



前回、第一回では統計的な異常検出の特徴と応用例について述べた。統計的な異常検出とは数学的・統計的な手法を用いた自動的な異常検出のことであり、クレジットカードの不正利用から機械の故障、ネットワークの異常、不正行為など幅広い応用が考えられる。今回は、その統計的な異常検出を実現する為の手法とはどのようなものかについて解説する。

### 3. 統計的異常検出手法

#### 【3.1 統計的な異常検出手法と機械学習】

統計的に異常を検出する為に、機械学習の技術がよく用いられる。従って、統計的異常検出手法は「機械学習を用いた異常検出手法」とも言い換えられる。機械学習とは「人間の持つ学習能力を機械(コンピューター)に持たせる」ことを目的とした研究分野である。簡単に言うと「人間ができることをコンピューターで自動的にやらせる方法」を考えるとということである。歴史的な経緯を抜きにすれば、機械学習という用語は統計・パターン認識・データマイニング等のアナリティクス技術・数理的データ分析技術を表す用語とほぼ同じである。

自動的な異常検出を実行するには、何らかのデータを用いてコンピューターに「学習」をさせる必要がある。異常検出を扱う場合、コンピューターに学習させるのは入力データと出力との関係性である。入力とは異常検出に用いる観測データであり、出力は異常かどうかの判定結果もしくは異常度を表すスコアのような異常かどうかの判断材料である。言い換えると、「学習」とは入力から出力を導く為の数理的な方法もしくは数理的なモデルを構築することである。

前回に述べたことの繰り返しになるが、統計的な異常検出の利点はビッグデータに対する自動的な分析が可能となる点にある。言い換えれば、機械学習を用いることの利点は、複雑且つ大量なデータを入力として自動的な分析を実現できることにある。

#### 【3.2 教師あり学習と教師なし学習】

機械学習の手法(もしくは統計的異常検出手法の構築方法)は、学習の方法の視点から教師なし学習と教師あり学習の二つに分類できる。以下では、これら二つの概要と異常検出への応用について説明する。

##### 【3.2.1 教師あり学習】

教師あり学習では、出力の正解例と入力との関係を学習・モデル化する。入力に対するあるべき出力(正解例)を学習させるという意味で「教師あり」学習と呼ばれる。例えば、信用格付けの推定問題(ローンの貸し出しにおける貸し倒れの発生を検知する異常検出問題)を考える。個人のプロフィール、所得、ローン残高等を入力として、貸し倒れの発生有無という出力の正解例(過去の事例)を予測するモデルを作るのは教師あり学習に含まれる。

教師あり学習の手法を用いるには、異常か否か(もしくは正常か異常かと、異常である場合には異常の種類)のラベルを含む過去の事例という、正解例付きのデータを用意しておく必要がある。例えば、信用格付けならば格付けの為の情報と貸し倒れの有無のフラグ、機械の故障検知ならばセンサーデータと故障の有無のフラグ、というデータが必要となる。

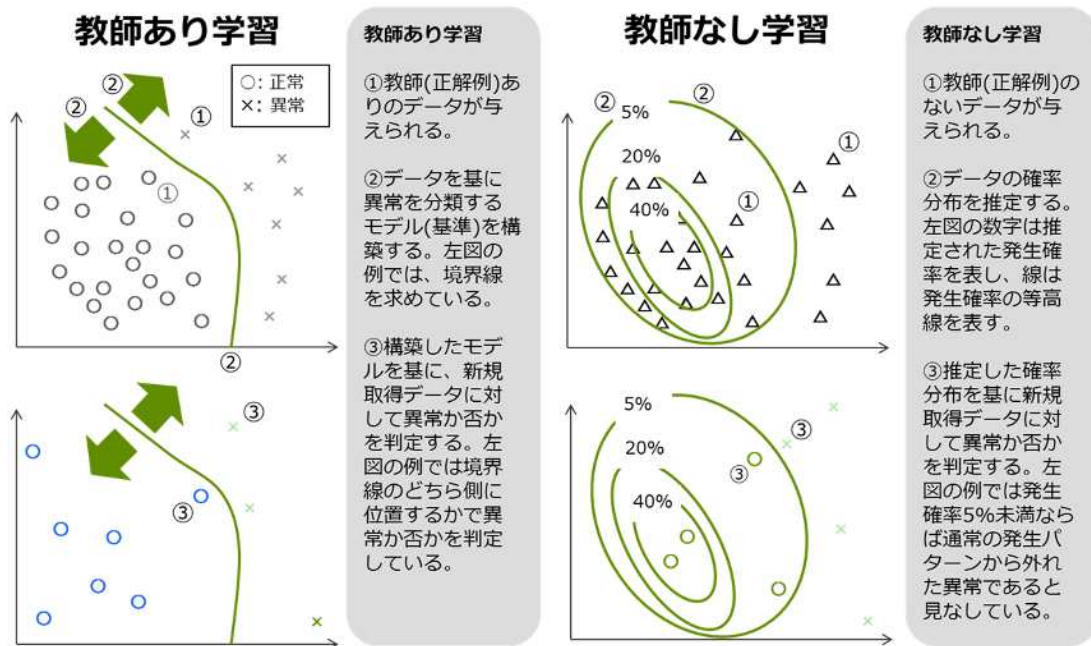
教師あり学習の結果として得られたモデルは、入力が与えられたときの事象(異常)の発生確率を出力する。例えば、上記の信用格付けの問題では、個人のプロフィールやローン残高等を入力としてモデルに与えると対象者が貸し倒れを発生させる確率が出力される。

教師あり学習に分類される上記のようなモデルは、入力から正解例を予測するという意味で「予測モデル」とも呼ばれる。予測モデルには様々な種類がある。例えば、決定木、回帰、フィッシャー判別分析、ニューラルネットワーク、SVM(サポートベクターマシン)、ランダムフォレスト等の手法は予測モデルに含まれる。

教師あり学習の利点は、過去の事例に類似するものを精度よく検出できることにある。異常検出の場合で言えば、過去に発生した異常と類似する異常を精度良く検出できる。これは、過去の事例を基にモデルを構築(関係性を学習)することによる。

一方で、教師あり学習の欠点として新規の事例(未知の異常)を検出できないことが挙げられる。過去の事例を基にモデルを構築するので、その事例に含まれない異常は検出できないかもしくは検出できたとしても偶然としか言えない。

図3. 機械学習：教師あり学習と教師なし学習



### 【3.2.2 教師なし学習】

教師なし学習では、入力に含まれるパターンや特徴を学習・モデル化する。教師あり学習とは異なり入力に対するあるべき出力(正解例)がない(学習させない)という意味で「教師なし」学習と呼ばれる。例えば、機械の故障検知を考える。センサーデータを入力として、センサーデータの値の確率分布を推定するのは教師なし学習に含まれる。今の例では、この結果に対して「確率が低いものは通常の挙動と異なるので異常である」という判断基準を適用して異常を検出することで機械の故障を検知する。

教師なし学習の手法を用いるには、教師あり学習の場合と異なり正解例付きのデータを用意する必要が無い。正解例を作成する時間とコストが不要であるという点で、教師あり学習よりも適用のハードルは低い。

一方、正解例付きのデータが無い為、教師なし学習を異常検出に用いるには得られた結果を異常と見なすかどうかの判定基準が必要となる。機械の故障検知の例では、教師なし学習の結果として得られるのはセンサーデータの確率分布である。ここからは観測されたセンサーデータが珍しいか否かわからない。これを異常か否かという出力に変換する為には、「確率が低いものは通常の挙動と異なるので異常である」等という判断基準を導入する必要がある。このように、結果に対して解釈を与える必要があるという点では、教師あり学習よりも適用のハードルは高い。

教師なし学習は入力に含まれるパターンや特徴を学習・モデル化するので、例えば、クラスタリング、アソシエーション分析、確率分布の推定、主成分分析、対応分析、正準相関分析、独立成分分析等の手法が含まれる。

教師なし学習を異常検出に用いるには、得られた結果を異常か否かの結論に結び付ける方法が必要となる。この為、異常検出には確率分布の推定がよく用いられる。なぜならば、データの発生確率が分かれば「確率が低ければ異常」という簡単な判断基準を採用できるからである。確率分布の推定を用いた異常検出は多くの場合、入力データの確率分布を推定し、確立分布を用いて新規の入力の発生確率を導き、発生確率が一定以下ならば「典型的な挙動から大きく外れる異常なデータ」と見なす、という手順で実行される。

確率分布の推定に基づく異常検出は外れ値検出と変化点検出の二つに分類される。外れ値検出では、他のデータから値が大きく外れた珍しいものを検出する。一方、変化点検出では、データの時間的な変動を追跡して通常の変化とは異なる珍しい(大きな)変化を検出する。例えば、機械の故障検知において、センサーデータの値の分布を推定して通常値から外れる場合に異常と見なすのが外れ値検出、センサーデータに時間的な因果関係がある(機械は物理的な法則に従う)と仮定して観測値の時間変動を確率モデル化し、通常と異なる時間遷移を検出して「変動パターンが崩れた」という異常を見つけるのが変化点検出である。変化点検出では観測値の大小ではなく変化の仕方を見て異常を検出する。この為、正常時と同じ範囲の値をとるが変化の仕方が正常時とは異なるという、外れ値検出では見つけられない異常を検出できる。ただし、変化点検出は観測値の時間的な順序に意味がある場合にしか使えないので注意が必要である。

教師なし学習の利点として、新規の事例(未知の異常)を検出できることが挙げられる。それまでに蓄積されたデータに対して、そこから大きく外れるものを異常として検出するので、過去の観測に含まれない異常を検出できる。

教師なし学習のもう一つの利点は、正解例を付与したデータを用意する必要がないことである。多くの場合に、収集されたデータに異常かどうかのラベルは自動的に付与されず、教師あり学習を実行するには人手でのラベルの付与が必要となる。この作業には時間とコストが掛かるので、ラベルの付与が不要であることは大きな利点である。

一方、教師なし学習の欠点は、過去の事例に類似するものの検出精度が教師あり学習よりも低いことである。教師あり学習では過去のデータの異常事例と正常事例の分類精度が最大となるように予測モデルを構築する。これに対して、教師なし学習は分類精度最大化を標準としていない為、過去の事例と類似するものの検出精度では教師あり学習に劣る。

### 【3.2.3 教師あり学習と教師なし学習】

統計的異常検出と人手による異常検出の比較と同様に、教師あり学習と教師なし学習は互いに一長一短があり、どちらか片方のみを用いるのが良いというわけではない。両者を相補的に用いるのが望ましい使い方である。特に、未知の異常を検出するには教師なし学習が適しており既知の異常を検出するには教師あり学習が適しているという点で、両者を併用するのが望ましいと言える。

### 【3.3 統計的異常検出の適用先と手法との関係】

統計的異常検出手法は、その数学的な定式化の違いによって分類される。この為、異常検出の適用先や目的に応じて別々の統計的異常検出手法があるわけではない。異なる手法が同じ対象に適用できることもあるし、逆に同じ手法を異なる対象に適用できることもある。従って、何らかの対象について統計的異常検出の枠組を構築した場合に、それがそのまま他の対象にも使えないかを検討してみることに意味がある。

## 4. まとめ

本稿では、アナリティクスの有用な応用先である異常検出について、統計的な異常検出を軸として説明した。統計的異常検出は、大量且つ複雑な観測値から構成されるビッグデータから洞察を得る為の方法論の一つであり、機械学習の手法を用いることで複雑且つ大量なデータに対する自動的な分析を可能とする。

本稿では触れなかった異常検出の先進的な話題として、異常の予兆の検知が挙げられる。何らかの損害を回避するという異常検出の目的からすると異常の発生よりも異常の予兆を検知できる方が望ましい。とは言え、予兆とは何かを定義すること自体が難しく、異常の予兆を検知する一般的な方法はまだない。しかし、予兆検知の問題には多くの研究者が精力的に取り組んでおり、今後の発展が期待される。

Deloitte Analytics 広瀬 俊亮

(注)当該記事は執筆者の私見であり、トーマツグループの公式見解ではありません。

## お問い合わせ先

有限責任監査法人 トーマツ デロイト アナリティクス  
〒100-0005 東京都千代田区丸の内3-3-1 新東京ビル

Tel: 03-6213-1112

e-mail: [tohatsu.analytics@tohatsu.co.jp](mailto:tohatsu.analytics@tohatsu.co.jp) URL: <http://www.tohatsu.com/jp/da>



トーマツグループは日本におけるデロイト トウシュートーマツ リミテッド(英国の法令に基づく(保証有限責任会社)のメンバーファームおよびそれらの関係会社(有限責任監査法人トーマツ、デロイト トーマツ コンサルティング株式会社、デロイト トーマツ ファイナンシャルアドバイザー株式会社および税理士法人トーマツを含む)の総称です。トーマツグループは日本で最大級のビジネスプロフェッショナルグループのひとつであり、各社がそれぞれの適用法令に従い、監査、税務、コンサルティング、ファイナンシャルアドバイザー等を提供しています。また、国内約40都市に約7,100名の専門家(公認会計士、税理士、コンサルタントなど)を擁し、多国籍企業や主要な日本企業をクライアントとしています。詳細はトーマツグループWebサイト([www.tohatsu.com](http://www.tohatsu.com))をご覧ください。

Deloitte(デロイト)は、監査、税務、コンサルティングおよびファイナンシャル アドバイザリーサービスを、さまざまな業種にわたる上場・非上場のクライアントに提供しています。全世界150か国を超えるメンバーファームのネットワークを通じ、デロイトは、高度に複合化されたビジネスに取り組むクライアントに向けて、深い洞察に基づき、世界最高水準の陣容をもって高品質なサービスを提供しています。デロイトの約200,000人におよぶ人材は、"standard of excellence"となることを目指しています。

Deloitte(デロイト)とは、デロイト トウシュートーマツ リミテッド(英国の法令に基づく(保証有限責任会社)およびそのネットワーク組織を構成するメンバーファームのひとつあるいは複数)を指します。デロイト トウシュートーマツ リミテッドおよび各メンバーファームはそれぞれ法的に独立した別個の組織体です。その法的な構成についての詳細は [www.tohatsu.com/deloitte/](http://www.tohatsu.com/deloitte/) をご覧ください。

© 2015. For information, contact Deloitte Touche Tohmatsu LLC

Member of  
Deloitte Touche Tohmatsu Limited