

Ridge 回帰の実装

正則化パラメータ λ の値によって推定される曲線がどのように変わるのか、Fortran90 でプログラムを組んで実験した。ここでは正則化項を $\lambda \|\mathbf{w}\|_2^2$ とした (Ridge 回帰)。

1 正則化最小二乗法 (Ridge 回帰)

トレーニングデータを $(x_i, y_i)_{i=1, \dots, m}$ と与える。これを n 次関数 $\hat{y} = w_0 + w_1x + w_2x^2 + \dots + w_nx^n = \sum_{k=0}^n w_k x^k$ で近似する。 x の次数順に並べたベクトルを $\mathbf{x} = [1, x, x^2, \dots, x^n]^\top$, 係数を並べたベクトルを $\mathbf{w} = [w_0, w_1, \dots, w_n]^\top$ とすれば (これらは $n+1$ 次の列ベクトル)、この関数は $\hat{y} = \mathbf{w}^\top \mathbf{x}$ となる。各トレーニングデータ (x_i, y_i) を使って単純にこの n 次多項式を計算した $\hat{y}_i = \sum_{k=0}^n w_k x_i^k$ を並べた m 次列ベクトルを $\hat{\mathbf{y}} = [\hat{y}_1, \dots, \hat{y}_m]^\top$ とする。これは行ベクトル $\mathbf{x}_i^\top = [1, x_i, x_i^2, \dots, x_i^n]$ を縦に並べた行列

$$X = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_m^\top \end{bmatrix} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & \dots & x_2^n \\ & & \vdots & & \\ 1 & x_m & x_m^2 & \dots & x_m^n \end{bmatrix}$$

を用いて $\hat{\mathbf{y}} = X\mathbf{w}$ とかける。ただしこの X は $m \times (n+1)$ の行列である。平均二乗誤差 $\text{MSE}(\mathbf{w})$ を行列で表現して展開すると

$$\begin{aligned} \text{MSE}(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 = \frac{1}{m} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2 = \frac{1}{m} \|X\mathbf{w} - \mathbf{y}\|_2^2 \\ &= \frac{1}{m} (X\mathbf{w} - \mathbf{y})^\top (X\mathbf{w} - \mathbf{y}) = \frac{1}{m} (\mathbf{w}^\top X^\top - \mathbf{y}^\top) (X\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{m} (\mathbf{w}^\top X^\top X\mathbf{w} - \mathbf{y}^\top X\mathbf{w} - \mathbf{w}^\top X^\top \mathbf{y} + \mathbf{y}^\top \mathbf{y}) \end{aligned}$$

となる。ただし $\mathbf{y} = [y_1, \dots, y_m]^\top$ とした。目下の目的は平均二乗誤差 $\text{MSE}(\mathbf{w})$ を最小化する \mathbf{w} を求めることである。行列の微分の公式

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^\top D \mathbf{w}) = (D + D^\top) \mathbf{w}, \quad \frac{\partial}{\partial \mathbf{w}} (\mathbf{a}^\top \mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^\top \mathbf{a}) = \mathbf{a}$$

を用いて平均二乗誤差 $\text{MSE}(\mathbf{w})$ を微分すると

$$\frac{\partial \text{MSE}(\mathbf{w})}{\partial \mathbf{w}} = \frac{2}{m} (X^\top X \mathbf{w} - X^\top \mathbf{y})$$

となるので、これが $\mathbf{0}$ となるのは $\mathbf{w} = (X^\top X)^{-1} X^\top \mathbf{y}$ のときである。

つぎに平均二乗誤差 $\text{MSE}(\mathbf{w})$ を最小化するのではなく、重み減衰 $\lambda \|\mathbf{w}\|_2^2 = \lambda \mathbf{w}^\top \mathbf{w}$ を加えたもの $J(\mathbf{w}) = \text{MSE}(\mathbf{w}) + \lambda \mathbf{w}^\top \mathbf{w}$ を最小化することを考える。 $\partial(\mathbf{w}^\top \mathbf{w}) / \partial \mathbf{w} = 2\mathbf{w}$ より

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \frac{\partial \text{MSE}(\mathbf{w})}{\partial \mathbf{w}} + \frac{\partial (\lambda \mathbf{w}^\top \mathbf{w})}{\partial \mathbf{w}} = \frac{2}{m} (X^\top X \mathbf{w} - X^\top \mathbf{y}) + 2\lambda \mathbf{w},$$

したがってこれが $\mathbf{0}$ となるのは $\mathbf{w} = (X^\top X + m\lambda I)^{-1} X^\top \mathbf{y}$ のときである。

プログラムを書くにあたって記号をいろいろ定義する。 $A := X^\top X$, $B := A + m\lambda I$, $\mathbf{v} := X^\top \mathbf{y}$ とおく。すると目的は $(n+1)$ 本の 1 次連立方程式 $B\mathbf{w} = \mathbf{v}$ を \mathbf{w} について解くことに帰着される。この目的は、さらに拡大係数行列 $C := \begin{bmatrix} B & \mathbf{v} \end{bmatrix}$ にガウスの掃き出し法を行うことに帰着される。ただし

$$X^\top = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_m \\ x_1^2 & x_2^2 & \cdots & x_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & \cdots & x_m^n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^n \\ 1 & x_2 & x_2^2 & \cdots & x_2^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \cdots & x_m^n \end{bmatrix}$$

より A の成分を具体的に書くと

$$A = X^\top X = \begin{bmatrix} s_0 & s_1 & s_2 & \cdots & s_m \\ s_1 & s_2 & s_3 & \cdots & s_{n+1} \\ s_2 & s_3 & \cdots & \cdots & s_{n+2} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ s_n & s_{n+1} & s_{n+2} & \cdots & s_{n+n} \end{bmatrix}, \quad s_k = \sum_{i=1}^m x_i^k, \quad k = 0, \dots, 2n$$

となる。つまり A は対角線に垂直な方向の成分がすべて同じであるような、 $(n+1)$ 次の対角行列になっている。その成分は $a_{ij} = s_{i+j-2}$, $i, j = 1, \dots, n+1$ と定式化できる。 B はその対角成分に $m\lambda$ を加えた行列である。ただし λ は任意の変数であるから、 $m\lambda$ ではなく単に λ を加えることにする。

$$b_{ij} = \begin{cases} a_{ij} + \lambda, & \text{if } i = j, \\ a_{ij}, & \text{if } i \neq j. \end{cases}$$

さらに \mathbf{v} を具体的に示すと

$$\mathbf{v} = X^\top \mathbf{y} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_m \\ x_1^2 & x_2^2 & \cdots & x_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^n & x_2^n & \cdots & x_m^n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} v_0 \\ v_1 \\ \vdots \\ v_n \end{bmatrix}, \quad v_l = \sum_{i=1}^m x_i^l y_i, \quad l = 0, \dots, n$$

である。拡大係数行列 C は

$$C = \begin{bmatrix} B & \mathbf{v} \end{bmatrix} = \begin{bmatrix} b_{11} & b_{21} & \cdots & b_{n+1,1} & v_0 \\ b_{12} & \ddots & & \vdots & v_1 \\ \vdots & & \ddots & \vdots & \vdots \\ b_{1,n+1} & \cdots & \cdots & b_{n+1,n+1} & v_n \end{bmatrix} \quad (1)$$

よりその成分 c_{ij} , $i = 1, \dots, n+1$, $j = 1, \dots, n+2$ は

$$c_{ij} = \begin{cases} b_{ij}, & \text{for } i, j = 1, \dots, n+1, \\ v_{i-1}, & \text{for } i = 1, \dots, n+1, j = n+2 \end{cases}$$

とかける。これは $(n+1) \times (n+2)$ の行列になる。

2 実験方法とソースコード

まずデータを読み込んで、推定したい多項式の係数 $\boldsymbol{w} = [w_0, \dots, w_n]^T$ と、推定曲線の (x, y) の値を出力するプログラムを作成した。それをソースコード 1 に示す。 $\boldsymbol{w} = [w_0, \dots, w_n]^T$ を求めるのには、式 (1) の拡大係数行列 C にガウスの消去法を適用した。データの数 m 、多項式の次数 n 、正則化パラメータ λ 、プロットの分割幅はプログラム作成者が入力できるようにした。 m 個の訓練データは、区間 $[0, 2)$ の一様乱数を関数 $2 \cos(2x) + x$ に足すことで生成した。それをソースコード 2 に示す。

正則化パラメータを $\lambda = 10^4, 10^2, 10^0, 10^{-2}, 10^{-4}, 0$ としたときに、 $m = 30$ 個のデータを $n = 10, 20, 30$ 次の多項式で推定した。

ソースコード 1 リッジ回帰のソースコード

```
1 program ridge_regression
2 implicit none
3
4 real(8), allocatable :: x(:), y(:), data(:), s(:), v(:)
5 real(8), allocatable :: a(:, :), b(:, :), c(:, :), w(:)
6 real(8) :: lambda, p, q, width, delta, xmin, y_hat, x1
7 integer :: i, j, n, m, k, l, i1, i2
8
9 ! ***入力*****
10 ! データの数
11 m = 30
12
13 ! 近似したい曲線の次数
14 n = 10
15
16 ! 正則化パラメータ
17 lambda = 1.0_8
18
19 !***データの読み込み*****
20 allocate(x(m), y(m), data(2 * m))
21 ! mx2 のデータ (「data.txt」など) をサイズ 2m の配列 data(2*m) として読み込む。
22 ! (「./a.out < data.txt」)
23 ! data(2*m) は x, y, x, y, ... のような並びになっている。
24
25 read(*,*) data
26
27 ! data(2*m) の奇数番目を x(j), 偶数番目を y(j) とする。
28 do i = 1, 2*m
29   if ( mod(i, 2) == 0 ) then
30     j = i / 2
31     y(j) = data(i)
32   else if ( mod(i, 2) == 1 ) then
33     j = (i + 1) / 2
```

```

34      x(j) = data(i)
35      end if
36  end do
37
38  ! 確認済み
39  ! do j = 1, m
40  ! write(*,*) j, x(j), y(j)
41  ! end do
42
43  !***拡大係数行列の準備*****
44  ! 解きたい連立方程式 $Bw=v$  の B と v の成分を計算しておく。
45  ! 拡大係数行列 $C=[B \ v]$ を構築する。
46  allocate(s(0:2 * n), a(n + 1, n + 1), b(n + 1, n + 1), v(0 : n), c(n + 1, n + 2))
47
48  !  $s_k = \sum_{i=1}^m \{x_i\}^k$ ,  $k=0 \sim 2n$  の計算
49  do k = 0, 2 * n
50      s(k) = 0.0_8
51      do i = 1, m
52          s(k) = s(k) + (x(i)) ** (real(k))
53      end do
54  ! write(*,*) k, s(k)
55  end do
56
57  ! 行列A の計算。その成分は  $a_{ij}=s_{i+j-2}$ ,  $i,j=1 \sim n+1$ .
58  do j = 1, n + 1
59      do i = 1, n + 1
60          a(i, j) = s(i + j - 2)
61  ! write(*,*) i, j, i + j - 2, s(i + j - 2), a(i, j) ! 確認済み
62      end do
63  end do
64
65  ! 行列B の計算。行列 A の対角成分に lambda を足すだけ。
66  do j = 1, n + 1
67      do i = 1, n + 1
68          if ( i == j ) then
69              b(i, j) = a(i, j) + lambda
70          else if ( i /= j ) then
71              b(i, j) = a(i, j)
72          end if
73  ! write(*,*) i, j, b(i, j) - a(i, j) ! 確認済み。
74      end do
75  end do
76
77  !  $v_l = \sum_{i=1}^m y_i \{x_i\}^l$ ,  $l=0 \sim n$  の計算
78  do l = 0, n
79      v(l) = 0.0_8

```

```

80      do i = 1, m
81          v(1) = v(1) + y(i) * (x(i)) ** real(1)
82      end do
83  ! write(*,*) 1, v(1) ! 確認済み
84      end do
85
86  ! 拡大係数行列C=[B v]を作る
87  do i = 0, n
88      do j = 1, n + 2
89          if (j == n + 2) then
90              c(i + 1, j) = v(i)
91          else if (j /= n + 2) then
92              c(i + 1, j) = b(i + 1, j)
93          end if
94      end do
95  end do
96  ! 確認済み
97  ! do i = 0, n
98  ! do j = 1, n + 2
99  ! write(*,*) i + 1, j, v(i), c(i + 1, j)
100 ! end do
101 ! end do
102
103 ! ***ガウスの掃き出し法*****
104 ! 上で作った (n+1)x(n+2)の拡大係数行列c にガウスの掃き出し法を適用する。
105 ! あるk 行目の対角成分 c(k,k)をp とおき、その行すべてを p で割る。
106 ! そうすればその対角成分c(k,k)は 1になる。
107 ! k 行目以外の行について (これを i 行目とする)、c(k,k)と同じ列にあるもの (c(i,k))をq とおく。
108 ! そのi 行目の k 列目以降の成分 c(i,j), j=k~n+2 からq*c(k,j)を引く。
109 ! これでc(k,k)の上下はゼロになる。
110 ! このことをすべての行 k=1~n+1 について行う。
111
112 do k = 1, n + 1
113     p = c(k, k)
114     do j = k, n + 2
115         c(k, j) = c(k, j) / p
116     end do
117     do i = 1, n + 1
118         if (i /= k) then
119             q = c(i, k)
120             do j = k, n + 2
121                 c(i, j) = c(i, j) - q * c(k, j)
122             end do
123         end if
124     end do
125 end do

```

```

126
127 ! 確認済み
128 ! do i = 1, n + 1
129 ! do j = 1, n + 2
130 ! write(*,*) i, j, c(i, j)
131 ! end do
132 ! end do
133
134 ! こうして得られた行列の最後の列c(i,n+2), i=1~n+1がwの解である。
135 allocate(w(n + 1))
136 do i = 1, n + 1
137     w(i) = c(i, n + 2)
138     write(*,*) i - 1, w(i) ! 確認済み
139 end do
140
141 ! ***曲線のプロット*****
142 ! 分割数
143 i1 = 200
144
145 ! xの幅 width と xの最小値 xmin を計算。
146 width = maxval(x) - minval(x)
147 xmin = minval(x)
148 ! write(*,*) width, xmin, maxval(x) ! 確認済み
149
150 ! 曲線のプロットの刻み幅
151 delta = width / real( i1 )
152
153 open (18, file='line.txt', status='replace')
154 do i2 = 0, i1
155     y_hat = 0.0_8
156     x1 = xmin + real(i2) * delta
157     do i = 1, n + 1
158         y_hat = y_hat + w(i) * x1 ** (real(i - 1))
159     end do
160     write(18, *) x1, y_hat ! ここでxに対する近似曲線  $y_{\text{hat}} = \sum_{k=0}^n w_k x^k$  の値を出力
161 end do
162 close(18)
163
164 deallocate(x, y, data, s, v, a, b, c, w)
165
166 stop
167 end program

```

ソースコード 2 データ生成のためのプログラム

```

1 program data_generator
2 implicit none

```

```

3
4  integer :: m, imax, i
5  real(8) :: xmin, xmax, width, delta, x, y
6  real(8), allocatable :: random(:)
7
8
9  ! データの個数
10 m = 30
11
12  allocate(random(m + 1))
13  call random_number(random)
14  ! 確認
15 ! do i = 0, m
16 ! write(*,*) i, 2.0_8 * random(i + 1)
17 ! end do
18
19 ! データの最小値・最大値の設定、幅、分割幅
20 xmin = 0.0_8
21 xmax = 6.0_8
22 width = xmax - xmin
23 delta = width / real(m)
24 ! write(*,*) width
25
26  open(18, file='generated_data.txt', status='replace')
27  do i = 0, m - 1
28      x = xmin + delta * real(i)
29      y = 2.0_8 * cos(2.0_8 * x) + x + 2.0_8 * random(i + 1)
30      write(18,*) x, y
31  end do
32  close(18)
33
34  deallocate(random)
35
36  stop
37  end program

```

3 出力結果

$n = 10, 20, 30$ に対して正則化パラメータを $\lambda = 10^4, 10^2, 10^0, 10^{-2}, 10^{-4}, 0$ としたときの推定結果を図1から図6に示す。与えた30個のデータは点で表し、推定した曲線を実線で表している。ただし推定した曲線は200個の点を線で繋いだものなので、鋭いピークがあるとそれが削れて表示されている可能性がある。

$n = 10$ など n が比較的小さいとき、 λ の値が大きいきちんとデータ点を推定しておらず、過少適合をおこなっているのが観察できる。逆に $n = 30$ など n が大きく λ が小さいときには過剰適合を起こしている。

$n = 30$, $\lambda = 10^2$ のときの推定曲線が過剰に振動している。 λ の値を少し変えて $\lambda = 99, 101$ として比較したが、そのようなことは起こっていなかった (図 7)。なぜそうなるのかはよくわからない。

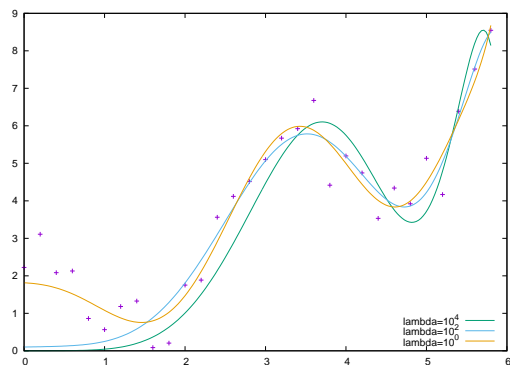


図 1 $n = 10$, $\lambda = 10^4, 10^2, 10^0$

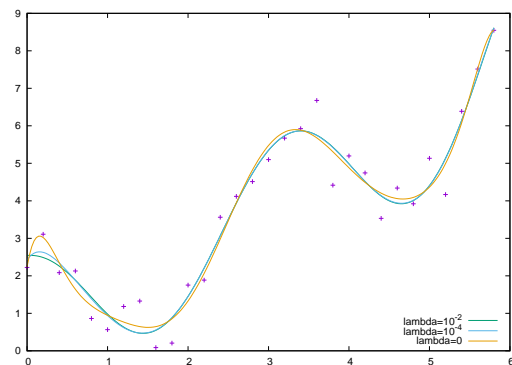


図 2 $n = 10$, $\lambda = 10^{-2}, 10^{-4}, 0$

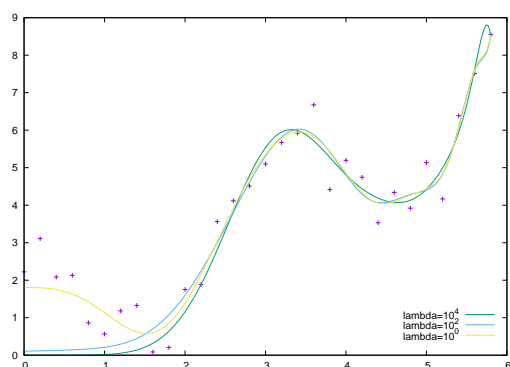


図 3 $n = 20$, $\lambda = 10^4, 10^2, 10^0$

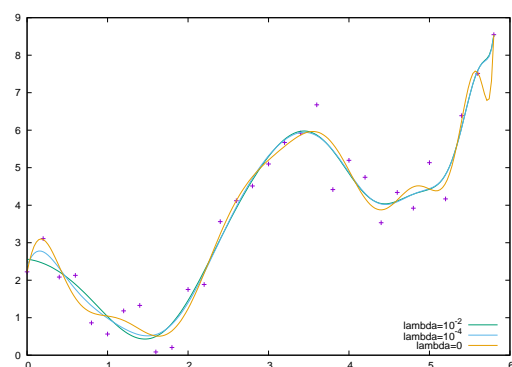


図 4 $n = 20$, $\lambda = 10^{-2}, 10^{-4}, 0$

4 今後の課題

なぜ $n = 30$, $\lambda = 10^2$ のときに異常な出力を得たのかを解決したい。

また適切な n と λ をどう決めるのかわからなかった。そこで n と λ を少しずつ変えて出力させ、テストデータを与えて汎化誤差が一番小さい n と λ を適切なものとする、というようなプログラムも書けるだろう。

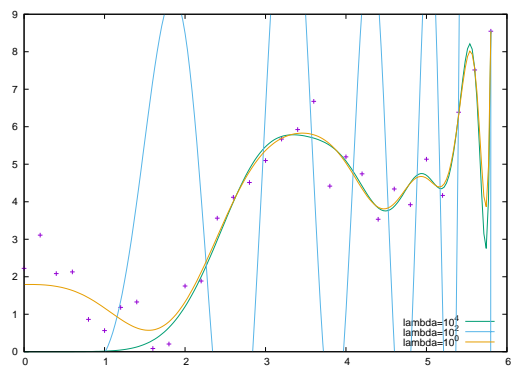


図5 $n = 30, \lambda = 10^4, 10^2, 10^0$

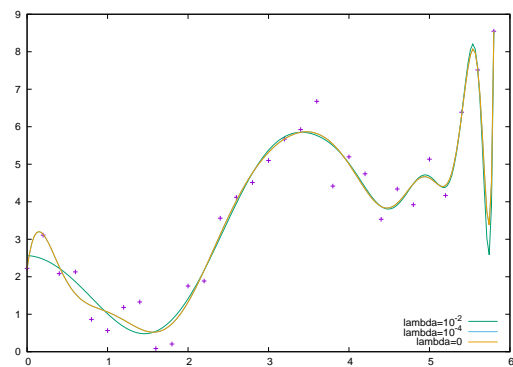


図6 $n = 30, \lambda = 10^{-2}, 10^{-4}, 0$

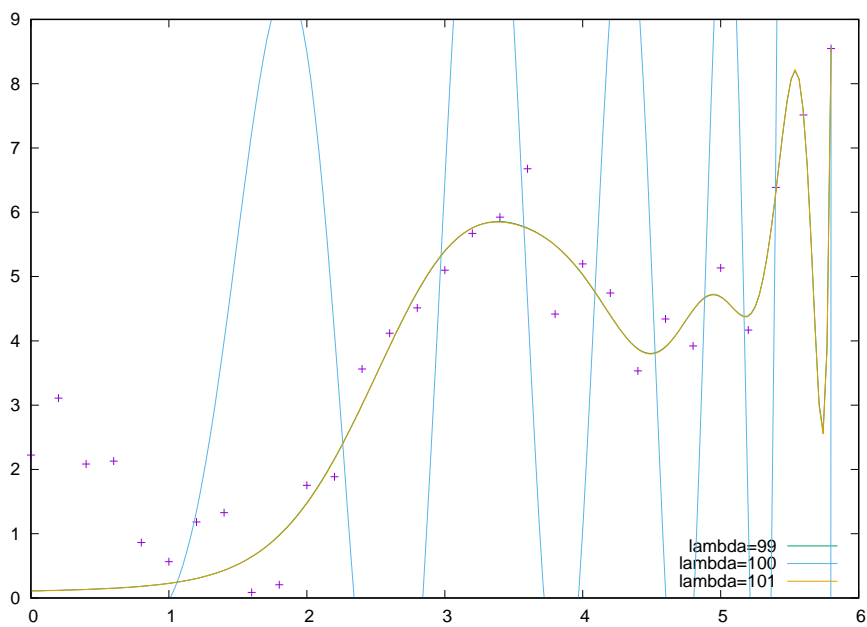


図7 $n = 30$ のときに $\lambda = 99, 100, 101$ を比較したもの