

サンプリング探索法

1. 概要

ソートされたデータの探索法としては、ほとんど二分探索法が用いられますが、他に
いか自分で考えてみました。

昇順にソートされたデータを対象として、x軸にデータの順番 ($x = 1 \sim n$) をとり、
y軸にデータの値をとった平面を考えます。このデータの集合 (x_i, y_i) ($i = 1 \sim n$)
の中から適当にm個をサンプリングして、最小二乗法により3次近似式 $y = a x^3 + b$
 $x^2 + c x + d$ の係数を求めます。

ここで探索したいデータ y_0 を3次近似式に代入して根を求めます。この根が求めるデ
ータの順番になりデータが発見されたこととなります。

ただし近似式なので、根も近似解となり根の周辺をさらに探索して真の解を求めます。
データの分布状況により近似解の精度が変わってきます。

おおむねランダムや正規分布のデータでは精度は比較的良くなります。また、サンプリ
ングの個数を増やせば精度は良くなります。

2. 最小二乗法により3次近似式を求める

データの集合 (x_i, y_i) ($i = 1 \sim n$) が与えられたとき、この各点の近くを通る3
次近似式 $f(x) = a x^3 + b x^2 + c x + d$ の係数を求めます。

データと3次近似式との残差の2乗和、

$$J = \sum_{i=1}^n \{ y_i - f(x_i) \}^2$$

においてJが最小になるように係数を定めます。

$$\frac{\partial J}{\partial a_j} = 0 \quad (j = 1 \sim 4)$$

これより未知数 a、b、c、d に関する4元連立1次方程式 (正規方程
式) を解いて未知数を求めます。

$$\begin{bmatrix} & x_i & x_i^2 & x_i^3 \\ x_i & x_i^2 & x_i^3 & x_i^4 \\ x_i^2 & x_i^3 & x_i^4 & x_i^5 \\ x_i^3 & x_i^4 & x_i^5 & x_i^6 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} f(x_i) \\ f(x_i) x_i \\ f(x_i) x_i^2 \\ f(x_i) x_i^3 \end{bmatrix}$$

3. 3次方程式の根を求める

3次方程式 $ax^3 + bx^2 + cx + d = 0$ において、

$$x = t - \frac{b}{3a} \quad \text{とおくと 式は}$$

$$t^3 + pt + q = 0$$

$$p = \frac{1}{a} \left(c - \frac{b^2}{3a} \right) \quad q = \frac{1}{a} \left(\frac{2b^3}{27a^2} - \frac{bc}{3a} + d \right)$$

となります。

式を解いて(カルダノの解法) 3つの根を求めます。

$$t = \sqrt[3]{r} - \sqrt[3]{s}$$

$$t = -\sqrt[3]{r} - \sqrt[3]{s}$$

$$t = \sqrt[3]{r} - \sqrt[3]{s}$$

となる。ただし、

$$r = \frac{q + \sqrt{q^2 + (4/27)p^3}}{2}$$

$$s = \frac{q - \sqrt{q^2 + (4/27)p^3}}{2}$$

$$= \frac{-1 + i\sqrt{3}}{2}$$

です。

ここで判別式 $D = -(4p^3 + 27q^3)$ において、

$D < 0$ のとき 1つの実根と2つの虚根となります。

$$\text{実根は } t = \sqrt[3]{r} - \sqrt[3]{s}$$

$D > 0$ のとき 3つの実根となります。

$$t = \sqrt[3]{2} \sqrt[3]{a^2 + b^2} \cos \frac{\theta}{3}$$

$$t = -\sqrt[3]{2} \sqrt[3]{a^2 + b^2} \cos \left(\frac{\theta}{3} - \frac{2\pi}{3} \right)$$

$$t = -\sqrt[3]{2} \sqrt[3]{a^2 + b^2} \cos \left(\frac{\theta}{3} + \frac{2\pi}{3} \right)$$

$$= \tan^{-1} \left(\frac{b}{a} \right)$$

$D = 0$ のとき 1つの3重根となります。

$$t = -\sqrt[3]{\frac{q}{2}}$$

です。

4. 探索のアルゴリズム

与えられたデータ (i, y_i) ($i = 1 \sim n$) に最小二乗法を適用することによって3次近似式が求められました。

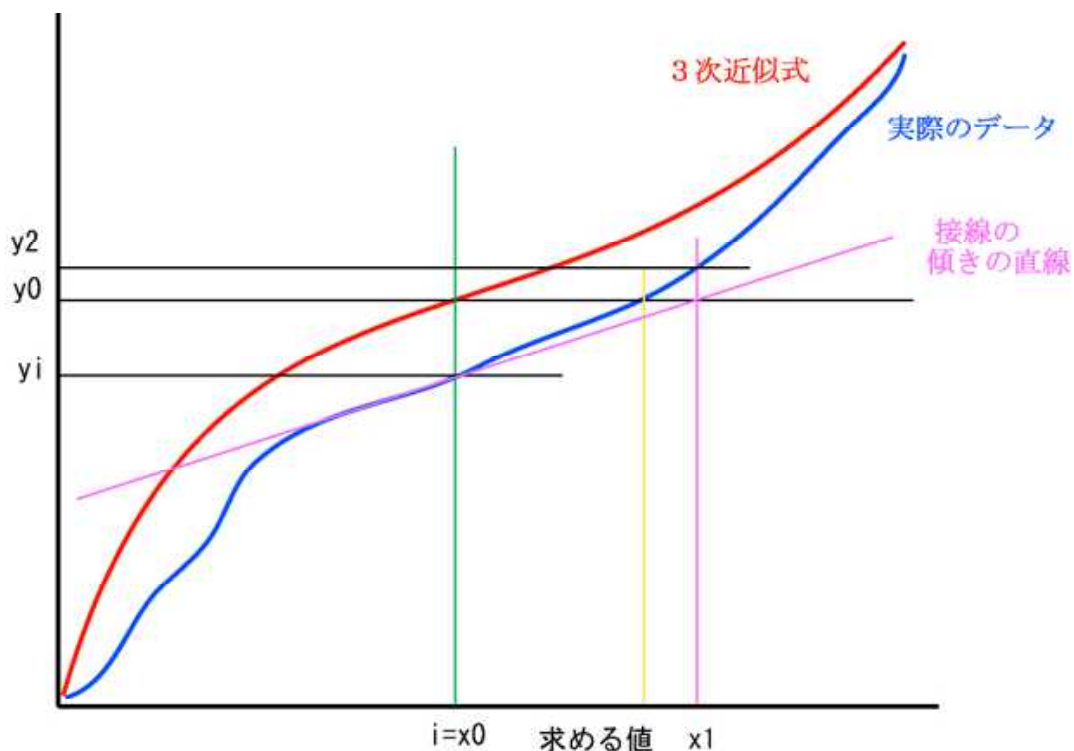
アルゴリズムAは、判別式 $D \leq 0$ の場合 (1実根) に用います。

(1) 調べたいデータ y_0 を3次近似式に代入してその根を求めます。この根の値 x_0 が y_0 の番号 i になります。

$y_i = y_0$ ならば、データを発見。

(2) それ以外ならば、 x_0 において3次近似式の接線の傾き k_1 を求めて、点 (x_0, y_i) を通る傾き k_1 の直線 $y = k_1 x + m_1$ を求めます。この直線の式に y_0 を代入して x_1 を求めます。

$i = x_1$ のときの y_i を求めて y_2 とします。



(3) $y_0 < y_2$ ならば、

$$x_0 = \frac{x_0 + x_1}{2} \quad \text{として(1)へ行きます。}$$

$y_0 > y_2$ ならば、 x_1 において3次近似式の接線の傾き k_2 を求めて、点 (x_1, y_2) を通る傾き k_2 の直線 $y = k_2 x + m_2$ を求めます。この直線の式に y_0 を代入して x_2 を求めます。

$i = x_2$ のときの y_i を求めて y_2 として、(3)へ戻ります。

$y_0 = y_2$ ならば、データを発見。

以上がアルゴリズムであるが、判別式 $D > 0$ になった場合には3つの実根が存在するので、アルゴリズムBを用います。

(4) 3次近似式の2つの極値を求めます。

3次近似式を微分した $y' = 3ax^2 + 2bx + c = 0$ の2根 x_1, x_2 ($x_1 < x_2$) が極値となります。

$$x = \frac{-b \pm \sqrt{b^2 - 3ac}}{3a}$$

$i = x_1$ のときの y_i を求めて y_1 とする。 $i = x_2$ のときの y_i を求めて y_2 とします。

(5) $y_0 > y_2$ または $y_0 < y_1$ ならば、アルゴリズムAを用います。

(6) $y_1 < y_0 < y_2$ ならば、区間 $x_1 \sim x_2$ を直線の式で近似します。

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1} (x - x_1) \quad -$$

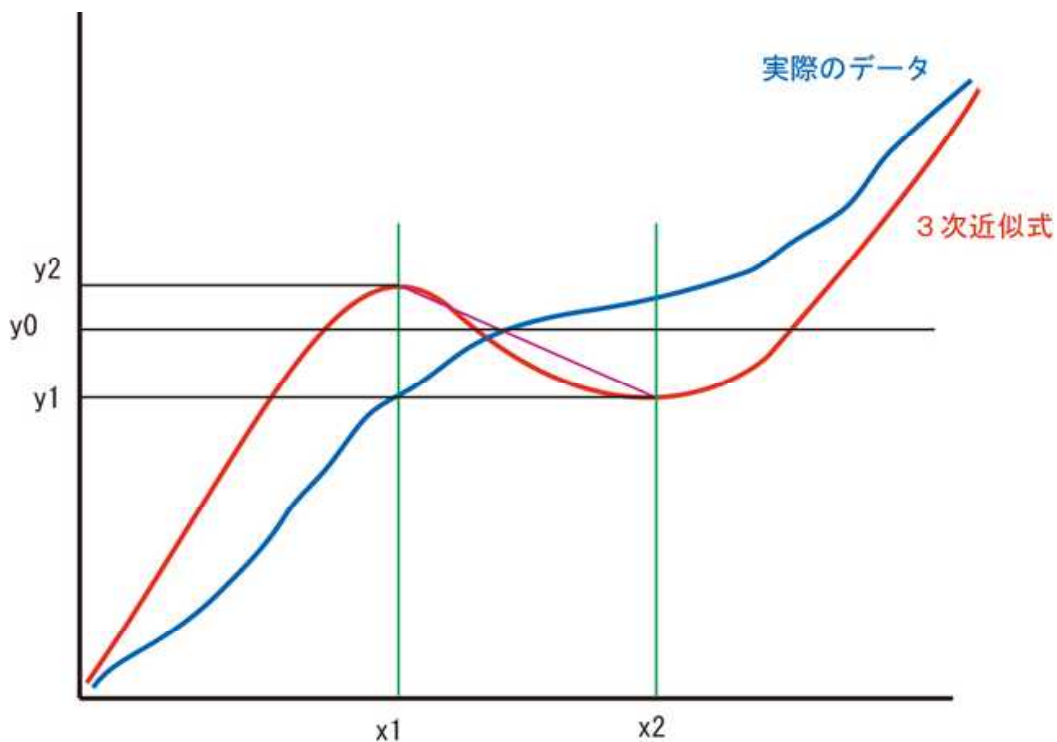
(7)

$x_0 = \frac{x_1 + x_2}{2}$ として、式に代入して y_3 を求めます。

$y_0 = y_3$ ならば、データを発見。

$i = x_0$ のときの y_i を求めて y_4 とします。

(8) $y_0 > y_4$ ならば $x_1 = x_0$ として、 $y_0 < y_4$ ならば $x_2 = x_0$ として(7)に戻ります。



5 . その他

サンプリング探索法では、データの分布が偏ると近似が悪くなり探索により多くの回数がかかるようになります。

データの値が大きくなるとオーバーフローする可能性があり、最小二乗法を使うときに行列を用いるのでメモリアオーバーになるかもしれません。

また文字列の探索には、あまり向いていないように思われます。

近似式として3次式を用いたのは、必ず実根が1つはあり、ほぼ(3実根のときは一部で減少)単調増加であるという理由によります。