

計画数理学特論

～第7回：勾配法の基礎・応用～

担当：蓮池隆（経営システム工学科）
e-mail: thasuike@waseda.jp

本日の講義について

- 勾配法の基礎として, 凸2次関数を例とした最急降下法の収束速度を考察
- 勾配法の応用として, アダブーストと座標降下法との関わり, 多層パーセプトロンと勾配法との関わりを紹介(詳細までは踏み込まない予定. 詳細を知りたい方は, それぞれの専門書を参照すること)

(再掲)直線探索を用いる反復法

(反復法の流れ)

STEP0 : 初期点 : $\mathbf{x}_0 \in \mathbb{R}^n$ を設定し, $k \leftarrow 0$ とする.

STEP1 : 停止条件が満たされるならば, \mathbf{x}_k を解として出力し, 終了.

STEP2 : 降下方向として \mathbf{d}_k を設定.

STEP3 : $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, $\alpha \geq 0$ に対する直線探索により, ステップ幅 α_k を計算.

STEP4 : $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$ と更新.

STEP5 : $k \leftarrow k + 1$ として, STEP1へ戻る.

(再掲)座標降下法のアルゴリズム

STEP0 : 初期点 : $\mathbf{x}_0 \in \mathbb{R}^n$ を設定し, $k \leftarrow 0$ とする.

STEP1 : 停止条件が満たされるならば, \mathbf{x}_k を解として出力し, 終了.

STEP2 : 降下方向として \mathbf{d}_k を, $\pm \mathbf{e}_1, \pm \mathbf{e}_2, \dots, \pm \mathbf{e}_n$ の中から選択する.

STEP3 : $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, $\alpha \geq 0$ に対する直線探索により, ステップ幅 α_k を計算.

STEP4 : $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$ と更新.

STEP5 : $k \leftarrow k + 1$ として, STEP1へ戻る.

反復法で, STEP 2 の内容が座標降下法に合った内容に変わっただけ

(再掲)最急降下法のアルゴリズム

STEP0 : 初期点 : $\mathbf{x}_0 \in \mathbb{R}^n$ を設定し, $k \leftarrow 0$ とする.

STEP1 : 停止条件が満たされるならば, \mathbf{x}_k を解として出力し, 終了.

STEP2 : 降下方向として $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$ を設定する.

STEP3 : $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, $\alpha \geq 0$ に対する直線探索により, ステップ幅 α_k を計算.

STEP4 : $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \alpha_k \mathbf{d}_k$ と更新.

STEP5 : $k \leftarrow k + 1$ として, STEP1へ戻る.

反復法で, STEP 2 の内容が最急降下法に合った内容に変わっただけ

最急降下法の収束速度

(例) 凸2次関数： $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^t Q(\mathbf{x} - \mathbf{x}^*)$ (ただし, Q は正定値対称行列)

→ $\mathbf{x} = \mathbf{x}^*$ が最適解となり, $f(\mathbf{x}^*) = 0$

・これを最急降下法で求めていくことを考えると…

・ $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ の最適解 α を厳密に求めると, 最急降下法の \mathbf{d}_k は

$$\mathbf{d}_k = -\nabla f(\mathbf{x}_k) = -Q(\mathbf{x}_k - \mathbf{x}^*) \text{より}$$

最急降下法の収束速度

(例) 凸2次関数： $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^t Q(\mathbf{x} - \mathbf{x}^*)$ (ただし, Q は正定値対称行列)

→ $\mathbf{x} = \mathbf{x}^*$ が最適解となり, $f(\mathbf{x}^*) = 0$

• $f(\mathbf{x}_{k+1})$ を求めると

$$f(\mathbf{x}_{k+1}) = \left\{ 1 - \frac{\|\nabla f(\mathbf{x}_k)\|^4}{\nabla f(\mathbf{x}_k)^t Q \nabla f(\mathbf{x}_k) \cdot \nabla f(\mathbf{x}_k)^t Q^{-1} \nabla f(\mathbf{x}_k)} \right\} f(\mathbf{x}_k)$$

• ここで, Q の固有値： $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ とすると

$$\begin{aligned} & \nabla f(\mathbf{x}_k)^t Q \nabla f(\mathbf{x}_k) \cdot \nabla f(\mathbf{x}_k)^t Q^{-1} \nabla f(\mathbf{x}_k) \\ & \leq \lambda_n \|\nabla f(\mathbf{x}_k)\|^2 \cdot \frac{1}{\lambda_1} \|\nabla f(\mathbf{x}_k)\|^2 = \frac{\lambda_n}{\lambda_1} \|\nabla f(\mathbf{x}_k)\|^4 \end{aligned}$$

最急降下法の収束速度

(例) 凸2次関数： $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^*)^t Q(\mathbf{x} - \mathbf{x}^*)$ (ただし, Q は正定値対称行列)

→ $\mathbf{x} = \mathbf{x}^*$ が最適解となり, $f(\mathbf{x}^*) = 0$

・つまり, Q の最小固有値を λ_1 , 最大固有値を λ_n とするとき,

最急降下法の収束速度の意味

- 最急降下法の収束速度と等高線・点列生成との関連性

機械学習への応用①

座標降下法とブースティング

- T個のデータ： $(\mathbf{a}_1, b_1), (\mathbf{a}_2, b_2), \dots, (\mathbf{a}_T, b_T)$, $\mathbf{a}_t \in \mathbb{R}^d, b_t \in \{+1, -1\}$
(つまり, 入力 \mathbf{a}_t が与えられたときに, 1か-1を b_t の値として出力する)
- 新たなデータ： $\mathbf{a} \in \mathbb{R}^d$ が得られたとき, 出力 $b \in \{+1, -1\}$ を以下の関数
 $H_x(\mathbf{a}) = \sum_{i=1}^n x_i h_i(\mathbf{a})$ の符号で判定 $\rightarrow H_x(\mathbf{a}) > 0$ なら $b = 1, H_x(\mathbf{a}) < 0$ なら $b = -1$
(ここで, $h_1(\mathbf{a}), h_2(\mathbf{a}), \dots, h_n(\mathbf{a})$ は $\{+1, -1\}$ の値をとる関数(基底関数とよぶ))
- 単純な基底関数を逐次的に複数組み合わせ、高い予測精度を達成する判別器を構成する方法：**ブースティング**

機械学習への応用①

- 新たなデータ： $\mathbf{a} \in \mathbb{R}^d$ が得られたとき，出力 $b \in \{+1, -1\}$ を以下の関数
 $H_x(\mathbf{a}) = \sum_{i=1}^n x_i h_i(\mathbf{a})$ の符号で判定 $\rightarrow H_x(\mathbf{a}) > 0$ なら $b = 1$, $H_x(\mathbf{a}) < 0$ なら $b = -1$
- Q： $H_x(\mathbf{a})$ の係数 $\mathbf{x} = (x_1, x_2, \dots, x_n)$ を観測データからどのように定めるか？
 - $\rightarrow b_t H_x(\mathbf{a}_t) > 0$ なら，関数 $H_x(\mathbf{a})$ で正しく判定できているので，できるだけ多くの (\mathbf{a}_t, b_t) に対して， $b_t H_x(\mathbf{a}_t) > 0$ が成り立つことが重要
- 損失関数を導入： $f(\mathbf{x}) = \sum_{t=1}^T e^{-b_t H_x(\mathbf{a}_t)}$
 - \rightarrow これを最小化するとき，座標降下法を学習アルゴリズムとして用いる
(アダブースト(Adaboost))

機械学習への応用①

Q : どのようにして座標降下法を用いるか？

A : $f(\mathbf{x}) = \sum_{t=1}^T e^{-b_t H_{\mathbf{x}}(\mathbf{a}_t)}$ の勾配を計算する. (このとき $w_t = e^{-b_t H_{\mathbf{x}}(\mathbf{a}_t)}$ とする)

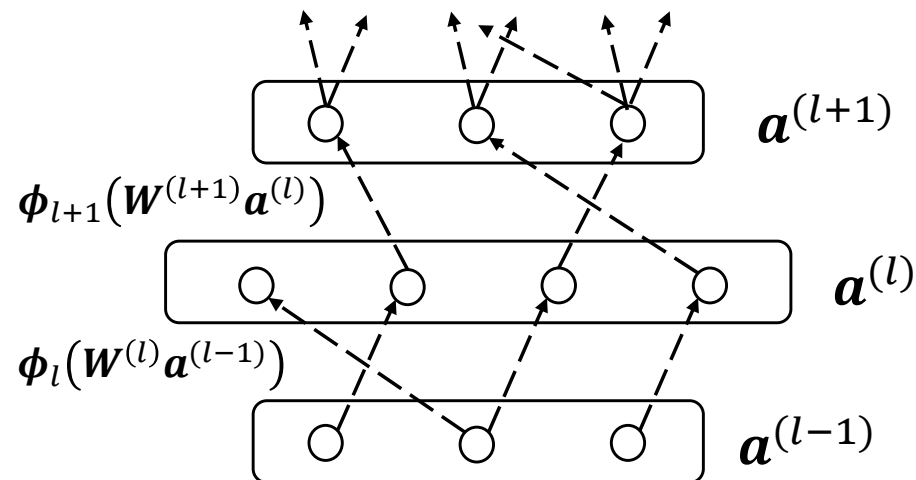
$$\frac{\partial f}{\partial x_i}(\mathbf{x}) = - \sum_{t=1}^T w_t b_t h_i(\mathbf{a}_t) = 2 \sum_{t=1}^T w_t \mathbf{1}[b_t \neq \pm h_i(\mathbf{a}_t)] - \sum_{t=1}^T w_t$$

(ここで, $\mathbf{1}[A]$ は A が真なら 1, 偽なら 0 をとる定義関数)

- よって, 関数値が減少する座標軸の方向は, $\pm h_1, \pm h_2, \dots, \pm h_n$ の中で, $\sum_{t=1}^T w_t \mathbf{1}[b_t \neq \pm h_i(\mathbf{a}_t)]$ を最小にする方向に対応(重み付き誤り数と解釈可能)
(= 座標降下法と相性が良い)

機械学習への応用②

- **ニューラルネットワーク**(今回は**多層パーセプトロン**)
- N個のデータ： $(\mathbf{a}_1, \mathbf{b}_1), (\mathbf{a}_2, \mathbf{b}_2), \dots, (\mathbf{a}_N, \mathbf{b}_N)$ において, \mathbf{a}_i と \mathbf{b}_i の関係を学習
→ 新たなデータ \mathbf{a} が入力されたときに, 出力 \mathbf{b} を精度よく予測できるか？



機械学習への応用②

多層パーセプトロン

・ 学習するパラメータ： $W^{(l)}, l = 1, 2, \dots, L$

・ 入力 a , 出力 b の関係 ($a^{(0)} = a$ として)

$$a^{(l)} = \phi_l(W^{(l)} a^{(l-1)}), l = 1, 2, \dots, L,$$

$$b = a^{(L)} = \phi_L \left(W^{(L)} \phi_{L-1} \left(W^{(L-1)} \dots W^{(2)} \phi_1 \left(W^{(1)} a^{(0)} \right) \dots \right) \right) = \Phi(a, W)$$

$$(W = (W^{(1)}, \dots, W^{(L)}))$$

ここで, $\phi_l(z) = (\psi(z_1), \psi(z_2), \dots, \psi(z_{D_l})) = (\tanh(z_1), \tanh(z_2), \dots, \tanh(z_{D_l}))$
 などが用いられる.

機械学習への応用②

- 入出力における2乗誤差を最小化する

$$f(\mathbf{W}) = \frac{1}{2N} \sum_{m=1}^N \|\Phi(\mathbf{a}_m, \mathbf{W}) - \mathbf{b}_m\|^2$$

- 探索方向を求めるためには、勾配の計算が必要

→ 誤差を出力層(最後の方)から入力層に伝搬させているように解釈できる

(**誤差逆伝搬法**(バック・プロパゲーション))

- $\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)} \mathbf{a}^{(l)} = \mathbf{W}^{(l+1)} \phi_l(\mathbf{u}^{(l)})$ を設定

機械学習への応用②

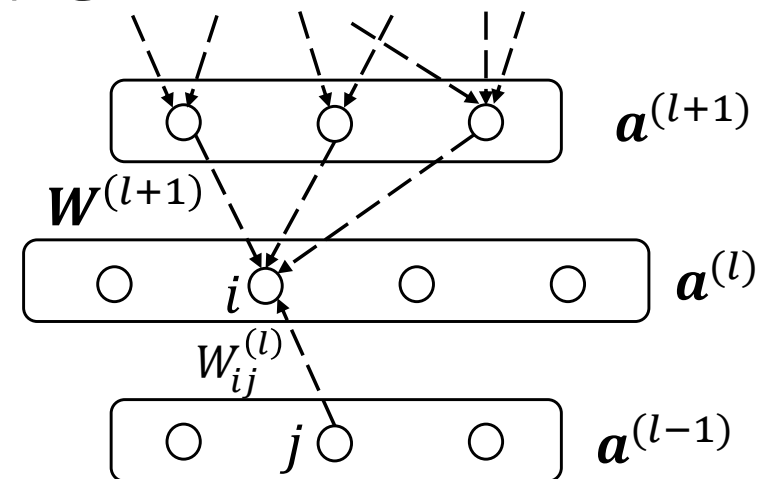
- 次の微分を考える ($\mathbf{u}^{(l+1)} = \mathbf{W}^{(l+1)} \boldsymbol{\phi}_l(\mathbf{u}^{(l)})$ を利用して)

$$\frac{\partial \Phi_d}{\partial u_i^{(l)}} = \sum_k \frac{\partial \Phi_d}{\partial u_k^{(l+1)}} \frac{\partial u_k^{(l+1)}}{\partial u_i^{(l)}} = \sum_k \frac{\partial \Phi_d}{\partial u_k^{(l+1)}} W_{ki}^{(l+1)} \psi'(u_i^{(l)})$$

つまり, $l+1$ 層目の微分から l 層目の微分が求められる。

- また, $\mathbf{u}^{(l)} = \mathbf{W}^{(l)} \mathbf{a}^{(l-1)}$ も利用すると

$$\frac{\partial \Phi_d}{\partial W_{ij}^{(l)}} = \sum_k \frac{\partial \Phi_d}{\partial u_k^{(l)}} \frac{\partial u_k^{(l)}}{\partial W_{ij}^{(l)}} = \frac{\partial \Phi_d}{\partial u_i^{(l)}} a_j^{(l-1)}$$



機械学習への応用②

- 以上の結果から

$$\frac{\partial f}{\partial W_{ij}^{(l)}}(\mathbf{W}) = \frac{1}{N} \sum_{m=1}^N \sum_{d=1}^D (\Phi_d(\mathbf{a}_m, \mathbf{W}) - b_{m,d}) \frac{\partial \Phi_d}{\partial W_{ij}^{(l)}}(\mathbf{a}_m, \mathbf{W})$$

を f の最適化に用いる.

- 確率的最適化のアイデアを用いると, $W_{ij}^{(l)}$ は以下のように表記できる

$$W_{ij}^{(l)} \leftarrow W_{ij}^{(l)} - \varepsilon_t \cdot \sum_{d=1}^D (\Phi_d(\bar{\mathbf{a}}_t, \mathbf{W}) - \bar{b}_{t,d}) \frac{\partial \Phi_d}{\partial W_{ij}^{(l)}}(\bar{\mathbf{a}}_t, \mathbf{W})$$

(ここで, $(\bar{\mathbf{a}}_t, \bar{\mathbf{b}}_t)$ は t ステップ目で与えられたデータ, ε_t は適当な正の数)

今回のまとめ

- 勾配法の基礎として, 凸2次関数を例とした最急降下法の収束速度を考察
 - 最急降下法を1次収束であるため, 最適解付近ではやはり収束が鈍化
 - ニュートン法や準ニュートン法の利用
- 勾配法の応用として, アダブーストと座標降下法との関わり, 多層パーセプトロンと勾配法との関わりを紹介(詳細までは踏み込まない予定. 詳細を知りたい方は, それぞれの専門書を参照すること)
 - 先端手法の影の主演は, 実は最適化