

外れ値を検出する方法とその特性に関する研究

2017SS004 青木勇太

指導教員：松田真一

1 はじめに

本研究では、外れ値を検出する手法に対してシミュレーションを行い、その外れ値検出手法の特性を理解し、性能を比較することを目的としている。

2 外れ値検出手法

本研究では R で利用可能な以下の手法を用いた。1次元では、ホテリング理論、箱ひげ図、スミルノフ・グラブス検定。多次元では、ホテリング T^2 理論、One Class Support Vector Machine (SVM)、Local Outlier Factor (LOF) である。(Web[2], 高島 [5] 参照)

なお、本研究では各手法の実現に R のパッケージを利用したが、ホテリング理論、多変量のホテリング T^2 理論は井出 [3] を、スミルノフ・グラブス検定は青木 [1] を参考に作成した。

3 シミュレーションについて

本研究では、以下の3つのシミュレーションを行う。

- 1次元の手法の特性を見るシミュレーション
- 1次元の手法の性能比較のシミュレーション
- 多次元の手法の特性比較のシミュレーション

以降はこれらをシミュレーション1等と呼ぶ。

4 シミュレーション1

1次元の各手法の特性を調べるためにシミュレーションを行った。外れ値を含みやすいという理由から自由度3の t 分布をデータに用いて、各手法がどのように外れ値を検出するかを確かめた。サンプルサイズ100個、試行回数1000回を行った。

それぞれの手法で検出された外れ値の総数は、ホテリング理論が2610個、箱ひげ図が5617個、スミルノフ・グラブス検定が2304個であった。

箱ひげ図が検出した外れ値の総数がほかの2倍以上あることから箱ひげ図は、計算段階で改善する必要があると考えられる。

4.1 箱ひげ図の改善について

まず、検出された総数の5617個という値が理論値に従ったものなのかを確かめる。自由度3の t 分布における外れる確率は5.6%であった。本シミュレーションは値を100000個発生させているから、その中の5.6%は5600個であり5617個は理論値に従っていることがわかった。

よって、他の手法の総数が約2500個だとすると、信頼区間は97.5%にすればよい。その場合に四分位範囲にかける値は2.2であった。そうしてもう一度シミュレーショ

ンした結果、外れ値の総数は2640個となった。よって改善できた。

5 シミュレーション2

本シミュレーションは、用いた1次元の手法の性能比較のために行った。外れる個数を二項分布によって決定して、二項分布の外れる確率を指定することにより行うシミュレーションである。外れ値を含まない方の山を平均0分散1の正規分布に従うように、外れ値を含む方の山を平均5分散0.3の正規分布に従うように乱数で発生させた。サンプルサイズ100個、試行回数1000回を行った。

なお、箱ひげ図は4.1節で改善したものをを用いた。

以下表1が結果であり、外れ値を含む確率それぞれに対して外れ値を正確に検出した回数を表している。

表1 各手法が外れ値の数を正確に検出した回数

| | ホテリング | 箱ひげ図 | スミルノフ |
|-----|-------|------|-------|
| 1% | 625 | 941 | 954 |
| 3% | 854 | 913 | 740 |
| 5% | 823 | 857 | 389 |
| 10% | 226 | 598 | 24 |

表1から、スミルノフ・グラブス検定とホテリング理論は、悪いときの値がかなり悪いが、箱ひげ図は、全体的にいい結果だとわかる。スミルノフ・グラブス検定は確率が小さいときに他より性能がいいと言える。ただし、それは確率が低いときのみに限ったことであり、1%から10%に向けて箱ひげ図の落ち幅が他より少なく、全体を通して言えば箱ひげ図のほうがいいと言える。

また、ホテリング理論だけが外れ値を含む確率が1%の時に一番良い結果にならないことについて、外れ値を含まない山の分布が原因なのではないかと考え、外れ値を含まない平均0分散1の正規分布のシミュレーションを行った。その結果が図1であり、ホテリング理論は外れ値を含まない場合に他より外れ値を検出しやすいことがわかった。したがって、ホテリング理論が外れ値を含む確率が1%のときに一番良い結果にならない原因は、理論ベースである正規分布に完全に従う場合に限界値の設定から1%の外れ値を検出しようとしてしまうことにある。すなわち外れ値ではないのにサンプルサイズが100なので図1で1個見つける場合が一番多くなっている。

6 シミュレーション3

多次元の各手法の性能を比較するために行うシミュレーションである。データは、外れ値を検出するのが難しいデータとして Peña and Prieto[4], 和田 [6] が使用した

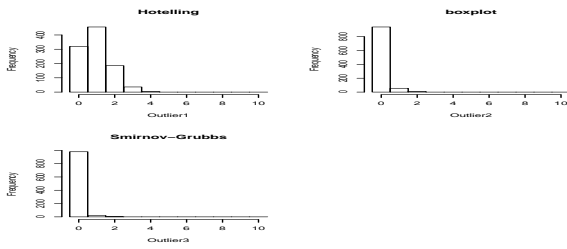


図 1 外れ値のない正規分布

多次元混合正規分布

$$(1 - \alpha)N_M(\mathbf{0}, I) + \alpha N_M(\delta \mathbf{e}_1, \lambda I) \quad (1)$$

を用いる。 α が外れ値の割合、 M が次元数、 δ が原点からの距離、 $\mathbf{e}_1 = \{1, 0, \dots, 0\}$ 、 λ が外れ値の分散を表している。パラメータは和田 [6] を参考に $\delta = 10$ 、 $\lambda = 0.01$ に定めた。

シミュレーションは、2次元データで外れ値の個数を5個から1個まで変化させシミュレーションを行い、その後次元数を5次元まで変化させて行う。紙面上の都合上、2次元で外れ値5個、2次元で外れ値1個、4次元で外れ値5個の結果のみを図2, 3, 4で示す。

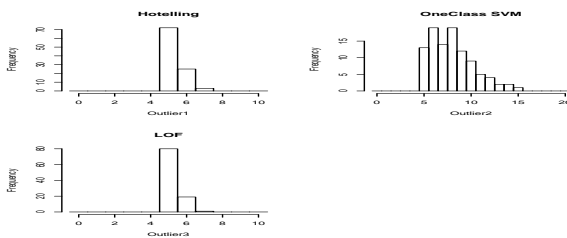


図 2 2次元データで外れ値5個

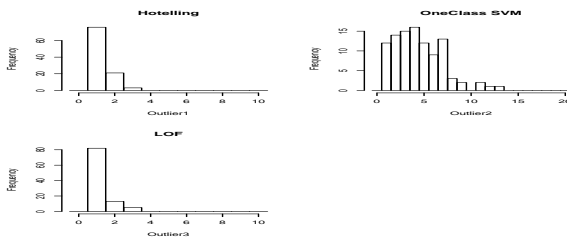


図 3 2次元データで外れ値1個

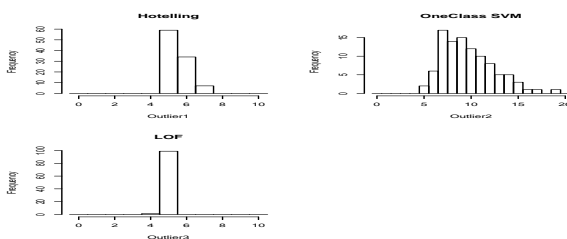


図 4 4次元データで外れ値5個

One Class SVM は、与えた混合分布自身で検定したら

うまくいかなかったので、外れ値を含まないデータを学習させて外れ値を含むデータの検定を行ったが、ほかの手法のような結果にはならなかった。

2次元のシミュレーションでは、LOFが一番精度がよく、その次にホテリング理論というように見える。次に次元数を大きくしたときは、ホテリング理論は横に広がる結果となり、反対にLOFは、精度がよくなっていることがわかる。この結果からも次元が大きいとLOFの精度がよくなっていて、2次元の時点でも精度としてが高いことから、LOFが一番いいと言える。

7 まとめ

1次元データに対する手法に関しては、外れ値が含まれる確率が低いときのスミルノフ・グラブス検定の精度の高さが目立ったが、実際に外れ値を検出する場面や、外れ値が含まれる確率が高いときのことを考えると改善後の箱ひげ図が一番性能がいいとわかった。多次元データに対する手法に関しては、次元数が大きくなって結果がよくなったLOFが一番いいとわかった。しかし、5次元データでLOFが稀に全く外れ値を検出できていないことがあった。この理由は、データの間隔が広がりすぎてうまくいかないせいだと考えられる。サンプルサイズがもっと大きければ改善する可能性がある。

8 おわりに

本研究を通して、用いた外れ値検出手法の使用方法や、各手法の特性や性能について理解することができた。外れ値を検出した時、外れ値を除くことが必要なのかわかりと判断したしながら、本研究で学んだことを活かしたい。

参考文献

- [1] 青木繁伸：「スミルノフ・グラブス検定」。
<http://aoki2.si.gunma-u.ac.jp/lecture/Grubbs/Grubbs.html/>, 2015. (2020/6 閲覧)
- [2] BellCurve：「4-3. 外れ値検出のある箱ひげ図」,
<https://bellcurve.jp/statistics/course/5222.html/>, (2020/6 閲覧).
- [3] 井出剛：「入門 機械学習による異常検知—Rによる実践ガイド」, コロナ社, 2015.
- [4] Peña,D. and F.J.Prieto：「Multivariate Outlier Detection and Robust Covariance Matrix Estimation」, *Technometrics*, Vol.43, pp.286-300, 2001..
- [5] 高島泰斗：「密度推定法に基づくカーネル判別機械」, 筑波大学大学院博士過程システム情報工学研究科修士論文, <https://commons.sk.tsukuba.ac.jp/wp-content/uploads/sites/13/2016/08/200520847.pdf/>, 2007. (2020/11 閲覧)
- [6] 和田かず美：「多変量外れ値の検出—MSD法とその改良手法について」, 統計研究彙報 第67号, pp.89-157, 2010.