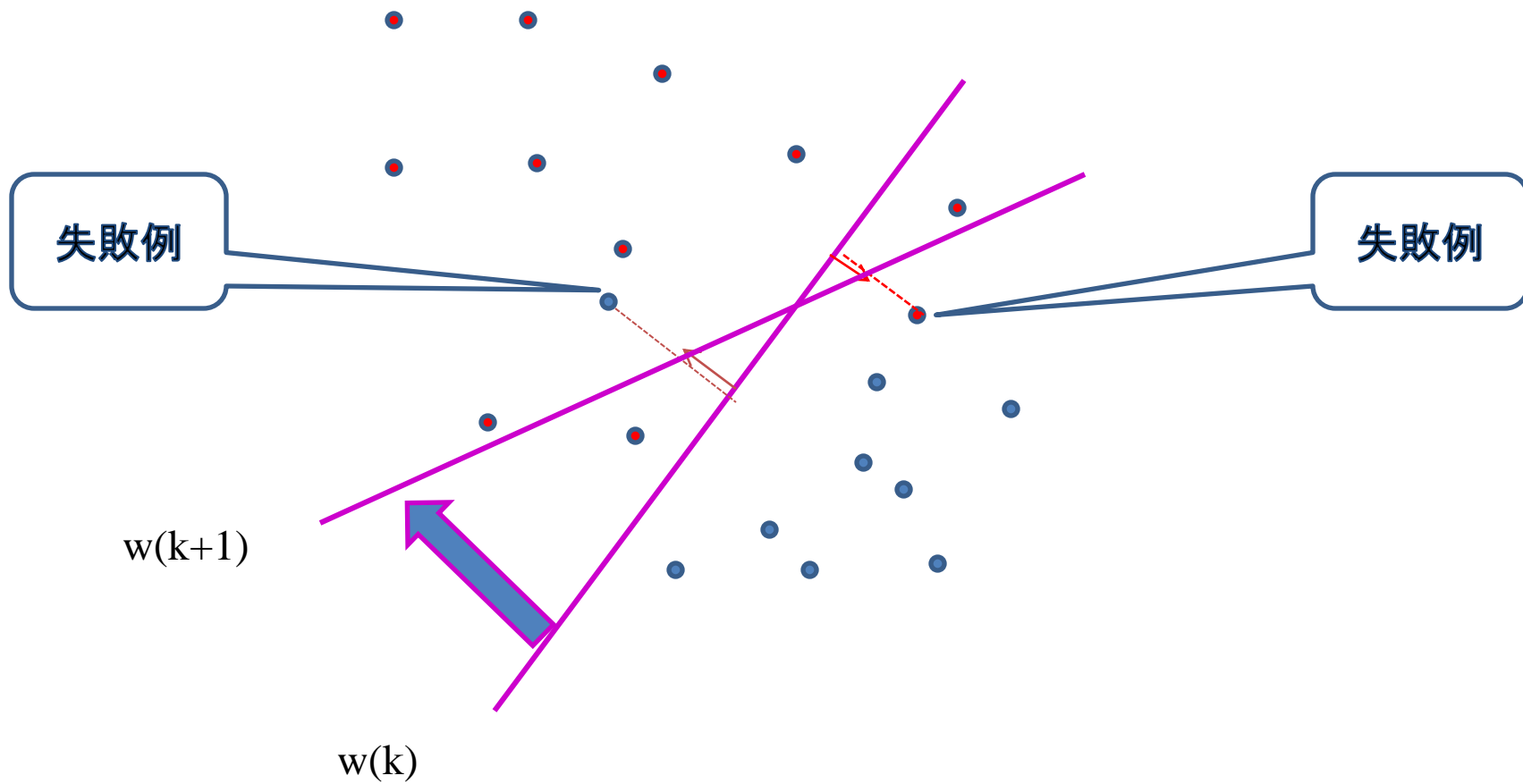
更新後の判定式を試みる

$$\begin{aligned} & y_i(w(k+1)^T x_i + b(k+1)) \\ &= y_i(w(k)^T x_i + b(k)) + y_i(\eta y_i x_i + \eta y_i R^2) \\ &= y_i(w(k)^T x_i + b(k)) + \eta y_i^2 (x_i^2 + R^2) \\ &= y_i(w(k)^T x_i + b(k)) + \eta y_i^2 (x_i^2 + \max \|x_i\|^2) \end{aligned}$$

第2項は必ず正

よって、
$$y_i(w(k+1)^T x_i + b(k+1)) > y_i(w(k)^T x_i + b(k))$$

したがって、失敗しない方向に $w$ は改善していく。



# パーセプトロンは有限回で収束

→ mistakeのupper bound

## Novikoffの定理(バイアスのある場合)

$$R = \max_{1 \leq i \leq N} \|x_i\| \quad (0)$$

$$\|w_{opt}\| = 1 \quad \text{かつ} \quad y_i(w_{opt}^T x_i + b_{opt}) \geq \gamma : \text{マージン} \quad (1)$$

である $w_{opt}$ が存在するなら、パーセプトロン

アルゴリズムが失敗する回数はたかだか

$$\left(\frac{2R}{\gamma}\right)^2 \text{回である}$$

# 証明

入力ベクトルに値 $R$ となる1次元加え、 $\hat{x}_i = (x_i^T, R)^T$ とする。

$\hat{w}_i = \left( w_i^T, \frac{b_i}{R} \right)^T$ とする。

$t$ 回目の失敗 に先立つ重みを $\hat{w}_{t-1}$

更新は、 $y_i(\hat{w}_{t-1} \cdot \hat{x}_i) = y_i(w_{t-1} \cdot x_i + b_{t-1}) \leq 0$

のとき起こる。  $x \cdot y$ は内積

$$\text{このとき } \hat{w}_t = \left( w_t^T, \frac{b_t}{R} \right)^T = \left( w_{t-1}^T, \frac{b_{t-1}}{R} \right)^T + \eta y_i (x_i^T, R)^T = \hat{w}_{t-1} + \eta y_i \hat{x}_i \quad (2)$$

$$\text{なぜなら、 } b_t = b_{t-1} + \eta y_i R^2 \quad \Rightarrow \quad \frac{b_t}{R} = \frac{b_{t-1}}{R} + \eta y_i R$$

# 証明 つづき

$$\hat{w}_t = \left( w_t^T, \frac{b_t}{R} \right)^T = \left( w_{t-1}^T, \frac{b_{t-1}}{R} \right)^T + \eta y_i (x_i^T, R)^T = \hat{w}_{t-1} + \eta y_i \hat{x}_i \quad (2)$$

$$(1)より \quad \hat{w}_t \cdot \hat{w}_{opt} = \hat{w}_{t-1} \cdot \hat{w}_{opt} + \eta y_i (\hat{x}_i \cdot \hat{w}_{opt}) \geq \hat{w}_{t-1} \cdot \hat{w}_{opt} + \eta \gamma \quad (3)$$

$$\hat{w}_0 = 0とすれば(3)を繰り返して用いて \hat{w}_t \cdot \hat{w}_{opt} \geq t \eta \gamma \quad (4)$$

$$(2)より \|\hat{w}_t\|^2 = \|\hat{w}_{t-1}\|^2 + 2\eta y_i (\hat{w}_{t-1} \cdot \hat{x}_i) + \eta^2 \|\hat{x}_i\|^2$$

←  $\hat{x}_i$ は負例なので第2項は負

$$\leq \|\hat{w}_{t-1}\|^2 + \eta^2 \|\hat{x}_i\|^2 = \|\hat{w}_{t-1}\|^2 + \eta^2 (\|x_i\|^2 + R^2) \leq \|\hat{w}_{t-1}\|^2 + 2\eta^2 R^2$$

$$\Rightarrow \|\hat{w}_t\|^2 \leq 2t \eta^2 R^2 \quad \Rightarrow \|\hat{w}_t\| \leq \sqrt{2t} \eta R \quad (5)$$

$$(4)(5)より \|\hat{w}_{opt}\| \sqrt{2t} \eta R \geq \|\hat{w}_{opt}\| \|\hat{w}_t\| \geq \hat{w}_t \cdot \hat{w}_{opt} \geq t \eta \gamma$$

$$\Rightarrow t \leq 2 \left( \frac{R}{\gamma} \right)^2 \|\hat{w}_{opt}\|^2 \leq 2 \left( \frac{R}{\gamma} \right)^2 (\|\hat{w}_{opt}\|^2 + 1) = \left( \frac{2R}{\gamma} \right)^2 \quad \blacksquare$$