

回帰分析結果検証のための3次元散布図可視化と変数選択

鈴木 千絵[†] 伊藤 貴之[†] 梅津 圭介^{††} 本橋 洋介^{††} 高塚 正浩^{†††}

[†] お茶の水女子大学大学院 人間文化創成科学研究科 〒112-8610 東京都文京区大塚二丁目1番1号

^{††} 日本電気株式会社 〒211-8666 神奈川県川崎市中原区下沼部1753

^{†††} Faculty of Engineering & IT, The University of Sydney School of Information Technologies, J12,
University of Sydney, NSW 2006 Australia

E-mail: [†]chie@itolab.ocha.ac.jp, ^{††}itot@is.ocha.ac.jp, ^{†††}k-umezu@ak.jp.nec.com,

^{††††}motohashi@bk.jp.nec.com, ^{†††††}masa.takatsuka@sydney.edu.au

あらまし 過去の購買データから将来の販売予測をする研究は活発に発表されており、その中には回帰分析を活用した研究が多い。一方で購買データには予測結果に大きく関与する情報と関連性が薄い情報が混在しており、関連性の薄い情報が回帰分析の精度を悪化させることもありえる。本報告では、実測値と回帰分析による予測値との誤差を3次元散布図を用いて可視化する手法を提案し、購買データに適用した事例を示す。本手法では、赤池情報量基準を適用して販売個数や気温などの数値変数を順位付けし、散布図上の点の分布を評価基準として月や曜日などの非数値情報として順位付けする。この順位を反映してユーザインタフェースを構築することで、対話操作による変数の取捨選択を容易にする。

キーワード 可視化, 回帰分析, 高次元データ, ステップワイズ法

1. はじめに

回帰分析は、自然科学や社会科学に関する学術分野や産業分野で活用されている統計的手法の1つである。1つ以上の説明変数と1つの目的変数の関係を数式化し、説明変数から目的変数を推測、予測する。回帰分析を用いた予測は健康状態予測などの医療問題、天災予測やエネルギー需要予測などの環境問題、経済予測や販売予測などの社会問題などが対象となり、その用途は非常に広い。

しかし、複数の説明変数を入力情報とする重回帰分析や、複数の回帰式を導入した混合モデルなどの導入により、その分析工程はより複雑になってきている。特に重回帰分析において、予測に大きく寄与する説明変数と大きく寄与しない説明変数、また予測値と実測値の誤差につながる説明変数を特定することが、回帰分析の性能を向上させるために重要である。

以下、小売店での商品販売を例題として議論する。日常的に販売される商品の売り上げは、その日の気温や曜日、また周辺でのイベント開催など、様々な要因に左右される。販売競争の激しい近年において、過剰発注による廃棄・処分を減らしたり、過剰在庫や完売を防いだりするため、適切な在庫数を保つことが必要不可欠である。そのために商品の販売数やその日の気象情報等の販売データを毎日入力・蓄積している企業は少なくない。取得したデータを解析することで将来の販売数がある程度予測することができるからである。そしてその解析方法の一つが、回帰分析である。

例として予測対象となる売り上げ (f) を目的変数とし、最高気温 (x_1) と最低気温 (x_2) と湿度 (x_3) を説明変数として回帰分析すると、

$$f = a \cdot x_1 + b \cdot x_2 + c \cdot x_3 + d \quad (1)$$

と表せるとする。この式に別日の気象情報を代入すればその日の売り上げを予測できる、というのが回帰分析の考え方である。

しかし取得するデータは膨大になってきており、それらの中には予測結果にほとんど影響しない要因や、逆に予測するためには不可欠な要因が同時に存在している。予測結果に影響しない要因を予測に用いることが、逆にノイズを生むことになり、予測値と実測値との誤差の要因になることがある。予測のための入力情報と予測値に寄与度を理解することは重要であるが、情報の複雑化によってその理解が難しい場合も多い。本報告では各変数が予測値にどの程度影響を与えるのかを理解するため、回帰分析による予測値と実績値との誤差を可視化する一手法を提案する。

2. 回帰分析結果の既存の可視化手法

回帰分析や予測問題の性能を定性的に分析し、その誤差の原因と対策を議論するためのツールとして、可視化は有用であると考えられる。現実の問題におけるデータセットと回帰分析結果との関係を理解することが重要であるにも関わらず、それを目的として新しい可視化システムを開発した研究事例はまだ少ない。

代表的な例として Thomas ら [1] は、回帰モデルの双方向的な構築を支援するため、複雑さが最小限に抑えられる数学的モデルと相関の高い説明変数の組み合わせをより効果的に推薦し、精度の高い回帰分析につながることを表現する可視化ツールを提案した。Krause ら [2] は、回帰モデルを含む予測モデルのための対話型の特徴選択を目的として、グリフベースの可視化ツールを提案した。

それに対して本手法では、回帰分析対象の可視化に3次元散布図を採用している。さらに予測値と実測値の誤差をプロットの色に割り当てることで、誤差が大きくなる標本が集中する3次元空間中の位置を視認しやすくする。

3. 前処理とユーザインタフェース

3.1 データ構造

本手法では以下の説明変数が混在したデータを想定する。

実数値型説明変数: 販売個数や販売日の気温など、実数で表現される情報。

カテゴリ型説明変数: 曜日や物品属性など、実数値で表されない情報。

そして入力データについて以下のデータ構造を想定する。

$$X = \{x_1, x_2, \dots, x_n\}$$

$$x_i = \{v_{i1}, \dots, v_{im}, c_{i1}, \dots, c_{il}, p_i, a_i\}$$

ここで X は標本群, n は標本数, x_i は i 番目の標本を表す。また m は実数値型説明変数の個数, v_{ij} は i 番目の標本における j 番目の実数値型説明変数の変数値, l はカテゴリ型説明変数の個数, c_{ij} は i 番目の標本における j 番目のカテゴリ型説明変数の変数値, p_i は i 番目の標本における予測値, a_i は i 番目の標本における実測値である。

3.2 3次元散布図による可視化

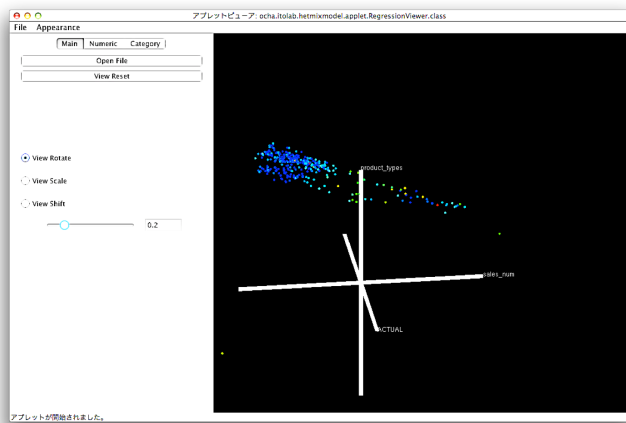


図1 本報告で提案する可視化ツール

本手法では前節で示したデータ構造の可視化に3次元散布図を採用している(図1)。この可視化では、 m 個の実数値型説明変数群の中から2個を選んでx軸およびy軸に割り当て、実績値または予測値をz軸に割り当てる。また、各標本における予測値と実測値の差の絶対値を色で表現している。差の絶対値が大きい標本を赤に近い暖色系の色相で、誤差の小さい標本を青に近い寒色系の色相で描画する。

可視化ツールの画面左側には4つのタブがある。1つ目のタブはファイル操作や描画調節をサポートする。2つ目のタブにはx,y,z軸に割り当てる変数選択のためのラジオボタンを搭載する(図2)。また、後節で示す手法による評価値の高い変数が

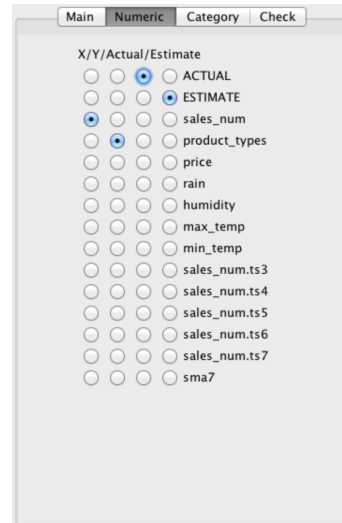


図2 実数値型説明変数の選択タブ

リストの上位に配置される。

3つ目および4つ目のタブではカテゴリ型説明変数の選択をサポートする。3つ目のタブでユーザーが変数を1個を選択すると、そのカテゴリ変数の選択肢となりえる変数値を選択するための4つ目のタブが表示される。例えば3つ目のタブ(図3)で「曜日」というカテゴリ変数を選択すると、4つ目のタブには「日曜」から「土曜」までの7個のチェックボックスが搭載される(図4)。4つ目のタブに搭載された選択肢群のうち、チェックされているカテゴリ変数値をもつ標本は彩度の高い色で描画され、チェックされていないカテゴリ変数値を持つ標本は灰色で描画される。この機能により、誤差分布とカテゴリ変数値の関係を表現可能にしている。

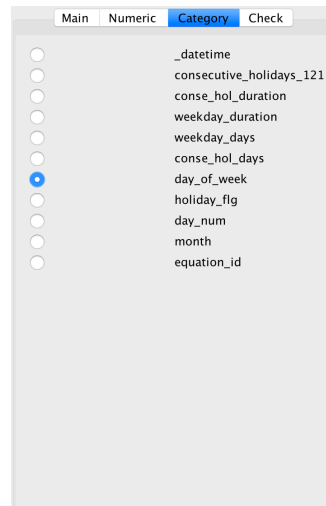


図3 カテゴリ型説明変数の選択タブ

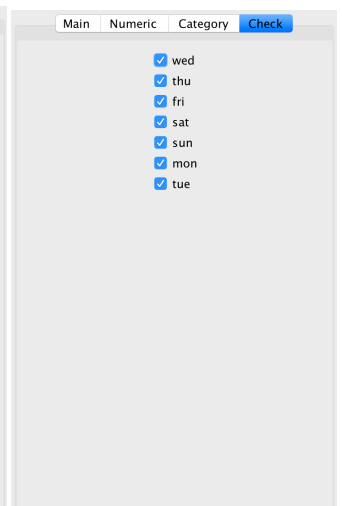


図4 選択したカテゴリ型説明変数の変数値選択タブ

3.3 説明変数の選択

説明変数が非常に多いデータを3次元散布図で可視化する場合に、どの説明変数をx,y軸に割り当てるかによって可視化の効果は大きく変わる。よって効果的な可視化を実現する説明変

数の選択が重要となる。また回帰分析の性能を向上する一手段として、重要でない説明変数を削除することがあげられる。一般的に回帰式は、目的関数への影響がある説明変数を多数導入したほうが、訓練データへの適合率は高まる。しかしノイズなどの異常値にも影響されるため、過剰適合が生じる場合がある。変数を削減することで過剰適合を防ぐことができるが、過度な変数削減はそのデータへの適合率の低下を招く。そのため、説明変数の各々について削除したほうがいかに議論するためにも、3次元散布図での可視化における説明変数の選択は重要である。

そこで前処理として、各説明変数および変数群について予測値への寄与と誤差への要因を評価し、各説明変数または変数群の興味深さを定量的にユーザに提示することが有用である。本報告では赤池情報量基準 (Akaike's Information Criterion; AIC) [3] [4] をもとに各実数値型説明変数を評価し、標本の分布をもとに各カテゴリ型説明変数を評価する。

AIC は訓練データへの当てはまりの悪さと複雑さを数値化したもので、以下の公式で表される。

$$AIC = -2 \cdot \log L + 2 \cdot k \quad (2)$$

ここで L は最大尤度、 k は自由パラメータの数である。AIC 値が最小となる変数を選択することで、多くの場合、良質な予測を実現できるモデルを選択できることが知られている。

カテゴリ型説明変数は、実績値と予測値の誤差が大きい点の乱雑さをもとに評価する。ここで、標準偏差を基準として任意の数値以上の誤差がある標本を「誤差が大きい」とする。x-y 平面を格子状に分割して各領域における「誤差が大きい標本」の個数を集計し、そのエントロピーが小さいカテゴリを「誤差が大きい標本が画面上で局所集中している特徴的な可視化結果を導く変数」とみなし、高く評価する。

3.4 説明変数の順位付け

本手法では前節で示した評価基準をもとに、実数値型説明変数およびカテゴリ型説明変数組を順位付ける。以下に手順を示す。

【step1】 全ての変数から変数を1つ除いて評価する。

【step2】 最も評価値の高い変数組を決定する。

【step3】 step2 の変数組から変数を1つ除いて評価する。

変数を順位付ける場合、

【step4-1】 step2, 3 を繰り返し、変数が1個になったら終了する。削除した変数の逆順を、実数値型説明変数の順位とする。変数組を順位付ける場合、

【step4-2】 step2, 3 を繰り返し、評価値が上がらなくなったら終了する。残った変数組の逆順を、カテゴリ型説明変数組の順位とする。

4. 販売情報の回帰分析結果への適用

販売情報の回帰分析結果について本手法を適用した。入力データには、実測値と予測値を含む344サンプル、12個の説明変数、および8個のカテゴリ変数が含まれている。回帰分析には異種混合モデル [5] [6] が適用されている。説明変数の名前

は機密情報であるため、本報告では各実数値型説明変数を A から L のアルファベットで示す。カテゴリ変数には、取得した月、日、曜日、および式番号が含まれている。入力データに対して AIC 値を求め、AIC 値が小さくなる2つの実数値型説明変数 (説明変数 A, B) を x 軸および y 軸に割り当て (図 5)、誤差の大きい点群のエントロピーが小さくなるカテゴリを選択して可視化した (図 6)。

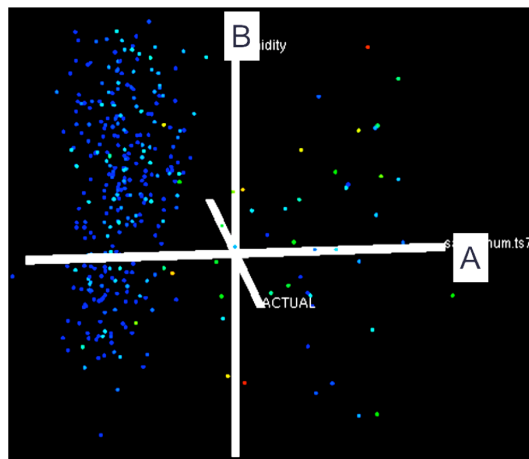


図 5 実数値型説明変数 (説明変数 A, B) の選択後の可視化結果

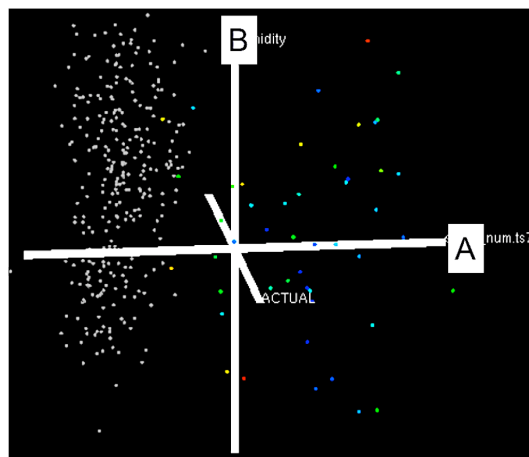


図 6 曜日カテゴリの選択後の可視化結果

続いて曜日カテゴリを選択し、月曜日のみを選択した (図 6)。図 6 より、A と B が大きいときに暖色点が多いことがわかる。

また月カテゴリを選択し、2月と9月を選択した可視化結果を (図 7) に示す。図 7 では、A が小さいときの寒色点は図 6 と比較して多い。しかし図 6 と図 7 のどちらにも高彩度で描画されている点に着目すると、寒色点の数が減少し、暖色点の割合が増加したことがわかる。

また説明変数 A, C を x 軸および y 軸に割り当て (図 8)、誤差の大きい点群のエントロピーが小さくなるカテゴリを選択して可視化した (図 9)。

説明変数 A, B を x 軸および y 軸に割り当てた場合と同じカテゴリを選択した。曜日カテゴリから月曜日のみを選択した可視化結果を図 9 に示し、また月カテゴリを選択し、2月と9月

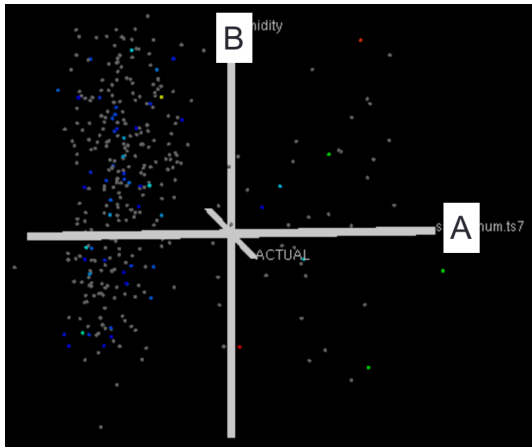


図 7 月カテゴリの選択後の可視化結果

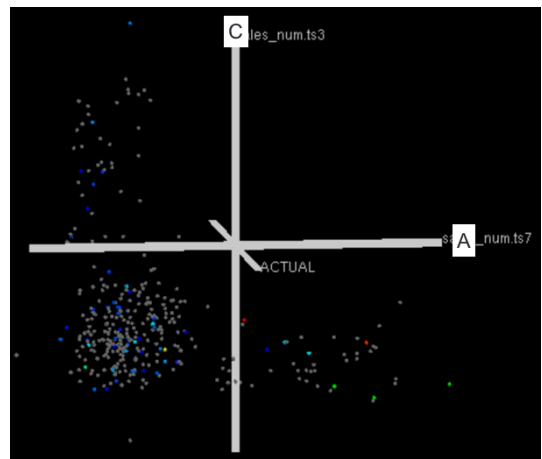


図 10 月カテゴリの選択後の可視化結果

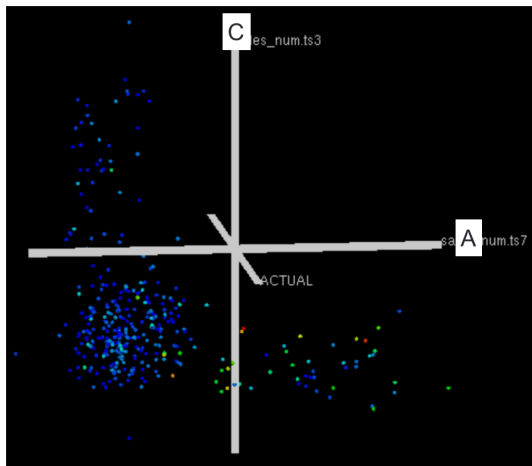


図 8 実数値型説明変数 (説明変数 A,C) の選択後の可視化結果

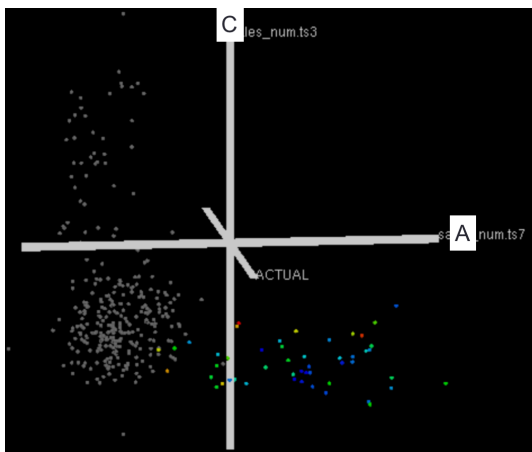


図 9 曜日カテゴリの選択後の可視化結果

のみを選択した可視化結果を図 10 に示す。

図 8 より、A が大きく、C が小さいときに暖色点が多いことがわかる。また各カテゴリ選択結果 (図 9, 図 10) より、寒色点の数が減少したことがわかる。

月曜日に誤差が生じる原因として、月曜日は祝日の代休となることが多く、平日と祝日を同様に扱っていることが考えられる。また 2 月と 9 月に誤差が生じる原因として、小売店ではこれらの月に商品の入れ替えを行うことが多いことが考えられ

る。その日が祝日か否かによって使用する回帰式、および季節の変わり目であるこれらの月に使用する回帰式を使い分けることで、精度向上につながる事が示唆される。

5. まとめ

本報告では、回帰分析の結果検証のための 3 次元散布図ツールと説明変数選択の評価方法を提案した。本手法では実数値型説明変数のうち 2 つを選んで x 軸と y 軸に割り当て、実測値を z 軸に割り当てることで、3 次元散布図を実現する。そして予測値と実測値の誤差を色で表現することで、誤差の大きい標本がどのように分布するかを視覚的に観察できる。またカテゴリ型説明変数のうち選択した変数値を持つ標本を高彩度で描画することで、誤差の要因をなり得るカテゴリ型説明変数と変数値を特定できる。ここで説明変数を選択するための一手段として、我々は AIC を用いて実数値型説明変数の評価を実施し、標本の分布をもとに各カテゴリ型説明変数の評価を実施した。その結果として、有用な実数値型説明変数を選択して 3 次元散布図の 2 軸に割り当て、有用なカテゴリ型説明変数の変数値を選択したことで、誤差の小さい要素と誤差の大きい要素を視覚的に分離できるような 3 次元散布図を実現でき、誤差が生じる詳細な要因の一つを特定できた。適用例では、実数値型説明変数 A,B が大きいときおよび実数値型説明変数 C が小さいときに誤差が大きく、曜日や月が誤差が生じる詳細な要因の一つであると特定した。

現在のツールでは同じカテゴリ内の変数値の組み合わせ (例えば曜日カテゴリの月曜と火曜) のみをカテゴリ型説明変数組の選択肢としている。今後は異なるカテゴリ間の変数値の組み合わせの選択 (例えば月カテゴリの 2 月と曜日カテゴリの月曜日) を可能にし、誤差が生じる要因の特定をより容易にするためのツールの改良を行う。また、3 次元散布図による可視化結果を最も効果的な視線方向から表示するための評価基準の検討と、その自動化を目指す。

文献

- [1] T. Muhlbacher, H. Piringer, “A Partition-Based Frame-

work for Building and Validating Regression Models,” IEEE Transactions on Visualization and Computer Graphics, 19(12), pp. 1962–1971, 2013.

- [2] J. Krause, A. Peter, E. Bertini, “INFUSE: Interactive Feature Selection for Predictive Modeling of High Dimensional Data,” IEEE Transactions on Visualization and Computer Graphics, 20(12), pp. 1614–1623, 2014.
- [3] H. Akaike, “Information Theory and an Extension of the Maximum Likelihood Principle,” Second International Symposium on Information Theory, B. N. Petrov & B. F. Csaki (Eds.), pp. 267–281, 1973.
- [4] H. Akaike, “Maximum likelihood identification of Gaussian autoregressive moving average models,” Biometrika, 60(2), pp. 255–265, 1973.
- [5] R. Fujimaki, S. Morinaga, “Factorized Asymptotic Bayesian Inference for Mixture Modeling,” International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 400–408, 2012.
- [6] R. Eto, R. Fujimaki, S. Morinaga, H. Tamano, “Fully-Automatic Bayesian Piecewise Sparse Linear Models,” International Conference on Artificial Intelligence and Statistics (AISTATS), pp. 238–246, 2014.