

# 最適化と学習アルゴリズム

勾配降下法

直線探索

劣勾配降下法

ニュートン法

確率的勾配降下法

# 学習における計算効率の問題

- 既に説明した回帰、識別(分類)の解法で得た閉じた式は逆行列計算が必要。例えば、線形回帰の場合の正規方程式の解である重みベクトルは
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$
- 学習データが高次元になると $(\mathbf{X}^T \mathbf{X})$ の次元も大きくなり、逆行列の計算はコストが高い(次元数の3乗)
- 要は損失(=Loss)を最小化したい
  - 正規方程式の場合は2乗誤差だった。
- 逆行列計算を避けて最適化する方法が実用上重要。

# 記法：損失(Loss)

入力データ： $\mathbf{x}_i$ , 重みベクトル： $\mathbf{w}$ は $D$ 次元ベクトルとする  
 $\mathbf{x}_i$ に対する分類の正解： $y_i$ (= -1 or 1) ただしデータは $m$ 個あるとする

$$L(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \text{loss}(\mathbf{w} \text{ on } (\mathbf{x}_i, y_i))$$

例：

2乗損失の場合 
$$L(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \|\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i\|_2^2 = \frac{1}{m} \sum_{i=1}^m (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

ヒンジ損失の場合 
$$L(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m [1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle]_+$$
  
ここで  $[z]_+ = \max(0, z)$

# 記法： 勾配(Gradient)とヘッセ行列(Hessian)

$\mathbf{w}$ は $D$ 次元ベクトルとする

$$pノルム : \|\mathbf{w}\|_p = (w_1^p + \dots + w_D^p)^{1/p} \quad \|\mathbf{w}\|_p^p = (w_1^p + \dots + w_D^p)$$

$\mathbf{w} = \mathbf{w}^{(k)}$ における $L$ の勾配

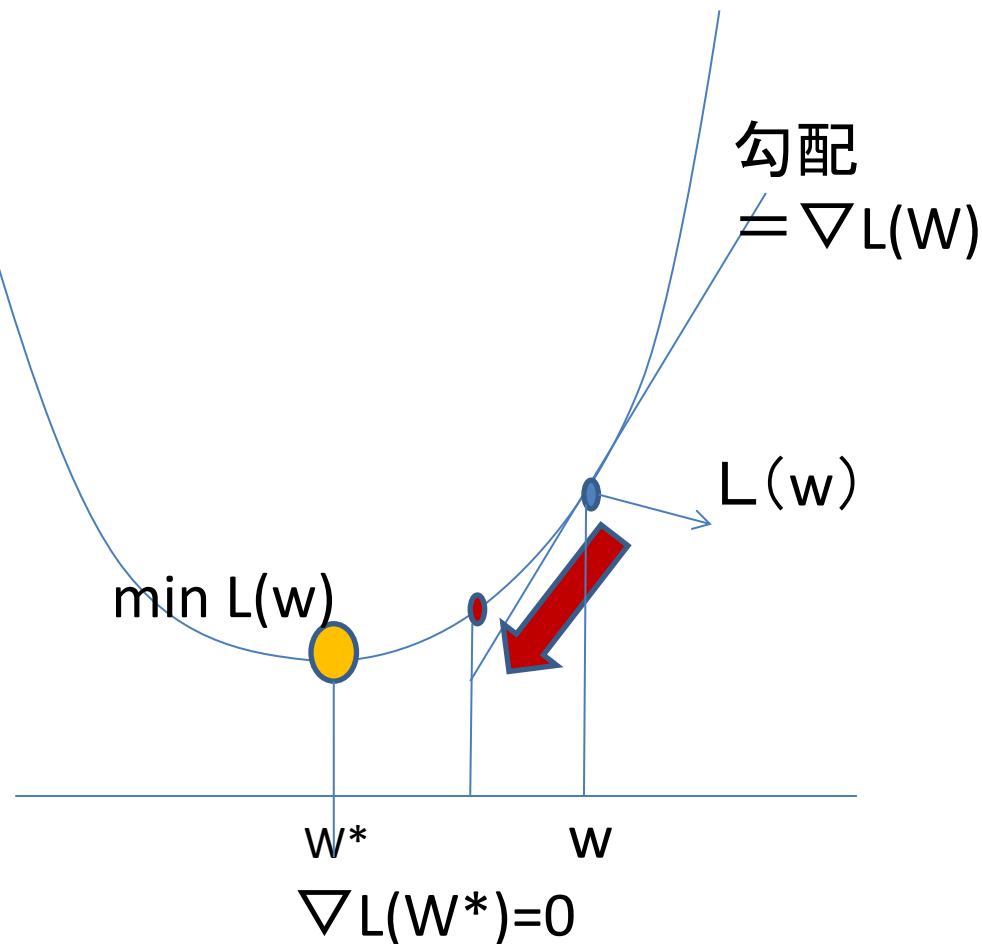
$$\nabla L(\mathbf{w}^{(k)}) = \left( \left. \frac{\partial L(\mathbf{w})}{\partial w_1} \right|_{w_1=w_1^{(k)}}, \dots, \left. \frac{\partial L(\mathbf{w})}{\partial w_n} \right|_{w_D=w_D^{(k)}} \right)^T$$

$\mathbf{w} = \mathbf{w}^{(k)}$ における $L$ のヘッセ行列

$$HL(\mathbf{w}^{(k)}) = \begin{pmatrix} \left. \frac{\partial^2 L(\mathbf{w})}{\partial w_1 \partial w_1} \right|_{w_1=w_1^{(k)}} & \dots & \left. \frac{\partial^2 L(\mathbf{w})}{\partial w_1 \partial w_D} \right|_{w_1=w_1^{(k)}, w_D=w_D^{(k)}} \\ \vdots & \ddots & \vdots \\ \left. \frac{\partial^2 L(\mathbf{w})}{\partial w_D \partial w_1} \right|_{w_D=w_D^{(k)}, w_1=w_1^{(k)}} & \dots & \left. \frac{\partial^2 L(\mathbf{w})}{\partial w_D \partial w_D} \right|_{w_D=w_D^{(k)}} \end{pmatrix}$$

# 勾配降下法

- 右図で分るように  $L(w)$  が凸の場合は、勾配  $\nabla L(w)$  の逆方向に進むと  $W^* = \operatorname{argmin} L(W)$  に向かって進む。
- 図は1次元だが実際は多次元なので、勾配の方向が大切なことに注意



# 定式化

- 勾配降下法(Gradient Descent)の定式化する

$\mathbf{w}$  : 重みベクトル      $L(\mathbf{w})$ : 重みベクトル $\mathbf{w}$ のときの損失

$\mathbf{w}^{(k)}$ :  $k$ 回目の繰り返し結果の更新された重みベクトル

最適な重みベクトル:  $\mathbf{w}^* = \arg \min_{\mathbf{w}} L(\mathbf{w})$

を求めるには勾配の逆方向に進むので下の式となる。

重みベクトルの更新式:  $\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{v}} L(\mathbf{v} - (\mathbf{w}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{w}^{(k)})))$

$L(*) = \|*\|_2^2$ : 2乗損失だと  $\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{v}} \|\mathbf{v} - (\mathbf{w}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{w}^{(k)}))\|_2^2$

損失を微分して0と置けば閉じた解が求まり

$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{w}^{(k)})$  となる

2乗損失でないとう簡単な式  
ではない

# 最急降下法アルゴリズム

- $L(\mathbf{w})$ を最小化するような $\mathbf{w}$ を求めるために、 $\mathbf{w}$ の現在の値 $\mathbf{w}^{(k)}$ から勾配 $\nabla L(\mathbf{w})$ の逆方向に進むアルゴリズム

初期値 =  $\mathbf{w}^{(0)}$ ;  $k = 0$ ;  $\mathbf{w}^* = \arg \min L(\mathbf{w})$ ;

*step1*:  $\mathbf{w}^{(k)}$ が $\mathbf{w}^*$ に十分近いなら終了

(具体的には $\mathbf{w}^{(k)}$ が $\mathbf{w}^{(k-1)}$ に十分近くて収束した)

*step2*:  $\nabla L(\mathbf{w}^{(k)})$ を計算

*step3*:  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{w}^{(k)})$        $\alpha^{(k)}$ は進む幅(step size)

*step4*:  $k = k + 1$     として*step1*へ

$L(\mathbf{w}) = \|y_n - \langle \mathbf{w}, \phi(\mathbf{x}_n) \rangle\|_2^2$       つまり2乗損失なら

$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} (y_n - \langle \mathbf{w}^{(k)}, \phi(\mathbf{x}_n) \rangle) \phi(\mathbf{x}_n)$

# 最急降下法の収束の評価

## ➤ 前提条件

➤ 凸 &  $\nabla^2 L(\mathbf{w})$  positive

➤ Lipschitz条件:  $|L(\mathbf{w}^{(k)}) - L(\mathbf{w}^{(k+1)})| \leq G \|\mathbf{w}^{(k)} - \mathbf{w}^{(k+1)}\|_2$   
 $\Rightarrow \|\nabla L(\mathbf{w}^{(k)})\|_2 < G$

## ➤ テイラー展開により近似

$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \nabla L(\mathbf{w}^{(k)})$  とし、 $L(\mathbf{w}^{(k+1)})$  を2次のテイラー近似すると

$$L(\mathbf{w}^{(k+1)}) - L(\mathbf{w}^{(k)}) \approx -\alpha^{(k)} \|\nabla L(\mathbf{w}^{(k)})\|_2 + \frac{(\alpha^{(k)})^2}{2} \|\nabla L(\mathbf{w}^{(k)})\|_2^2 \|\nabla^2 L(\mathbf{w}^{(k)})\|_2$$

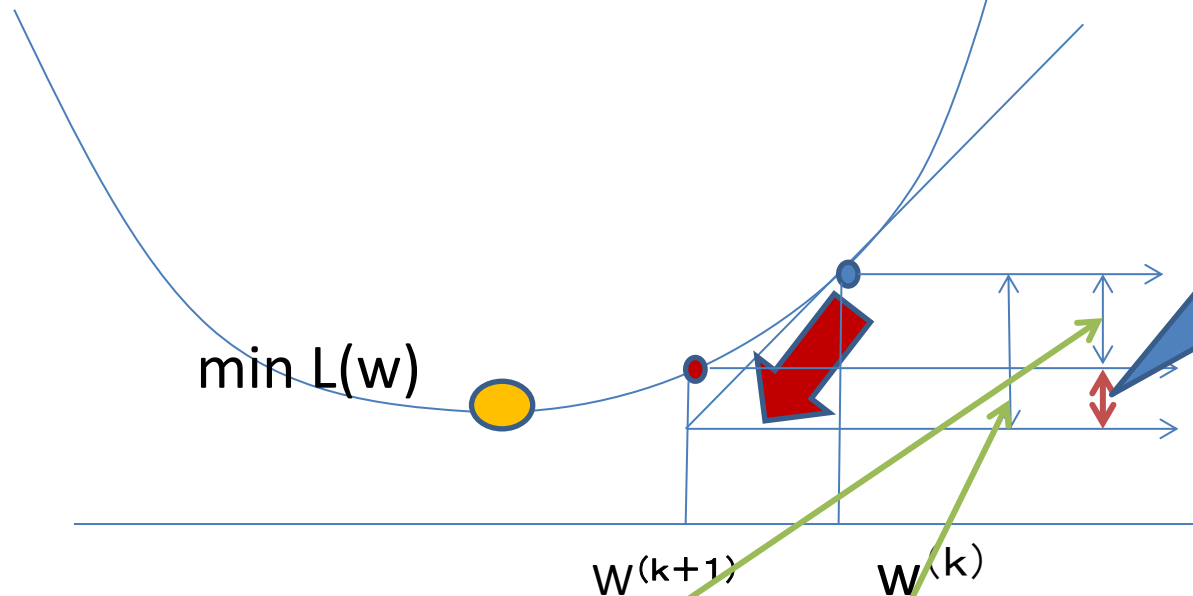
$\|\nabla^2 L(\mathbf{w}^{(k)})\|_2 \geq 0$  により

$$L(\mathbf{w}^{(k+1)}) - L(\mathbf{w}^{(k)}) \geq -\alpha^{(k)} \|\nabla L(\mathbf{w}^{(k)})\|_2^2$$

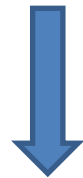
$$\alpha^{(k)} \|\nabla L(\mathbf{w}^{(k)})\|_2^2 \geq L(\mathbf{w}^{(k)}) - L(\mathbf{w}^{(k+1)}) \quad (1)$$



# 直観的説明



凸かつ $\nabla^2 L(w) \geq 0$ なので、 $\min L$ に近づくにつれてこの差は小さくなる



適当な $G$ で上から抑えられる

$$|L(\mathbf{w}^{(k)}) - L(\mathbf{w}^{(k+1)})| \leq G \|\mathbf{w}^{(k)} - \mathbf{w}^{(k+1)}\|_2$$

Lipshitz条件より以下のような $G$ がある

$$\Rightarrow |L(\mathbf{w}^{(k)}) - L(\mathbf{w}^{(k+1)})| \leq \|\nabla L(\mathbf{w}^{(k)})\|_2 \|\mathbf{w}^{(k)} - \mathbf{w}^{(k+1)}\|_2 \leq G \|\mathbf{w}^{(k)} - \mathbf{w}^{(k+1)}\|_2$$

$$\Rightarrow \|\nabla L(\mathbf{w}^{(k)})\|_2 < G$$

# 最急降下法の収束の評価: つづき

$$\alpha^{(k)} \left\| \left( \nabla L(\mathbf{w}^{(k)}) \right) \right\|_2^2 \geq L(\mathbf{w}^{(k)}) - L(\mathbf{w}^{(k+1)}) \quad (1)$$

$\mathbf{w}^{(k+1)}$ が $\mathbf{w}^*$ に十分近い ( $\mathbf{w}^{(k+1)} \approx \mathbf{w}^*$ )ないし等しいとすると

$$\alpha^{(k)} \left\| \left( \nabla L(\mathbf{w}^{(k)}) \right) \right\|_2^2 \geq L(\mathbf{w}^{(k)}) - L(\mathbf{w}^*)$$

ここで、 $B \leq \|\mathbf{w}^*\|_2$ という $B$ を選び、 $\alpha^{(k)} = \frac{B}{G\sqrt{k}}$ とおく

Lipschitz条件  $\|\nabla L(\mathbf{w}^{(k)})\|_2 < G$  より

$$\Rightarrow L(\mathbf{w}^{(k)}) - L(\mathbf{w}^*) \leq \frac{B}{G\sqrt{k}} \cdot G^2 = \frac{BG}{\sqrt{k}}$$

誤差を $\varepsilon$ 以下にするのはオーダーとして:  $O(\varepsilon) = O\left(\frac{BG}{\sqrt{k}}\right)$

$\Rightarrow O\left(\frac{B^2G^2}{\varepsilon^2}\right)$ 回の更新で誤差 $\varepsilon$ 以下に収束

# 捕捉

$\alpha^{(k)}$ をやたらと大きくすると無意味！

$$|L(\mathbf{w}^{(k)}) - L(\mathbf{w}^*)| \leq G \|\mathbf{w}^{(k)} - \mathbf{w}^*\|_2$$

$$B \leq \|\mathbf{w}^*\|_2 \quad \Rightarrow \quad B \leq \|\mathbf{w}^{(k)} - \mathbf{w}^*\|_2 \text{ くらいにはなるだろう}$$

$$\Rightarrow |L(\mathbf{w}^{(k)}) - L(\mathbf{w}^*)| \approx BG \leq G \|\mathbf{w}^{(k)} - \mathbf{w}^*\|_2$$

このあたりを想定してという  $B$  を選び、 $\alpha^{(k)} = \frac{B}{G\sqrt{k}}$  とおく

という雰囲気

# step size $\alpha$ を決める: 直線探索

➤  $\alpha^{(k)}$ が小さすぎると収束が遅々として進まず、  
大きすぎると最悪の場合最適値を飛び越す

➤ 最適な $\alpha^{(k)}$ を求める問題: 直線探索

$$\min_{\alpha^{(k)}} L(\mathbf{w}^{(k)} + \alpha^{(k)} \mathbf{d}) \quad \mathbf{d} \text{は降下方向(つまり } -\nabla L(\mathbf{w}^{(k)}) \text{方向)}$$

の単位ベクトル

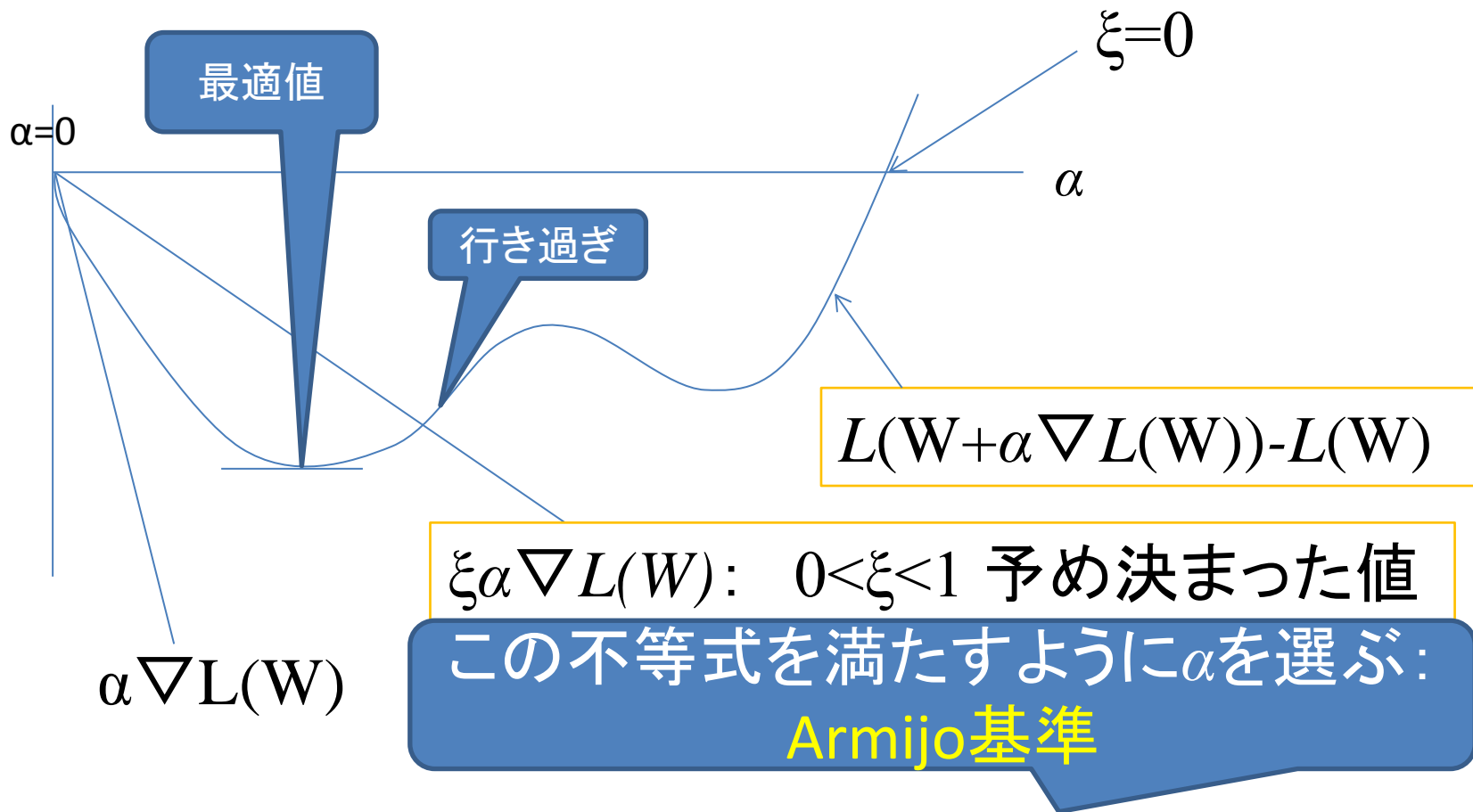
$$\text{subject to } \alpha^{(k)} > 0$$

➤ 厳密に解くのは困難

# 直線探索の代案

- $k$ の減少関数にする→ i.e.  $\alpha^{(k)} \propto 1/\sqrt{k}$
- しかし、あまりにもヒューリスティック
- もう少し客観的な基準として
  
- Armijo基準とWolfe基準がある

# Armijo基準とWolfe基準



## Wolfe基準:

$\alpha$  が小さくなりすぎないため。 $\alpha=0$ だと明らかに成立しない(! $\nabla L(W)^T \mathbf{d} < 0$ )

$$L(W + \alpha \mathbf{d}) - L(W) \leq \xi \alpha \nabla L(W)^T \mathbf{d}$$

$$\begin{aligned} \mu \nabla L(W)^T \mathbf{d} &\leq \nabla L(W + \alpha \mathbf{d})^T \mathbf{d} \\ \xi &< \mu < 1 \end{aligned}$$

# 直線探索の代案:2分割法

➤ k回目の繰り返しにおいて以下のアルゴリズムで $\alpha^{(k)}$ を減少させる

Step1  $t=1, \alpha_t^{(k)}=\alpha^{(k)}$

Step2 if 停止基準(Armijo基準等)が満たされる  
then return  $\alpha_t^{(k)}$   
else  $\alpha_t^{(k)}=\alpha^{(k)}/2$

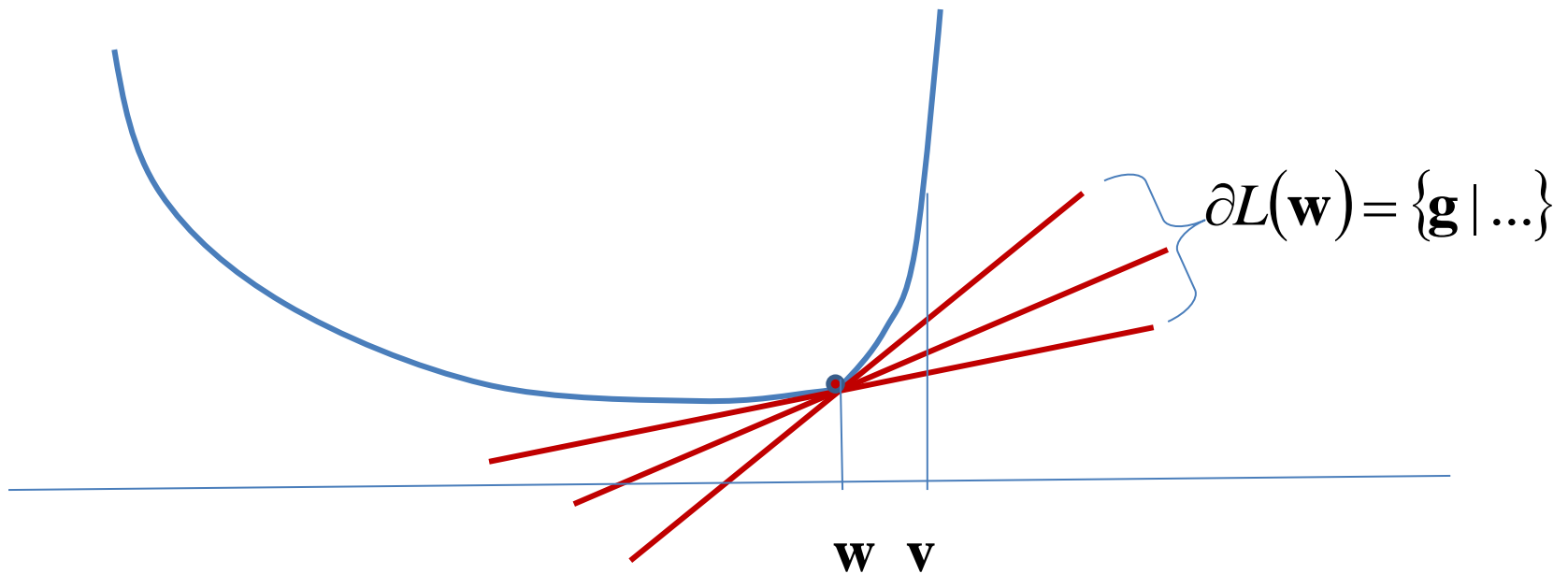
Step3  $k \leftarrow k+1$

Step4 Step 2へ

# 微分可能でない場合： 劣勾配(sub-gradient)の利用

凸関数 : convex :  $L$  の subgradient

$$\partial L(\mathbf{w}) = \left\{ \mathbf{g} \mid \forall \mathbf{v} \geq 0 \quad L(\mathbf{v}) \geq L(\mathbf{w}) + (\mathbf{v} - \mathbf{w})^T \mathbf{g} \right\}$$





# 例

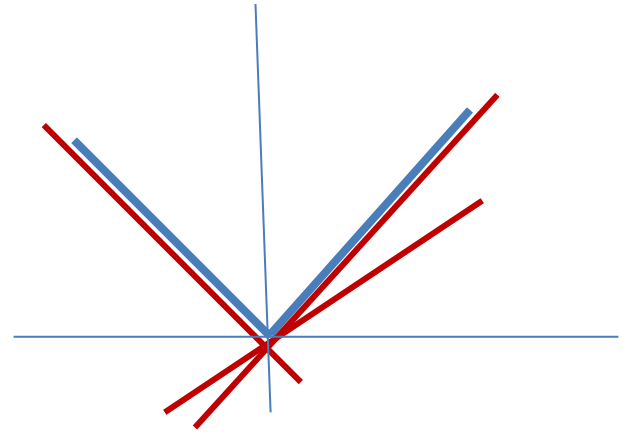
*L1 loss*

$$L(w) = \|w\|_1 = |w|$$

$$\nabla L(w) = -1 \quad w < 0$$

$$\nabla L(w) = [-1, 1]$$

$$\nabla L(w) = 1 \quad w > 0$$



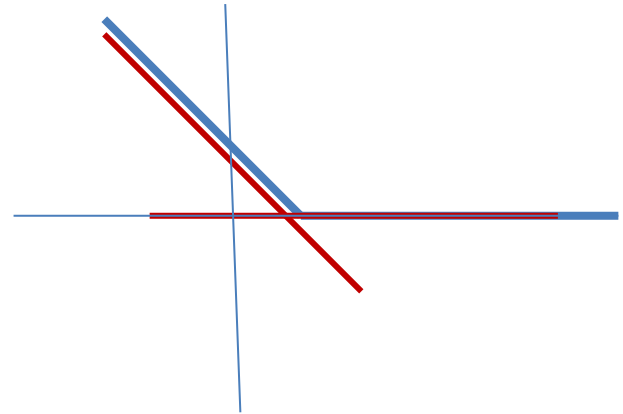
*hinge loss*

$$L(w) = [1 - w]_+$$

$$\nabla L(w) = -1 \quad w < 1$$

$$\nabla L(w) = [-1, 0]$$

$$\nabla L(w) = 0 \quad w > 1$$



# 劣勾配降下法

## Sub-Gradient Descent

初期値 =  $\mathbf{w}^{(0)}$ ;  $k = 0$ ;  $\mathbf{w}^* = \arg \min L(\mathbf{w})$ ;

*step1*:  $\mathbf{w}^{(k)}$ が $\mathbf{w}^*$ に十分近いなら終了

(具体的には $\mathbf{w}^{(k)}$ が $\mathbf{w}^{(k-1)}$ に十分近くて収束した)

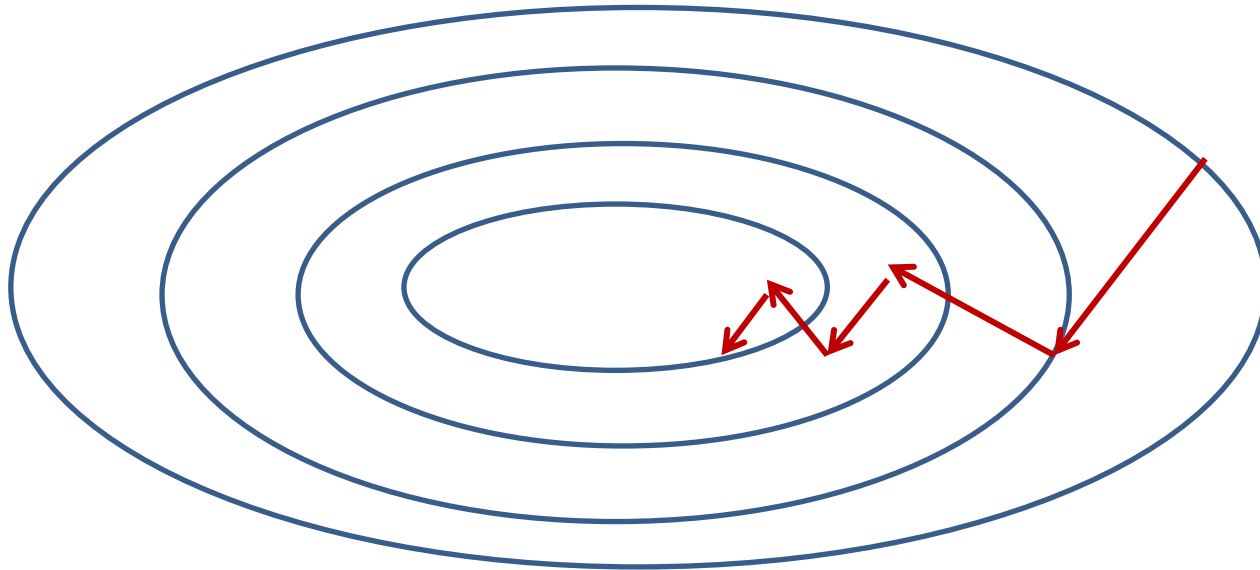
*step 2* sub - gradient :  $\partial L(\mathbf{w})$

=  $\{\mathbf{g} \mid \forall \mathbf{v} \geq 0 \quad L(\mathbf{v}) \geq L(\mathbf{w}) + (\mathbf{v} - \mathbf{w}) \cdot \mathbf{g}\}$  を計算

*step3*:  $\mathbf{w}^{(k+1)} = \arg \min_{\mathbf{v}, \mathbf{g} \in \partial L(\mathbf{w})} L(\mathbf{v} - (\mathbf{w}^{(k)} - \alpha^{(k)} \mathbf{g}))$

*step4*:  $k = k + 1$  として*step1*へ

# 最急降下法がうまくいかない場合



- ジグザグに動きながら最適解に近づくので収束する速度が遅い(効率が良くない)
- 勾配以外の情報も使ってみる→ニュートン法

# ニュートン法

- 最小化したいのは損失  $L(\mathbf{w})$ 
  - ただし、直接には最適化できない場合を考えている
- 今  $\mathbf{w}$  の現在値  $\mathbf{w}^{(k)}$  が与えられていて、 $d$  を加えることによってより最適値に近い  $\mathbf{w}^{(k+1)}$  を得たい。
- ここで  $\mathbf{w}^{(k)}$  が定数で  $d$  を変数と思って
- $L(\mathbf{w}^{(k+1)}) = L(\mathbf{w}^{(k)} + d)$ 

の左辺を直接最適化する代わりに右辺を最小にする  $d$  を求める。
- この結果を  $d^{(k)}$  とし、よりよい  $L$  を  $L(\mathbf{w}^{(k)} + d^{(k)})$  として、繰り返すことによって真の最適値に近づこうとする
- ここで  $L(\mathbf{w}^{(k)} + d)$  を2次のテーラー近似する場合はニュートン法

➤2次のテーラー展開

$$L(\mathbf{w}^{(k)} + \mathbf{d}) \approx L(\mathbf{w}^{(k)}) + \nabla L(\mathbf{w}^{(k)})^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{H} L(\mathbf{w}^{(k)}) \mathbf{d}$$

➤右辺を $\mathbf{d}$ の関数と見て最小化しニュートン方向 $\mathbf{d}^{(k)}$ を求める

$$\mathbf{d}^{(k)} = \arg \min_{\mathbf{d}} \text{右辺} = -\mathbf{H} L(\mathbf{w}^{(k)})^{-1} \nabla L(\mathbf{w}^{(k)})$$

## ニュートン法のアルゴリズム

*step0*:  $k = 0$ ; 初期値 =  $\mathbf{w}^{(0)}$ ;  $\mathbf{w}^* = \arg \min L(\mathbf{w})$ ;

*step1*:  $\mathbf{w}^{(k)}$ が $\mathbf{w}^*$ に十分近いなら終了

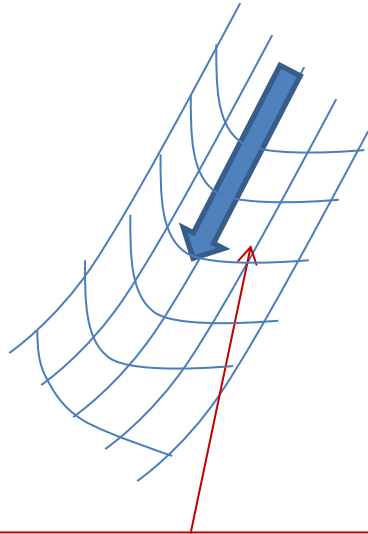
*step2*: ニュートン方向  $\mathbf{d}^{(k)} = -\mathbf{H} L(\mathbf{w}^{(k)})^{-1} \nabla L(\mathbf{w}^{(k)})$ を計算

*step3*:  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \mathbf{H} L(\mathbf{w}^{(k)})^{-1} \nabla L(\mathbf{w}^{(k)})$   $\alpha^{(k)}$ はstep size

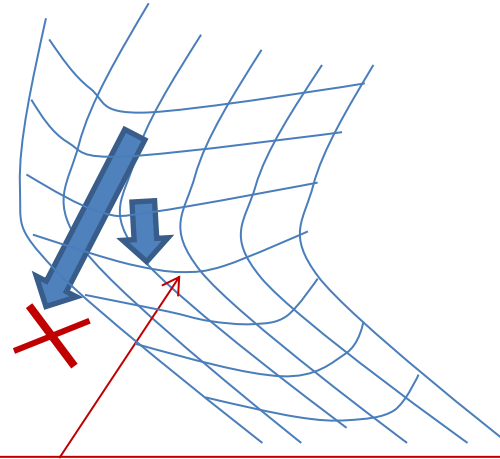
*step4*:  $k = k + 1$  として*step1*へ

$\alpha^{(k)}$ は直線探索で決めてもよい

# ニュートン法の直感的な解釈



この辺りは $\nabla$ の変化すなわちヘッセ行列  $HL(W)$  が小さいので大きく進む方が効率がよい  
→  $HL(W)^{-1}$  : 大



この辺りは $\nabla$ の変化すなわちヘッセ行列  $HL(W)$  が大きいので大きく進むと不味い方向に行ってしまう。  
→  $HL(W)^{-1}$  : 小

# ニュートン法の問題点

- ニュートン法は勾配降下法より多くの場合で性能がよいが、ヘッセ行列の逆行列を繰り返しの度に計算しなければならない
- $N \times N$ の行列の逆行列の計算には $O(N^3)$ のコストがかかる。 $N$ が大きくなると実用的でなくなる
- ヘッセ行列の逆行列を近似する行列を行列のかけ算だけによる更新式を使って $O(N^2)$ で計算する方法もある。
  - BFGS公式のヘッセ逆行列版 (cf. 共立出版「最適化法」p.122)

# 確率的勾配降下法

## Stochastic Gradient Descent(SGD)

- ここまで述べてきた方法は全データを一挙に学習に用いるバッチ処理
- 1データないし少数のデータ毎に重みベクトル $w$ 学習する方法: オンライン学習に近い
  - 最終結果としては、1ないし少数データ毎の学習で得られた $w$ の平均
  - メモリに全データを展開しなくてよい可能性もある→省メモリ→しかし、全データを何回も繰り返し使うなら省メモリにならない
  - →1つのデータは1回しか見ない→ストリームマイニング的
- 2種類のSGD
  - Sample Approximation(SA):1データ毎の学習



# SGD(SA) と バッチ学習

$$\mathbf{w}_0$$

$$\mathbf{w}_1 = \mathbf{w}_0$$

$$- \alpha^{(0)} \nabla \text{loss}(\mathbf{w}_0 \text{ on } (\mathbf{x}_1, y_1))$$

$$\mathbf{w}_2 = \mathbf{w}_1$$

$$- \alpha^{(1)} \nabla \text{loss}(\mathbf{w}_1 \text{ on } (\mathbf{x}_2, y_2))$$

$$\mathbf{w}_3 = \mathbf{w}_2$$

$$- \alpha^{(2)} \nabla \text{loss}(\mathbf{w}_2 \text{ on } (\mathbf{x}_3, y_3))$$

⋮

$$\mathbf{w}_i = \mathbf{w}_{i-1}$$

$$- \alpha^{(i-1)} \nabla \text{loss}(\mathbf{w}_{i-1} \text{ on } (\mathbf{x}_i, y_i))$$

⋮

$$\mathbf{w}_m = \mathbf{w}_{m-1}$$

$$- \alpha^{(m-1)} \nabla \text{loss}(\mathbf{w}_{m-1} \text{ on } (\mathbf{x}_m, y_m))$$

⇓

$$\mathbf{w} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}_i$$

このデータ生成が確率的

$\mathbf{x}_1, y_1$
$\mathbf{x}_2, y_2$
$\mathbf{x}_3, y_3$
⋮
$\mathbf{x}_i, y_i$
⋮
$\mathbf{x}_m, y_m$

$$\mathbf{w}^{(k)}$$

$$\nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_1, y_1))$$

$$\nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_2, y_2))$$

$$\nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_3, y_3))$$

⋮

$$\nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_i, y_i))$$

⋮

$$\nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_m, y_m))$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \left( \frac{1}{m} \sum_{i=1}^m \nabla \text{loss}(\mathbf{w}^{(k)} \text{ on } (\mathbf{x}_i, y_i)) \right)$$

さらに $k$ を+1して上の処理を繰り返す

# Sample Average Approximation(SAA)

## SA inside SAA

- 全データからランダムにm個のデータをサンプルしてm個のデータからなるデータ集合を作る。
- このようにして異なるデータ集合をk個作る。
- データ集合(i)にSGD(SA)を適用して学習し、重みベクトル $w(i)$ を計算
- 最終結果 
$$\mathbf{w} = \frac{1}{m} \sum_{i=1}^m \mathbf{w}(i)$$

# 例：SGDによるサポートベクターマシン

$$\min_{\|\mathbf{w}\|_2 < B} \frac{1}{m} \sum_{i=1}^m [1 - y_i(\mathbf{w}^T \mathbf{x}_i)]_+ = \min_{\mathbf{w}} \left( \frac{1}{m} \sum_{i=1}^m [1 - y_i(\mathbf{w}^T \mathbf{x}_i)]_+ + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right)$$

に注意するとSVDは以下のようなになる

初期値： $\mathbf{w}^{(0)}$ , 全データ数： $n$

以下を $k$ 回繰り返す

{ $i$ を $[1..n]$ からランダムに生成

if  $y_i(\mathbf{w}^T \mathbf{x}_i) < 1$  then  $\mathbf{w} \leftarrow \mathbf{w} + \alpha y_i \mathbf{x}_i$

if  $\|\mathbf{w}\|_2 \geq B$  then  $\mathbf{w} \leftarrow B\mathbf{w} / \|\mathbf{w}\|_2$

$\mathbf{w}_{sum} = \mathbf{w}_{sum} + \mathbf{w}$ }

結果： $\mathbf{w} = \mathbf{w}_{sum} / k$