

大規模データの線形識別

Recent Advances of Large-scale Linear
Classification

Guo-Xun Yuan, Chia-Hua Ho, and Chih-Jen Lin
を中心にまとめた

データの性質

- 入力データの次元と入力データ数
 - 入力データの各次元を素性とも言う
 - 次元の高い場合はsparse(有効な素性が少ない)
 - 次元が低い場合はdense(多数の素性が有効)

	入力データの次元 小 10 ³ 以下	入力データの次元 大 10 ⁴ 以上
入力データ数 小 Mega	新規理論をとりあえず試す toy data	正解タグ付きテキストコーパス など (場合によっては画像も)
入力データ数 大 Giga-Tera	計測される数値データ (センサーデータ、市場 データ、など)	正解タグなしの生テキストコー パス
ストリーム 時系列到着	センサーデータ	Twitter、ソーシャルメディアなど

データの性質

- 次元が大きい(10の4乗以上)の場合は線形識別も非線形識別の精度に大きな差はない。
- 次元が低いと非線形識別のほうが精度がかなりよい。非線形化すなわち高次元化による特徴抽出能力の改善のため
- 教師データ(タグ付けデータ)と生データ
 - 教師データが大きければ、それを相手に学習
 - 教師データがなければクラスタリングなど
 - 小さな教師データと多数の生データ(実際は多い)
 - Semi-supervised learning, Active learning などチャレンジング

記法

$$(y_i, \mathbf{x}_i) \in \{+1, -1\} \times \mathbb{R}^n, i = 1, \dots, l$$

decision function: $d(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) + b$ $\phi(\mathbf{x}) = \mathbf{x}$ なら *linear*

$$\mathbf{w} = \sum_{i=1}^l \alpha_i \phi(\mathbf{x}_i)$$

以下では $\mathbf{w}^T \leftarrow [\mathbf{w}^T \ b]$ $\mathbf{x}_i^T \leftarrow [\mathbf{x}_i^T \ 1]$ とし、 b は陽には書かない

線形識別 linear classification は以下の最適化問題を含む

$$\min_{\mathbf{w}} f(\mathbf{w}) \equiv r(\mathbf{w}) + C \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i)$$

$r(\mathbf{w})$: 正則化項 $\xi(\mathbf{w}; \mathbf{x}_i, y_i)$: 損失

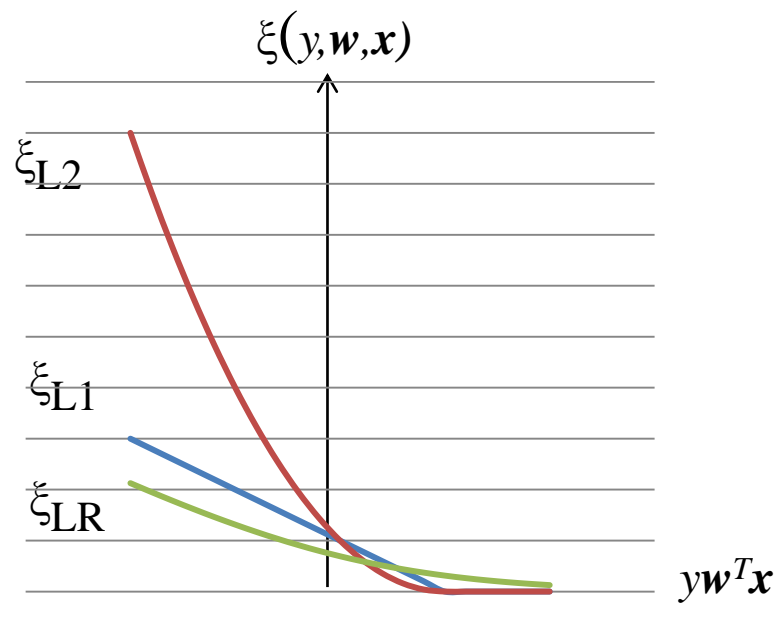
正則化項: r + 損失: ξ の最小化

➤ 各種の損失関数(右図)

$$\xi_{L1}(\mathbf{w}; \mathbf{x}, y) \equiv \max(0, 1 - y\mathbf{w}^T \mathbf{x})$$

$$\xi_{L2}(\mathbf{w}; \mathbf{x}, y) \equiv \max(0, 1 - y\mathbf{w}^T \mathbf{x})^2$$

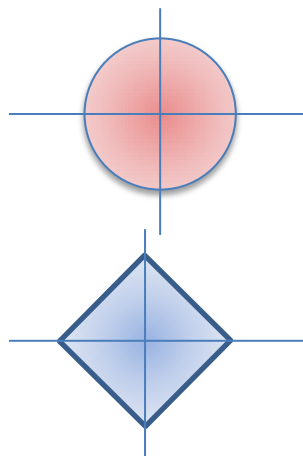
$$\xi_{LR}(\mathbf{w}; \mathbf{x}, y) \equiv \log(1 + \exp(-y\mathbf{w}^T \mathbf{x}))$$



➤ 各種の正則化項

$$r_{L2}(\mathbf{w}) \equiv \|\mathbf{w}\|_2^2 = \frac{1}{2} \sum_{i=1}^n w_i^2$$

$$r_{L1}(\mathbf{w}) \equiv \|\mathbf{w}\|_1^0 = \sum_{i=1}^n |w_i|$$



学習アルゴリズム

- 以下では種々の学習アルゴリズムについて述べる

双対化

➤ 双対問題を解くほうが良いこともある。双対化とは以下の通り:

Primal problem: minimize $f(x)$
 subject to $g_i(x) \leq 0 \quad i = 1, \dots, n$

Dual problem:

$q(\mu) = \inf L(x, \mu)$ where $L(x, \mu) = f(x) + \sum_{i=1}^n \mu_i g_i(x)$
 maximize $q(\mu)$
 subject to $\mu_i \geq 0 \quad i = 1, \dots, n$

双対化

$$SVM \text{ primal : } f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 \quad g_i(\mathbf{w}) = \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \leq 0$$

$$SVM \text{ dual : } \inf L(x, \mu) = \inf \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \mu_i \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i)$$

$$\Rightarrow \mathbf{w} = \sum_{i=1}^n y_i \mu_i \mathbf{x}_i$$

$$\Rightarrow \max -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \mu_i \boxed{y_i \mathbf{x}_i^T \mathbf{x}_j y_j} \mu_j + \sum_{i=1}^n \mu_i$$

$$\Rightarrow \min \frac{1}{2} \mu^T Q \mu - \mathbf{1}^T \mu \quad \text{subject to } 0 \leq \mu_i \leq C$$

この制約は $\max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \leq \eta_i$ and $\min \frac{1}{2} \|\mathbf{w}\|_2^2 + C \eta_i$

という線形分離できない場合の定式化から出る

(数理手法IV カーネル法 サポートベクターマシン 参照)

比較

- Gradient Descentのような低次の手法はupdateは簡単でメモリも少ない。収束が遅く、解は精密ではないが、実は性能に大差ない
- ニュートン法のような高次の手法は、収束は早いですが、updateのコストは高い（例えば、 Hessian^{-1} の計算とかLBFGSとか）。また、解空間がスムーズ（2階微分可能）でないといけない。精密な最適化には有効。
 - exp/logの計算はコストが高いことを忘れないように
 - 並列化する場合はコア間、ノード間の通信のコストも大きいことに注意

Pegasos : Shalev-Schwalz et al. ICML 2007

➤ L2正則化 + L1損失のSVM

$$\min_{\mathbf{w}} f(\mathbf{w}; B) \equiv \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{i \in B} \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \quad B \text{は全学習データ}$$

➤ Pegasosの更新は次式の劣微分による

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \partial f(\mathbf{w}; B)$$

$$\text{where } \partial f(\mathbf{w}; B) = \mathbf{w} - C \sum_{i \in B^+} y_i \mathbf{x}_i \quad B^+ \equiv \{i \mid i \in B, 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0\},$$

$$\eta = Cl / k, \quad l \text{はベクトル } \mathbf{w}, \mathbf{x}_i \text{の次元、} k \text{は繰り返し回数}$$

➤ 更新の後、 \mathbf{w} を半径 $(Cl)^{-1/2}$ の球にproject

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

➤ 以上を収束するまで繰り返す。

Trust Region Newton Method(TRON):導入

- データの次元: n が大きいと $\nabla^2 f(\mathbf{w})$ は $n \times n$ なのでメモリに置くのも逆行列の計算も大変。

- 逆行列計算はBFGSなどによる

- まずlogistic regressionについて

$$\min_{\mathbf{w}} f(\mathbf{w}) \equiv \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i))$$

$$\nabla f(\mathbf{w}) = \mathbf{w} + C \sum_{i=1}^l \left(\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} - 1 \right) y_i \mathbf{x}_i$$

$$\nabla^2 f(\mathbf{w}) = I + CX^T DX \quad D \text{は対角行列}$$

$$X = [\mathbf{x}_1^T \cdots \mathbf{x}_l^T]^T \quad D_{ii} = \left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} \left(1 - \left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} \right)$$

- するとNewton法は以下のようになる。

$$\mathbf{w}^{k+1} = \mathbf{w}^k + s^k$$

$$-\nabla f(\mathbf{w}^k) = \nabla^2 f(\mathbf{w}^k) s^k = (I + CX^T DX) s^k$$

$[D_{ij}]$ が対角行列であることの直観的説明

$$\min_{\mathbf{w}} f(\mathbf{w}) \equiv \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^l \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i))$$

$$\nabla f(\mathbf{w}) = \mathbf{w} + C \sum_{i=1}^l \left(\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} - 1 \right) y_i \mathbf{x}_i \quad \nabla^2 f(\mathbf{w}) = I + CX^T DX \quad D \text{は対角行列}$$

$$X = [\mathbf{x}_1^T \cdots \mathbf{x}_l^T]^T \quad D_{ii} = \left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} \left(1 - \left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} \right)$$

2×2で具体例

$$\begin{aligned} \nabla^2 \sum_i \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)) &= \nabla \sum_i \left(\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^{-1} - 1 \right) y_i \mathbf{x}_i \\ &= \sum_i \frac{\exp(-y_i \mathbf{w}^T \mathbf{x}_i) y_i \mathbf{x}_i^T y_i \mathbf{x}_i}{\left(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i) \right)^2} = \sum_i \mathbf{x}_i^T D_{ii} \mathbf{x}_i \end{aligned}$$

$$X^T = [\mathbf{x}_1 \quad \mathbf{x}_2] = \begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \end{bmatrix} \quad X = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}$$

$$\nabla^2 \sum_{i=1}^2 \log(1 + \exp(-y_i \mathbf{w}^T \mathbf{x}_i)) = \sum_{i=1}^2 \mathbf{x}_i^T D_{ii} \mathbf{x}_i = \begin{bmatrix} x_{11}^2 D_{11} + x_{21}^2 D_{22} & x_{11} x_{12} D_{11} + x_{21} x_{22} D_{22} \\ x_{11} x_{12} D_{11} + x_{21} x_{22} D_{22} & x_{12}^2 D_{11} + x_{22}^2 D_{22} \end{bmatrix}$$

$$= \begin{bmatrix} x_{11} D_{11} & x_{21} D_{22} \\ x_{12} D_{11} & x_{22} D_{22} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \end{bmatrix} \begin{bmatrix} D_{11} & 0 \\ 0 & D_{22} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}$$

Dは対角行列

➤ \mathbf{x}_i がsparse \rightarrow

$$DX = \begin{pmatrix} D_{11} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & D_{11} \end{pmatrix} \begin{pmatrix} 0 & \dots 0x0\dots & 0 \\ 0 & \dots 0x0. & \vdots \\ \vdots & .0x0\dots & \end{pmatrix} \approx \begin{pmatrix} 0 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & 0 \end{pmatrix}$$

➤ $\mathbf{w}^{k+1} = \mathbf{w}^k + s^k$

$$-\nabla f(\mathbf{w}^k) = \nabla^2 f(\mathbf{w}^k) s^k = (I + CX^T DX) s^k \approx s^k$$

$$\Rightarrow \mathbf{w}^{k+1} = \mathbf{w}^k - \nabla f(\mathbf{w}^k)$$

➤ $(I+\varepsilon)^{-1}$ は ε が小さいと近似計算もできそう。

➤ 1次近似なら $(I+\varepsilon)^{-1} = I - \varepsilon$

$$\approx \mathbf{w}^{k+1} = \mathbf{w}^k - \nabla f(\mathbf{w}^k) (I - CX^T DX)$$

➤ したがって、データが高次元Sparseならlogistic regressionはHessianの計算を避けられて、効率よく計算できる。

Trust Region Newton Method(TRON):

C.-J Lin et.al. SIAM J. Optimization 9 (1999)

➤ 以下の問題を f の2次の展開 q を用いて解く

$$\min_{\mathbf{w}} f(\mathbf{w}) \equiv r(\mathbf{w}) + C \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i) \quad r \text{ は } L2, \xi \text{ は微分可能}$$

$$f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w}) \approx q(\mathbf{d}) = \nabla f(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{w}) \mathbf{d}$$

➤ 最適化問題

予め決めた閾値

$$\text{step 1} \quad \min_{\mathbf{d}} q(\mathbf{d}) \quad \text{subject to } \|\mathbf{d}\|_2 \leq \Delta$$

$$\text{step 2} \quad \text{if } \rho \equiv \frac{f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})}{q(\mathbf{d})} > \eta_0 \text{ then } \mathbf{w} \leftarrow \mathbf{w} + \mathbf{d}$$

$$\text{step 3} \quad \text{Adjust } \Delta \text{ according to } \rho$$

Trust Region bound Δ_k の変え方

$\Delta_k \rightarrow \Delta_{k+1}$ の規則

$$\eta_1 < \eta_2 < 1 \quad \sigma_1 < \sigma_2 < 1 < \sigma_3$$

$$\Delta_{k+1} \in [\sigma_1 \min\{\|\mathbf{d}_k\|, \sigma_2 \Delta_k\}, \sigma_2 \Delta_k] \quad \text{if} \quad \rho_k \leq \eta_1$$

$$\Delta_{k+1} \in [\sigma_1 \Delta_k, \sigma_3 \Delta_k] \quad \text{if} \quad \rho_k \in (\eta_1, \eta_2)$$

$$\Delta_{k+1} \in [\Delta_k, \sigma_3 \Delta_k] \quad \text{if} \quad \rho_k \geq \eta_2$$

$$\frac{f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})}{q(\mathbf{d})} = \rho$$

大 \rightarrow Δ 大きく
小 \rightarrow Δ 小さく

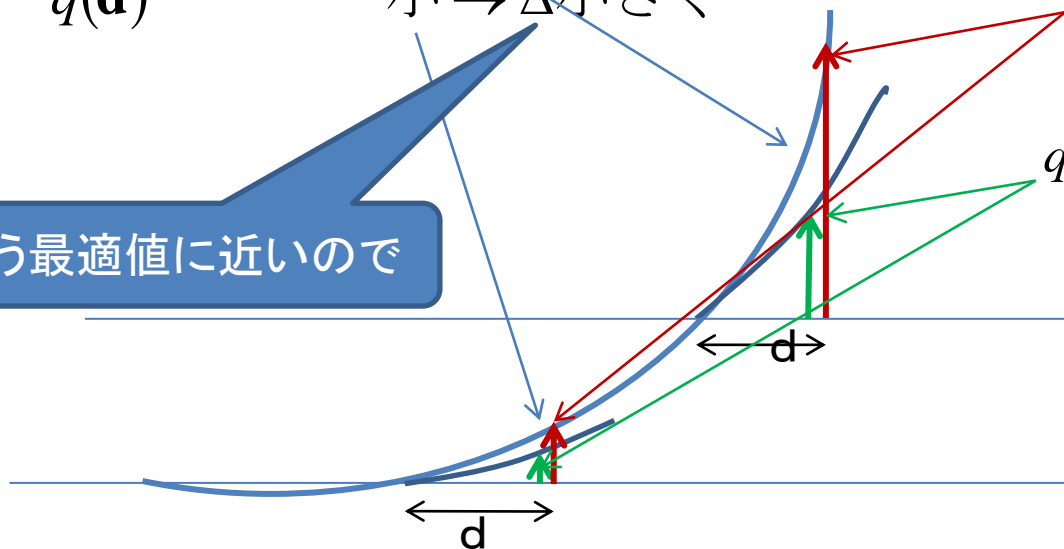
まだ最適値からは遠いので

もう最適値に近いので

$$f(\mathbf{w} + \mathbf{d}) - f(\mathbf{w})$$

$$q(\mathbf{d}) = \nabla f(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{w}) \mathbf{d}$$

このアイデアは目的関数が凸で微分可能なあたりにポイントがある！



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 $Q_{ij} \equiv y_i y_j \mathbf{x}_i^T \mathbf{x}_j$

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 α_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 α_i の最適化の中で1回計算すればよいが

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

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

➤ $\nabla_i f^D(\boldsymbol{\alpha}) = (Q\boldsymbol{\alpha})_i - 1 = y_i \mathbf{u}^T \mathbf{x}_i - 1$ (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、 α_i (CD10)の更新後

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

α 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 α is not optimal

For $i = 1, \dots, l$

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

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

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

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

Newton+ Coordinate Descent

J. H. Friedman, T. Hastie, and R. Tibshirani,
“Regularization paths for
generalized linear models via coordinate
descent,” Journal of Statistical
Software, vol. 33, no. 1, pp. 1–22, 2010.

➤ 最小化する目的関数 f は以下

L1正則化項

$$f(\mathbf{w}) = \|\mathbf{w}\|_1 + L(\mathbf{w}) \quad \text{where} \quad L(\mathbf{w}) \equiv C \sum_{i=1}^l \xi(\mathbf{w}; \mathbf{x}_i, y_i)$$

繰り返しごとに解くのは2次近似した以下の問題

$$\min_{\mathbf{d}} q(\mathbf{d}) \equiv \|\mathbf{w} + \mathbf{d}\|_1 - \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T H \mathbf{d} \quad (\text{NCD10})$$

where $H \equiv \nabla^2 L(\mathbf{w}) + \nu I$

小さな ν は H をpositive保つため導入

➤ ただし、L1正則化項(1-norm)のために全体を一度に解けないので、coordinate descentを導入

➤ (NCD10)を1変数化してみると以下のようになる

$$q(\mathbf{d} + ze_j) - q(\mathbf{d})$$

$$= \|\mathbf{w} + \mathbf{d} + ze_j\|_1 - \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T (\mathbf{d} + ze_j) + \frac{1}{2} (\mathbf{d} + ze_j)^T H (\mathbf{d} + ze_j)$$

$$- \|\mathbf{w} + \mathbf{d}\|_1 + \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T H \mathbf{d}$$

$$= |w_j + d_j + z| - |w_j + d_j| + G_j z + \frac{1}{2} H_{jj} z^2 \quad \text{where} \quad G \equiv \nabla L(\mathbf{w}) + H \mathbf{d}$$

➤ L1正則化における次元ごとの最適化の手法を利用すると

$$z = \begin{cases} -\frac{G_j + 1}{H_{jj}} & \text{if } G_j + 1 \leq H_{jj}(w_j + d_j) \\ -\frac{G_j - 1}{H_{jj}} & \text{if } G_j - 1 \geq H_{jj}(w_j + d_j) \\ -(w_j + d_j) & \text{otherwise} \end{cases} \quad (NCD20)$$

(NCD20)の直観的説明

$\min_z |w + d + z| - |w + d| + Gz + \frac{1}{2}Hz^2$ を求める。 $|w + d|$ は定数なので無視。

$$|w + d + z| + Gz + \frac{1}{2}Hz^2 = \begin{cases} g1(z) & \text{if } z \geq -(w + d) \\ g2(z) & \text{if } z \leq -(w + d) \end{cases}$$

$$g1(z) = w + d + z + Gz + \frac{1}{2}Hz^2, \quad g2(z) = -w - d - z + Gz + \frac{1}{2}Hz^2$$

$g1(z)$ の2次関数としての最小値が $-(w + d)$ より大きいのは

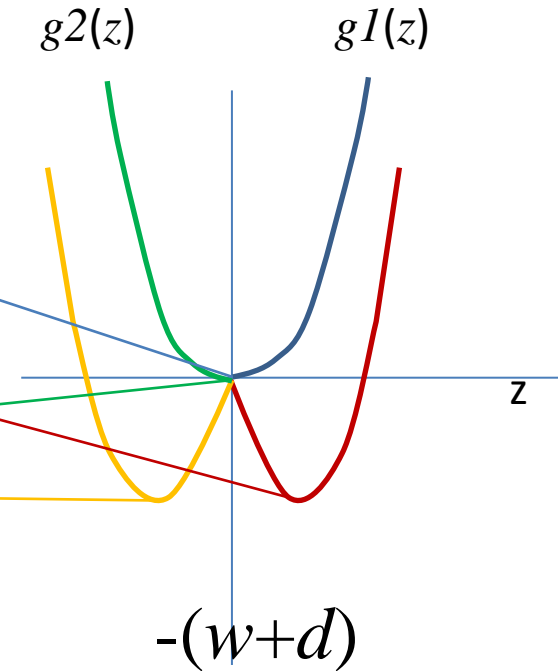
$$\frac{\partial g1(z)}{\partial z} = 0 \Rightarrow z = -\frac{G+1}{H} \geq -(w + d) \quad \text{すなわち } G+1 \leq H(w + d)$$

$$\text{then } z = -\frac{G+1}{H} \quad \text{otherwise } z = -(w + d)$$

$g2(z)$ の2次関数としての最小値が $-(w + d)$ より小さいのは

$$\frac{\partial g2(z)}{\partial z} = 0 \Rightarrow z = -\frac{G-1}{H} \leq -(w + d) \quad \text{すなわち } G-1 \geq H(w + d)$$

$$\text{then } z = -\frac{G-1}{H} \quad \text{otherwise } z = -(w + d)$$



- (NCD20)を各次元に対して行くと最適化問題 (NCD10)の1回の繰り返しができる。
- さて、(NCD10)を直線探索で繰り返して近似精度を上げる
- 直線探索による近似アルゴリズム

$\mathbf{w}, 0 < \beta, 0 < \sigma$ が与えられた

for $k = 1, 2, \dots$

(NCD10)の近似解 \mathbf{d} をcoordinate descent法(NCD20)で得る

$f(\mathbf{w} + \lambda \mathbf{d}) - f(\mathbf{w}) \leq \sigma \lambda (\|\mathbf{w} + \mathbf{d}\|_1 + \|\mathbf{w}\|_1 + \nabla L(\mathbf{w})^T \mathbf{d})$ を満たす

ような $\lambda = \max\{1, \beta, \beta^2, \dots\}$ を見つける

$\mathbf{w} \leftarrow \mathbf{w} + \lambda \mathbf{d}$

メモリに乗り切らないビッグデータの 処理に関して

- ディスクとメモリのデータ転送
 - 学習中にディスクとメモリのデータ転送回数が増えるのはもちろん問題だが
 - ビッグデータを(分割するにしても)メモリに1回転送する時間も膨大
- 分散環境(地理的な場合も含む)にデータがある場合
 - 分散環境での学習はもっと厄介
 - あまり公表された成果がない。Googleがちらっとブログで言った程度

オンライン学習による方法

- 元来が1データ毎の学習なので自然な方法
- L1-正則化&L2-損失のSVMは以下の更新式

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla^s \left(\frac{1}{2} \|\mathbf{w}\|_2^2 + C \max(0, 1 - y_i \mathbf{w}^T \mathbf{x}_i) \right) \quad \nabla^s \text{は sub - gradient}$$

$$\Rightarrow \text{if } 1 - y_i \mathbf{w}^T \mathbf{x}_i > 0 \text{ then } \mathbf{w} \leftarrow (1 - \eta) \mathbf{w} + \eta C y_i \mathbf{x}_i$$

収束速度が遅い

理論的決め方なし

- Coordinate Descentのように双対化して解く

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

η を決める必要なし

➤ 収束の遅さを改善するために高次情報を使う

$$\mathbf{w} \leftarrow \mathbf{w} - \eta H^{-1} \nabla^S \left(\frac{1}{2} \|\mathbf{w}\|_2^2 + C \max(0, 1 - y_i \mathbf{w}^T x_i) \right)$$

$$H = \nabla^2 \left(\frac{1}{2} \|\mathbf{w}\|_2^2 + C \max(0, 1 - y_i \mathbf{w}^T x_i) \right) \Rightarrow \text{対角行列になることが多い}$$

$\Rightarrow H^{-1}$ の計算が重くない

バッチ学習による方法

- ディスクアクセスを減らすために一度のあるままとまりのデータ(=ブロック)を読み込み学習
- ブロックBだけではなく、メモリにcacheデータも残す

$$\min_{\alpha} f^D(\alpha) \equiv \frac{1}{2} \alpha^T Q \alpha - \mathbf{e}^T \alpha \quad \text{subject to} \quad 0 \leq \alpha_i \leq C \quad \forall i = 1, \dots, l$$

where $Q_{ij} \equiv y_i y_j \mathbf{x}_i^T \mathbf{x}_j$ を解くために

ブロック B は全データ $= \{B_1, \dots, B_m\}$ から1個ずつ選んで最適化
条件 : $0 \leq \alpha_i + d_i \leq C$ for $\forall i \in B$ and $d_i = 0$ for $\forall i \notin B$

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

$$= \frac{1}{2} \mathbf{d}_B^T Q_{BB} \mathbf{d}_B + \sum_{i \in B} y_i d_i (\mathbf{u}^T \mathbf{x}_i) - \mathbf{d}_B^T \mathbf{e}_B \quad \text{where } \mathbf{u} = \sum_{i=1}^l y_i \alpha_i \mathbf{x}_i$$

Q_{BB} は Q のうち B に関する部分行列