

0. はじめに

このシケプリでは、 $\nabla f(\mathbf{x})^t$ のことを $f'(\mathbf{x})$ と表記してあります。これは、手持ちの数学の教科書でこのようになっていた（実際はベクトルが太字にさえなっていなかった）こともありますが、僕が妥協したことが一番の原因だったりします。ご了承下さい。また、1. は微分の復習なので読み飛ばしても問題ないと思います。

1. 微分の復習

1. では $f(\mathbf{x} + t\mathbf{d}) = f(\mathbf{x}) + tf'(\mathbf{x})\mathbf{d} + \frac{t^2}{2}\mathbf{d}^t\nabla^2 f(\mathbf{x})\mathbf{d} + O(t^3)$ の導出をメインとしています。これを導出するために、教科書の丸写しみたいなことをしてページ数を増やしてしまいました。ただ先生の授業とプリントが、シケプリなんていらなくらい丁寧で、序盤で唯一不親切な気がしたのが上の式の導出だったので… と言いつつ。

定義 関数 $f: \mathbb{R} \rightarrow \mathbb{R}$ が点 x で微分可能とは次式を満たす $c \in \mathbb{R}$ が存在することである。

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = c \quad (1)$$

この時 c を f の x における微分係数または導値といい、 $c = f'(a)$ と表す。

これを今更書く必要は無かったかもしれませんが、一応挙げておきました。次にこれを多変数関数に拡張します。

定義 関数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ が点 \mathbf{x} で微分可能とは次式を満たす定ベクトル $\mathbf{c} = (c_1, \dots, c_n) \in \mathbb{R}^n$ が存在することである。

$$\lim_{|\mathbf{h}| \rightarrow 0} \frac{f(\mathbf{x} + \mathbf{h}) - f(\mathbf{x}) - \mathbf{c}\mathbf{h}}{|\mathbf{h}|} = 0 \quad (2)$$

この時 \mathbf{c} を f の \mathbf{x} における微分係数または導値といい、 $\mathbf{c} = f'(\mathbf{x})$ (授業でいうところの $\nabla f(\mathbf{x})$ です) と表す。

\mathbf{c} だけ \mathbf{x} や \mathbf{h} と違って横ベクトルです。また \mathbf{h} の方向によらず、その大きさが 0 に向かってくれさえすれば \mathbf{c} がただ一つに決まることが、多変数関数の場合の微分のポイントです。

定理 f が \mathbf{x} で微分可能である時、 f は \mathbf{x} で各座標について偏微分可能で、

$$f'(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right) \quad (3)$$

が成り立つ。証明略。

数理科学 でやったと思いますが、 $f'(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$ のように「定義」しないのは、仮に $\left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)$ が存在したとしても、(2) 式の意味で微分可能かどうかは不明だからです。このことから、微分可能のことをあえて全微分可能と表現することもあるようです。ただし、最適化手法で出てくるような関数はほとんど全微分可能なので、いきなりこの定理によって微分係数を与えてもさほど問題ないものと思います。

次の定義と定理はあまり最適化手法とあまり関係ありませんが、その次の連鎖律のために必要なので乗せておきました。

定義 関数 $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ($f = (f_1, \dots, f_m)^t$) が点 x で微分可能とは次式を満たす行列 $M \in \mathbb{R}^{m \times n}$ が存在することである。

$$\lim_{|h| \rightarrow 0} \frac{f(x+h) - f(x) - Mh}{|h|} = 0 \quad (4)$$

この時 M を f の x における微分係数または導値といい、 $M = f'(x)$ と表す。

定理 f が x で微分可能である時、 f_i ($i = 1, \dots, m$) は x で各座標について偏微分可能で、

$$f'(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad (5)$$

が成り立つ。証明略。

定理 (連鎖律) 二つの関数 $g: \mathbb{R}^n \rightarrow \mathbb{R}^m$ と $f: \mathbb{R}^m \rightarrow \mathbb{R}^p$ が合成可能で両方とも微分可能であるとき、その合成関数の微分に関して次式が成り立つ。

$$(f \circ g)'(x) = f'(y)g'(x), \quad (y = g(x)) \quad (6)$$

証明略

合成関数の微分が合成される関数の微分 (この場合行列) の積で表されるということです。この定理は、名前こそあまり知られていないものの、非常に使用頻度が高いと思われます。

次に $f: \mathbb{R}^n \rightarrow \mathbb{R}$ に関して、 $f(x+td)$ を t の関数と見た時 (他を定ベクトルと考えた時) の微分を考えます。この関数は $g(t): \mathbb{R} \rightarrow \mathbb{R}$ ですから、形は非常になじみのある関数です。実際 t で微分してみると、

$$g'(t) = f'(x+td)d \quad (7)$$

となります。ここで既に連鎖律が使われています。(6) での $g = x+td$ とし、 t を引数にベクトルを返す関数だと考えています。そしてこれを t で微分すると (5) 式より n 行 1 列行列として d を得ます。次に二階の導関数を求めるわけですが、 $f'(x+td)$ は横ベクトルであるため、うまく微分できません。そこで転置します。 $g'(t) = d^t \nabla f(x+td)$ として、これを行列の積と考えて微分します。 d^t は変化しないのでその後ろだけ考えればよくて、

$$g''(t) = d^t \nabla^2 f(x+td)d \quad (8)$$

となります。 $\nabla^2 f(x)$ はヘシアンですが、これの定義式を `tex` で書くのは心が折れてしまうので、ノートの方を参照下さい。この式変形で解かるとおりヘシアンは、 $\nabla f'(x)$ という多変数ベクトル値関数の、(4) 式の意味での微分係数として現れます。二階までの導関数を求めたことにより、 f を 2 次まで展開することができます。その結果は

$$\begin{aligned} f(x+td) &= g(t) = g(0) + tg'(0) + \frac{t^2}{2}g''(0) + O(t^3) \\ &= f(x) + tf'(x)d + \frac{t^2}{2}d^t \nabla^2 f(x)d + O(t^3) \end{aligned} \quad (9)$$

となります。

$g'(t) = d^t \nabla f(\mathbf{x} + t\mathbf{d})$ をよく見ると、一階微分を表す演算子が、 $(d_1 \frac{\partial}{\partial x_1} + \dots, d_n \frac{\partial}{\partial x_n})$ となっていることがわかります。よって、

$$g^{(n)}(t) = (d_1 \frac{\partial}{\partial x_1} + \dots, d_n \frac{\partial}{\partial x_n})^n f(\mathbf{x} + t\mathbf{d}) \quad (10)$$

を得ます。これは n 回微分でも簡単に表せる利点を持ちますが、(9) の形式だと関数の極値問題を行列の正定値問題に帰着できるので主にこちらを使います。

2. 無制約最適化問題

一応ここからが最適化手法の本題です。だいたいノートの順番に則って書いたのですが、ノートを取ってない人には役に立つかもしれませんが、しっかりノートを取った人には無用かもしれません…。このシケプリはノートを取ってない人用の試験対策向けという感じで作ってあります。そのくせ、ノート参照とか出てきますが、そこらへんは友達に借りてください。

2.1 種々の言葉の定義

定義 関数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ に対して、 $\min f(\mathbf{x})$ を考えるような問題を無制約最適化問題という。 $f(\mathbf{x}) \geq f(\mathbf{x}^*), \forall \mathbf{x} \in \mathbb{R}^n$ であるとき、 \mathbf{x}^* を大域的最適解という。

$f(\mathbf{x}) \geq f(\mathbf{x}^*), \forall \mathbf{x} \in N(\mathbf{x}^*)$ であるとき、 \mathbf{x}^* を局所最適解という。ただし、 $N(\mathbf{x}^*)$ とは、 \mathbf{x}^* の近傍である。

$\min\{f(\mathbf{x}) | \mathbf{x} \in \mathbb{R}^n\}$ を f の最適値といい、 f^* と表す。

$f'(\mathbf{x}) = 0$ のとき、 \mathbf{x} を停留解という。

極小でも極大でもない停留解を鞍点という。

定理というほどのものではありませんが、 \mathbf{x}^* が大域的最適解 $\Leftrightarrow f(\mathbf{x}^*) = f^*$ です。

例 $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ に対して $g(\mathbf{x}) = 0$ の解を求めることは、無制約最適化問題 $\min_{\mathbf{x}} g(\mathbf{x})^t g(\mathbf{x})$ に帰着される。

例 ペナルティ法 ノート参照。

定義 対称行列 A に対して、 A の固有値が全て 0 より大きいとき、 A を正定値といい、 $A \succ 0$ と表す。 $(\mathbf{x}^t A \mathbf{x} > 0, (\forall \mathbf{x} \neq 0))$ と同値)

対称行列 A に対して、 A の固有値が全て 0 以上のとき、 A を半正定値といい、 $A \succeq 0$ と表す。 $(\mathbf{x}^t A \mathbf{x} \geq 0, (\forall \mathbf{x}))$ と同値)

括弧内のことは証明が必要ですが省略します。また、ノートに正定値行列に関する公式があるので見ておいて下さい。(例えば、「 A が正定値ならば逆行列が存在する」など)

定理 f が微分可能であるとき、 $\mathbf{x} \in \mathbb{R}^n$ が (局所) 最適解 $\Rightarrow f'(\mathbf{x}) = 0$

f が二回微分可能であるとき、 $\mathbf{x} \in \mathbb{R}^n$ が (局所) 最適解 $\Rightarrow f'(\mathbf{x}) = 0, \nabla^2 f(\mathbf{x})$ は半正定値
これらの定理は一般に逆は成り立ちません。

定理 f が二回微分可能で、 $f'(\mathbf{x}) = 0, \nabla^2 f(\mathbf{x}) \succ 0 \Leftrightarrow \mathbf{x}$ は f の局所最適解。

こちらは、最適解の十分条件です。

定義 $\mathbf{x}_{k+1} = \mathbf{g}_k(\mathbf{x}_k)$ として点列 $\{\mathbf{x}_k\}$ を作るようなアルゴリズムを反復法と呼ぶ。
基礎的な反復法は、 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ の形をしています。

2.2 最急降下法

以下では $f(\mathbf{x}_k) = f_k, f'(\mathbf{x}_k) = f'_k$ と置いています。アルゴリズムについては、プリント「最急降下法の概要」に詳しくでているのでそちらをご覧ください。最急降下法の \mathbf{d}_k は f の \mathbf{x}_k におけるテイラー展開による 1 次モデル

$$f_k^L(\mathbf{x}_k + \mathbf{d}) = f_k + f'_k \mathbf{d} \quad (11)$$

((9) 式で $t = 1$ と置いて 2 次以下を無視する。) を最小にする向きになります。また、 \mathbf{x}_k が解に近づくにつれて \mathbf{d}_k も勝手に小さくなってくれますが、だからといって α_k を定数と置いてしまうと収束しません。 α_k の求め方は「正確な直線探索」「Armijo の基準」「Wolfe の基準」の 3 種類あります。「正確な直線探索」はプリントに出ています。後の 2 つは特に覚えなくても多分問題は無いと思います。一応ノートに書いてあるはずですが。

2.3 Newton 法

テイラー展開による f の 2 次モデルは、

$$f_k^Q(\mathbf{x}_k + \mathbf{d}) = f_k + f'_k \mathbf{d} + \frac{1}{2} \mathbf{d}^t \nabla^2 \mathbf{f}_k \mathbf{d} \quad (12)$$

((9) 式で $t = 1$ と置いて 3 次以下を無視する。) となります。ここで、求めなくてはいけないのは \mathbf{d} なので、 \mathbf{d} を変数と見た勾配 ∇ を考えます。そしてこの微分のベクトルが 0 になったらそれを \mathbf{d}_k とします。まず、 f_k は \mathbf{d} に関係ありません。 $f'_k \mathbf{d}$ の \mathbf{d} に関する勾配ベクトルは $\nabla f'_k$ に、 $\frac{1}{2} \mathbf{d}^t \nabla^2 \mathbf{f}_k \mathbf{d}$ は $\nabla^2 \mathbf{f}_k \mathbf{d}$ となります (計算過程略)。これにより、 \mathbf{d} による勾配が 0 になることの条件として Newton 方程式

$$\nabla^2 \mathbf{f}_k \mathbf{d} = -\nabla f'_k \quad (13)$$

を得ます。Newton 方程式は解けるとは限らないし、減少しているかもわからないところが難点です。

定義 \mathbf{d} が点 \mathbf{x} における f の降下方向とは、 $f'(\mathbf{x})\mathbf{d} < 0$ が成り立つことである。

定理 $\nabla^2 \mathbf{f}_k \succ 0$ のとき、Newton 方程式は一意的に解けて、 \mathbf{d}_k は降下方向。(証明はノート参照)

2.4 収束の早さ

定義 $\{\mathbf{x}_k\}$ が \mathbf{x}^* に 1 次収束するとは、

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = M, 0 \leq M < 1 \quad (14)$$

が成り立つことをいう。超 1 次収束するとは

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0 \quad (15)$$

が成り立つことをいう。2 次収束するとは、1 次収束するうえでさらに

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} = M \geq 0 \quad (16)$$

が成り立つことをいう。

定理 最急降下法は 1 次収束 (プリント「最急降下法の収束速度について」参照)。Newton 法は 2 次収束 (ノート参照)

2.5 準 Newton 法

準 Newton 法では 2 次モデルを

$$\tilde{f}_k^Q(\mathbf{x}_k + \mathbf{d}) = f_k + f'_k \mathbf{d} + \frac{1}{2} \mathbf{d}^t B_k \mathbf{d} \quad (17)$$

とします。テイラー展開による 2 次モデルとは違うものなのでチルダがついています。ここで、 B_k の条件の 1 つは正定値となることです。 B_0 の条件はこれだけです。この条件が満たされれば、Newton 方程式によく似た式

$$B_k \mathbf{d} = -\nabla f_k \quad (18)$$

の解として、降下方向のベクトル \mathbf{d}_k を得ることができます。次に $k + 1$ での 2 次モデル

$$\tilde{f}_{k+1}^Q(\mathbf{x}_{k+1} + \mathbf{d}) = f_{k+1} + f'_{k+1} \mathbf{d} + \frac{1}{2} \mathbf{d}^t B_{k+1} \mathbf{d} \quad (19)$$

を考えます。これを変数を \mathbf{d} と考えて勾配を取ると、Newton 方程式を導出したときと同じ計算により、 $\nabla f_{k+1} + B_{k+1} \mathbf{d}$ を得ます。これが、 $\mathbf{d} = \mathbf{0}$ 、 $-\mathbf{x}_{k+1} + \mathbf{x}_k$ のときそれぞれ、 ∇f_{k+1} 、 ∇f_k に一致してくれれば 2 次モデルは精度の高いものといえます ($\nabla \tilde{f}_{k+1}^Q$ と ∇f が 2 点 \mathbf{x}_{k+1} 、 \mathbf{x}_k で等しくなるため)。前者は明らかに成り立ちますが、後者が成り立つことは

$$B_{k+1}(\mathbf{x}_{k+1} - \mathbf{x}_k) = \nabla f_{k+1} - \nabla f_k \quad (20)$$

と同値です。上の式であたえられる条件をセカント条件といいます。 B_1 以下では正定値であることの他に、この条件がつかます。これらの条件を満たすように $\{B_k\}$ を作るやりかたはいくつか知られていて、BFGS 公式、DFP 公式などと呼ばれています。公式はプリント「準 Newton 法の概要」に載っています。さらに、(18) 式を解くのは手間がかかるので、 $\{B_k\}$ を計算するのではなく、その逆行列 $\{H_k\}$ を計算しようという公式を H 公式と呼びます。多分こちらの方が実用的です。

定理 準 Newton 法は超 1 次収束する。(プリント「準 Newton 法の超 1 次収束性」参照)

2.6 信頼領域法

信頼領域法の2次モデルは

$$f_k^Q(\mathbf{x}_k + \mathbf{d}) = f_k + f'_k \mathbf{d} + \frac{1}{2} \mathbf{d}^t B_k \mathbf{d} \quad (21)$$

となります。(さっきはテイラー展開と違うとの理由でチルダをつけましたが今回は省略です。)そして更新公式は、 $\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{d}_k$ となります。 \mathbf{d}_k は制約条件付き最適化問題

$$\min f_k^Q(\mathbf{x}_k + \mathbf{d}) \text{ s.t. } \|\mathbf{d}\| \leq \Delta k \quad (22)$$

の最適解として与えられます。(条件式が閉空間を作っているため、必ず最適解を持ちます。) \mathbf{d}_* が部分問題の最適解であることの必要条件がノートに載っているため、そちらを参照してみてください。

$\{B_k\}$ の求め方は手持ちのノートに書いていなかったのですが、おそらく準Newton法と同じように求めるのだと思います(このようにすれば、 B_k は正定値です)。そして、まず領域の内部に解があると仮定して、 $-B_k^{-1} \nabla f_k$ を求め(B_k が正定値なら無制約下での最適解はこれ1つ)それがちゃんと条件の領域の内部にあったらそれを \mathbf{d}_k とし、領域から出てしまったら部分問題の解は領域上にあるので、ラグランジュ未定乗数法で領域上の解を求めるのだと思います。

また信頼領域法では、求めた \mathbf{d}_k が妥当かどうかチェックするようです。もし妥当じゃないと判定されたら、 Δk を小さくしてやり直すようです。これについてはノートを参照してください。

2.7 凸関数

定義 関数 $f: \mathbf{R}^n \rightarrow \mathbf{R}$ で次の条件をみたすものを凸関数という。

$$\forall \mathbf{x}, \mathbf{y} \in \mathbf{R}^n, \forall \lambda \in [0, 1] \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \geq f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \quad (23)$$

さらに次の条件をみたすものを狭義凸関数という。

$$\forall \mathbf{x}, \mathbf{y} \in \mathbf{R}^n, \forall \lambda \in (0, 1) \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) > f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \quad (24)$$

定理 f が微分可能であるとき、 f が凸であることの必要十分条件は

$$\forall \mathbf{x}, \mathbf{y} \in \mathbf{R}^n, f(\mathbf{y}) \geq f(\mathbf{x}) + f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \quad (25)$$

f が微分可能な凸関数のとき、 $\nabla f(\mathbf{x}) = \mathbf{0} \Rightarrow \mathbf{x}$ は f の大域的最適解

f が2回連続微分可能のとき、 f が凸関数 $\Leftrightarrow \nabla^2 f \succeq 0 (\forall \mathbf{x})$

例 2次関数 $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^t Q \mathbf{x} + \mathbf{p}^t \mathbf{x} + r$ において f が凸であることと Q が半正定値であることは同値

ところで、上の定理において凸を狭義凸に、 \geq, \succeq を $>, \succ$ に変えたものは少し定理が必要ではなく、2回連続微分可能で f が凸関数でも、 $\nabla^2 f \succ 0 (\forall \mathbf{x})$ は言えないはず。それ以外は問題ないと思います。