

制約なし最適化の解法②

山下信雄

講義内容

1. 具体的な直線探索法

- 最急降下法
- ニュートン法
- 準ニュートン法

2. 信頼領域法

3. 確率的勾配法

直線探索法

ステップ1： 降下方向 d^k を求める.

ステップ2： ステップ幅 t_k を定める.

ステップ3： $x^{k+1} = x^k + t_k d^k$ として、ステップ1へ

適当な仮定のもとで、停留点($\nabla f(x^*) = 0$)
に大域的収束する.

最急降下法

最急降下方向：

$$d^k = -\nabla f(x^k) = -\textcolor{red}{I}\nabla f(x^k)$$

$$\nabla f(x^k)^\top d^k = -\nabla f(x^k)^\top \nabla f(x^k) = -\|\nabla f(x^k)\|^2$$

$$\|d^k\| = \|\nabla f(x^k)\|$$



大域的収束の仮定を満たす

最急降下法の性質

良いところ:

- 各反復の計算が簡単
- 大域的収束する

悪いところ:

- 収束が遅い (高々一次収束)

ニュートン法

ニュートン方向：

$$d^k = -\nabla^2 f(x^k)^{-1} \nabla f(x^k)$$

ニュートン方向は次の最小化問題(部分問題)の最適解

$$\min m_k(d)$$

$$\text{s.t. } d \in R^n$$

$$0 = \nabla m_k(d) = \nabla f(x^k) + \nabla^2 f(x^k) d$$

$$d = -\nabla f(x^k)^{-1} \nabla f(x^k)$$

ここで $m_k(d) = f(x^k) + \nabla f(x^k)^\top d + \frac{1}{2} d^\top \nabla^2 f(x^k) d$

$f(x^k + d)$ の2次近似

ニュートン法の収束性

ニュートン方向は一般には降下方向にならない.

⇒ 直線探索法では大域的収束しないことがある.

大域的収束と2次収束するための十分条件:

次の不等式をみたす定数 $c_1, c_2 > 0$ が存在する.

$$c_1 \|v\|^2 \leq v^\top \nabla^2 f(x)^{-1} v \leq c_2 \|v\|^2 \quad \forall v, x \in R^n$$

➡

$$\begin{aligned} \nabla f(x^k)^\top d^k &= -\nabla f(x^k)^\top \nabla^2 f(x^k)^{-1} \nabla f(x^k) \leq -c_1 \|\nabla f(x^k)\|^2 \\ \|d^k\| &= \|\nabla^2 f(x^k)^{-1} \nabla f(x^k)\| \leq c_2 \|\nabla f(x^k)\| \end{aligned}$$

ニュートン法の2次収束性

ステップ幅は $t_k = 1$ とする.

$$\begin{aligned}\|x^{k+1} - x^*\| &= \|x^k - \nabla^2 f(x^k)^{-1} \nabla f(x^k) - x^*\| \\ &= \|\nabla^2 f(x^k)^{-1} (\nabla f(x^*) - \nabla f(x^k) - \nabla^2 f(x^k)(x^k - x^*))\| \\ &\leq c_2 \|\nabla f(x^*) - \nabla f(x^k) - \nabla^2 f(x^k)(x^k - x^*)\|\end{aligned}$$

f を2回連続的に微分可能とすると

$$\nabla f(x^*) = \nabla f(x^k) + \nabla^2 f(x^k)(x^* - x^k) + O(\|x^* - x^k\|^2)$$

このとき $\|x^{k+1} - x^*\| \leq O(\|x^k - x^*\|^2)$

ニュートン法の性質のまとめ

よいところ

- 収束すれば、収束は速い(2次収束)

わるいところ

- 直線探索法では大域的収束しない
- 各反復で線形方程式(ニュートン方程式)を解かなければならない.

$$\nabla f(x^k) + \nabla^2 f(x^k)d = 0$$

一般に計算量は $O(n^3)$

準ニュートン法

準ニュートン方向：

$$d^k = -(\mathbf{B}_k)^{-1} \nabla f(x^k) = -\mathbf{H}_k \nabla f(x^k)$$

$\mathbf{H}_k = \mathbf{B}_k^{-1}$ が正定値行列のとき

$$\nabla f(x^k)^\top d^k = -\nabla f(x^k)^\top \mathbf{H}_k \nabla f(x^k) < 0$$

\Rightarrow 降下方向

近似ヘッセ行列の望ましい条件

近似ヘッセ行列 B_k とその逆行列 H_k

- H_k が正定値 \Rightarrow 大域的収束
- $H_k \approx \nabla^2 f(x^k)^{-1} \Rightarrow$ 速い収束

近似ヘッセ行列のセカント条件

$\nabla f(x_k)$ のテーラー展開を考えると,

$$\nabla f(x_{k+1}) - \nabla f(x_k) \approx \nabla^2 f(x_{k+1})(x_{k+1} - x_k)$$

を得る.

$$s_k = x_{k+1} - x_k, \quad y_k = \nabla f(x_{k+1}) - \nabla f(x_k)$$

とする. 近似ヘッセ行列が

$$y_k = B_{k+1} s_k \quad \text{or} \quad H_{k+1} y_k = s_k$$

セカント条件

を満たせば, ヘッセ行列と似た性質をもつ.

⇒ 速い収束が期待できる.

Broyden–Fletcher–Goldfarb–Shanno (BFGS)更新

BFGS更新：

$$\begin{aligned} B_{k+1} &= B_k - \frac{B_k s_k (B_k s_k)^\top}{s_k^\top B_k s_k} + \frac{y_k y_k^\top}{s_k^\top y_k} \\ H_{k+1} &= H_k - \frac{H_k y_k s_k^\top + s_k (H_k y_k)^\top}{s_k^\top y_k} + \left(1 + \frac{y_k^\top H_k y_k}{s_k^\top y_k} \right) \frac{s_k s_k^\top}{s_k^\top y_k} \end{aligned}$$

- セカント条件を満たす.

$$B_{k+1} s_k = B_k s_k - \frac{B_k s_k (B_k s_k)^\top s_k}{s_k^\top B_k s_k} + \frac{y_k y_k^\top s_k}{s_k^\top y_k} = y_k$$

- H_k が正定値行列で, $s_k^\top y_k > 0$ であれば正定値行列となる.

準ニュートン法の性質

適当な条件のもとで

- ・ 大域的収束する.
- ・ 超一次収束する.
- ・ H_{k+1} の更新と, ベクトル $\nabla f(x^k)$ と H_{k+1} の掛け算は $O(n^2)$

$$\mathbf{H}_{k+1} = \mathbf{H}_k - \frac{\mathbf{H}_k \mathbf{y}_k \mathbf{s}_k^\top + \mathbf{s}_k (\mathbf{H}_k \mathbf{y}_k)^\top}{\mathbf{s}_k^\top \mathbf{y}_k} + \left(1 + \frac{\mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k}{\mathbf{s}_k^\top \mathbf{y}_k} \right) \frac{\mathbf{s}_k \mathbf{s}_k^\top}{\mathbf{s}_k^\top \mathbf{y}_k}$$

注意：行列と行列の掛け算は $O(n^3)$

準ニュートン法の問題点

BFGS更新：
$$H_{k+1} = H_k - \frac{H_k y_k s_k^\top + s_k (H_k y_k)^\top}{s_k^\top y_k} + \left(1 + \frac{y_k^\top H_k y_k}{s_k^\top y_k} \right) \frac{s_k s_k^\top}{s_k^\top y_k}$$

$s_k = (1, 1, \dots, 1)^\top$ のとき,

$$s_k s_k^\top = \begin{bmatrix} 1 & \cdots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \cdots & 1 \end{bmatrix}$$

ヘッセ行列が疎であっても、 H_{k+1} は密な行列となる.

⇒ 大規模な問題には使えない！

補足：スパース性(疎性)について

スパース性：ベクトルや行列の成分がほとんど0になる性質

* 0の成分は、足し算や掛け算で計算する必要がない。

例： $f(x) = \sum_{i=1}^{n-1} (x_i - x_{i+1})^2 + x_n^2$

ヘッセ行列は3重対角行列. $\nabla^2 f(x) = \begin{pmatrix} 2 & -2 & 0 & \cdots & 0 \\ -2 & 4 & -2 & \ddots & \vdots \\ 0 & -2 & 4 & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & -2 \\ 0 & \cdots & 0 & -2 & 4 \end{pmatrix}$

ニュートン方程式は $O(n)$ で計算できる。

準ニュートン法 + 直線探索

[良いところ]

- 大域的収束かつ超1次収束
- 実装が簡単

[悪いところ]

- H_k は非ゼロ要素がほとんどないため,
大規模な問題(数万変数以上)には適用できない.

講義内容

1. 具体的な直線探索法

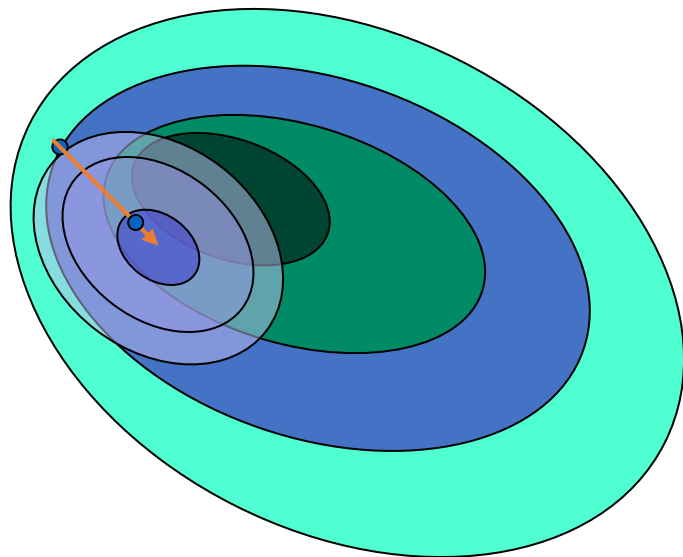
- 最急降下法
- ニュートン法
- 準ニュートン法

2. 信頼領域法

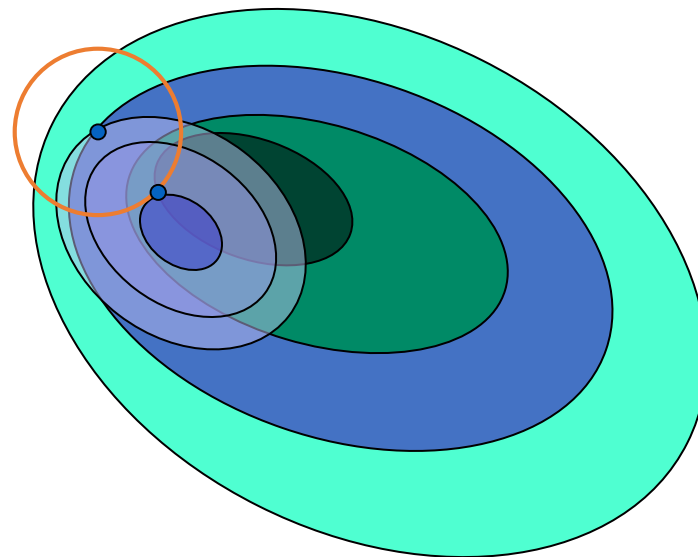
3. 確率的勾配法

大域的収束性を持たせる技術

- 直線探索法と信頼領域法 -



直線探索法



信頼領域法

信頼領域法

モデル関数のコンパクト集合上での最小解を探索報告とする.

$$\begin{aligned} \min \quad & m_k(d) \\ \text{s.t.} \quad & \|d\| \leq \Delta_k \end{aligned}$$

- この部分問題には最適解 d^k が存在する.
- Δ_k を信頼半径という.
- $B(x_k; \Delta) = \{x \mid \|x - x_k\| \leq \Delta\}$ を信頼領域という.

モデル関数は2次近似関数なので, 信頼半径が小さいときは,

$$f(x^{\text{tri}}) < f(x_k)$$

が成り立つ. ただし, $x^{\text{tri}} = x^k + d^k$

信頼半径の更新

モデルの信頼性を測る指標

$$\begin{aligned}\rho_k &= \frac{\text{目的関数の減少量}}{\text{モデル関数の減少量}} \\ &= \frac{f(x_k) - f(x^{\text{tri}})}{m_k(0) - m_k(d^k)}\end{aligned}$$

$$m_k(0) = f(x^k)$$

$$\begin{aligned}m_k(d^k) &= f(x^k) + \nabla f(x^k)^\top d^k + \frac{1}{2}(d^k)^\top \nabla^2 f(x^k) d^k \\ &= f(x^k) + \nabla f(x^k)^\top (x^{\text{tri}} - x^k) + \frac{1}{2}(x^{\text{tri}} - x^k)^\top \nabla^2 f(x^k)(x^{\text{tri}} - x^k) \\ &\approx f(x^{\text{tri}})\end{aligned}$$

【信頼半径の更新】

$$\Delta_{k+1} = \begin{cases} \alpha_1 \|x^{\text{tri}} - x_k\| & \text{if } \rho_k < \eta_1 \\ \Delta_k & \text{if } \eta_1 \leq \rho_k < \eta_2 \\ \max\{\alpha_2 \|x^{\text{tri}} - x_k\|, \Delta_k\} & \text{if } \eta_2 \leq \rho_k \end{cases}$$

$$\alpha_1 < 1 < \alpha_2$$

$$0 < \eta_1 < \eta_2 < 1$$

ニュートン法 + 信頼領域法

[問題点]

- ヘッセ行列の計算
⇒ 自動微分
- 非凸の部分問題の求解
⇒ 高速の近似解法の開発
(線形方程式を1回解く程度の計算時間)

大規模な問題では,
準ニュートン法 + 直線探索
よりも優れている.

講義内容

1. 具体的な直線探索法

- 最急降下法
- ニュートン法
- 準ニュートン法

2. 信頼領域法

3. 確率的勾配法

大規模な最適化問題

$$\begin{aligned} \min \quad & \frac{1}{M} \sum_{i=1}^M f_i(x) \\ \text{s.t.} \quad & x \in R^n \end{aligned}$$

大規模となるもの

- 変数の数 n
- 目的関数を構成する関数の数 M

応用例： データ解析

$$\min \frac{1}{M} \sum_{i=1}^M \theta(x, a^i)$$

- a^i はデータ. θ は損失関数
- $\frac{1}{M} \sum_{i=1}^M \theta(x, a^i)$ は経験的損失とよばれる.

想定する問題

- M が大きい, 大規模な問題
 - 関数値や勾配の計算が大変
- 似たような関数 f_i (データ a^i)が多い

例：ディープラーニング, SVM, L1正則化問題, etc

確率勾配法

以下では、目的関数を $f(x) = \frac{1}{M} \sum_{i=1}^M f_i(x)$ とする.

最急降下法 : $x^{k+1} = x^k - t_k \nabla f(x^k)$

$\nabla f_i(x^k)$ の計算量が $O(n)$ のとき,

1回の反復の計算量は $O(nM)$

確率的勾配法 : $i_k \in \{1, 2, \dots, M\}$ をランダムに選ぶ

$$x^{k+1} = x^k - t_k \nabla f_{i_k}(x^k)$$

1回の反復の計算量は $O(n)$

いろいろな呼び方

- Stochastic gradient descent method

i_k をランダムにとってくるとき

ニューラルネットワークの学習では誤差逆伝播法という

- Incremental gradient method

i_k を $\{1, 2, \dots, M\}$ から順番に取ってくるとき

- Online gradient descent method

逐次的に 暫定解 \mathbf{x}^k を ``実行(利用)'' するとき

確率的勾配 $g^k = \nabla f_{i_k}(x^k)$ の特徴

- i_k を確率 $p_i = \frac{1}{M}$ で取ってきたとき, 探索方向の期待値は目的関数の最急降下方向と一致する.

$$E_i[g^k] = E_i[\nabla f_{i_k}(x^k)] = \sum_{i=1}^M p_i \nabla f_i(x^k) = \frac{1}{M} \sum_{i=1}^M \nabla f_i(x^k) = \nabla f(x^k)$$

- 分散 $E_i\left[\left(g^k - E[\nabla f_i(x^k)]\right)^2\right]$ は 0 とならない.

→ 分散が小さいとき, 最急降下法に近づく

ステップ幅 t_k を固定したとき

例： $\min \frac{1}{2}(x-1)^2 + \frac{1}{2}(x+1)^2$

$$f_1(x) = \frac{1}{2}(x-1)^2, f_2(x) = \frac{1}{2}(x+1)^2$$

最適解は $x^* = 0$

$x^0 = -0.5, t_k = 0.5$ とした最急降下法

$$x^1 = x^0 - t_k f'(x^0) = -0.5 - 0.5 \times \{(-1.5) + 0.5\} = 0$$

$x^0 = -0.5, t_k = 0.5$ とした確率的勾配降下法

$$i_k=1: x^1 = x^0 - t_k f_1'(x^0) = -0.5 - 0.5 \times \{-1.5\} = 0.25$$

$$i_k=2: x^1 = x^0 - t_k f_2'(x^0) = -0.5 - 0.5 \times \{0.5\} = -0.75$$

ステップ幅を固定すると、一般に収束しない

ステップ幅の取り方と収束

Diminishing Rule:

$$t_k \rightarrow 0, \sum_{k=1}^{\infty} t_k = +\infty$$

Diminishing Ruleを用いた確率勾配降下法は大域的収束する.

- f_i が凸であれば最適解に収束. そうでないときは停留点に収束.