

Integrated Gradientsによる格付AR値変化の寄与度分解

賛助会員 (株)日経金融工学研究所 *木村和央 KIMURA Kazuo

1. はじめに

ロジットモデルに代表されるクラス分類を目的としたモデルの精度評価指標として、AUC, AR値といった順序統計量(序列精度)が重要視されている(山下・三浦[1]). このうち、金融機関においては伝統的にAR値の利用が多く、その時系列変化に大きな関心が寄せられている。木村ほか[2]では、格付AR値(以下AR値)を対象に、格付別の構成比と相対デフォルト率の変化差分を用いた寄与度分解を提案したが、相対デフォルト率には全体のデフォルト率の影響が入り込むため結果の解釈に困難な部分があった。この困難をクリアするべく、本稿では説明可能なAI(XAI)技術の1つであるIntegrated Gradients(IG)を利用して、AR値の変化要因をセグメント・格付別の構成比変化とデフォルト率変化に寄与度分解する手法を提案し、実際の債務者データに対して適用した。

2. Integrated Gradients(IG)の概略

詳細はSundararajan et al.[3]を参照。基準となるサンプルの入力値を \mathbf{x}' 、評価したいサンプルの入力値を \mathbf{x} とする。入力値が \mathbf{x}' から \mathbf{x} まで変化するとき、出力値は $F(\mathbf{x}')$ から $F(\mathbf{x})$ まで変化するが、その差である $F(\mathbf{x}) - F(\mathbf{x}')$ を次式により入力値の各要素 x_i の変化からの寄与度に分解する。

$$IG_i(\mathbf{x}) ::=$$

$$(x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(\mathbf{x}' + \alpha \times (\mathbf{x} - \mathbf{x}'))}{\partial x_i} d\alpha.$$

通常、解析的に計算するのは困難であるため、 m 個の区間に分割して近似した区分求積法にて行う。

$$IG_i^{\text{approx}}(\mathbf{x}) ::=$$

$$(x_i - x'_i) \times \sum_{k=1}^m \frac{\partial F(\mathbf{x}' + (k/m) \times (\mathbf{x} - \mathbf{x}'))}{\partial x_i} \cdot \frac{1}{m}.$$

なお、 F がほとんど至るところ微分可能であれば、寄与度の合計が出力値の差分に一致する。

$$\sum_i IG_i(\mathbf{x}) = F(\mathbf{x}) - F(\mathbf{x}').$$

*本稿の内容は筆者に属し、所属組織の見解ではない。

3. 格付AR値の定義

業種等のセグメント $s(= 1, \dots, S)$ 、格付 $k(= 1, \dots, K)$ の債務者数の全債務者数に対する構成比を a_{sk} 、デフォルト率を p_{sk} とすると、全体のデフォルト率は $P = \sum a_{sk} p_{sk}$ 。このとき、AR値(AR)は以下のとおり。

$$AR = \frac{1/2}{(1-P)P} \sum_{u,v} \sum_{i,j} a_{ui} a_{vj} (p_{vj} - p_{ui}) \text{sgn}(j - i).$$

4. IGに必要な微分計算とその解釈

構成比 a_{sk} とデフォルト率 p_{sk} が入力値 \mathbf{x} 、AR値(AR)が出力値 $F(\mathbf{x})$ と考えて、IGに必要な微分計算を行う。

$$\begin{aligned} \frac{\partial AR}{\partial a_{sk}} &= AR \cdot \frac{2P - 1}{(1 - P)P} \cdot p_{sk} \\ &\quad + \frac{1}{(1 - P)P} \sum_v \sum_j a_{vj} (p_{vj} - p_{sk}) \text{sgn}(j - k). \\ \frac{\partial AR}{\partial p_{sk}} &= AR \cdot \frac{2P - 1}{(1 - P)P} \cdot a_{sk} \\ &\quad - \frac{1}{(1 - P)P} \sum_v \sum_j a_{sk} a_{vj} \text{sgn}(j - k). \end{aligned}$$

右辺第1項は、全体デフォルト率 P の変化を通じてAR値に与える影響と解釈できる。通常、 $P < 0.5$ より、 $2P - 1 < 0$ であるから、 a_{sk} と p_{sk} の増加/減少はARの減少/増加に寄与する。

右辺第2項は、格付間の整序性の変化を通じてAR値に与える影響と解釈できる。格付別デフォルト率が $j \leq k$ のとき $p_{vj} \leq p_{sk}$ であれば、 a_{sk} 微分の式の右辺第2項は正となり、 a_{sk} の増加/減少はARの増加/減少に寄与する。他方、 p_{sk} 微分の式の右辺第2項は、 \sum 記号内の項別に、 $k > j$ ならば正、 $k < j$ ならば負となる。つまり、格付 k のデフォルト率上昇/下落は、上位格付とのデフォルト率差を拡大/縮小させるためARの増加/減少となり、下位格付とのデフォルト率差を縮小/拡大させるためARの減少/増加となる。このバランスにより右辺第2項の正負が決まる。

5. 実際の債務者統合データへの適用

債務者統合データ¹へ本手法を適用し、2004～17年度のAR値の時系列変化につき寄与度分解を実施した。セグメント（業種）は、1. 製造、2. 建設、3. 卸売、4. 小売、5. 不動産、6. サービス、7. その他である。格付は、決算データのスコアリング結果から付与したもので、最上位格が1、最下位格が7である。なお、区間分割数を1,000とした台形公式にて区分求積計算を実施した。以下、構成比変化とデフォルト率変化の寄与度は合算して扱った。

図1は業種別、図2は格付別の寄与度分解である（以下、数値を1,000倍して表示していることに留意）。全体のAR値は、2006年度とリーマンショックの08年度に大幅低下したが、その後は回復基調であった。データ秘匿のため図は省略したが、2008年度の業種別に計算されたAR値は、不動産業が他の業種の2倍程度低下しており、一般に、全体のAR値低下は不動産業の要因であると解釈されているが、図1によれば製造業の寄与度が上回った。これは債務者数に占める製造業の割合が不動産業よりも多いことに起因するものと考えられるが、本手法により定量的な比較・解釈が可能となった。また、前回の相対デフォルト率を用いた寄与度分解では、最下位格付であるRank7の寄与が重要と報告していたが、図2によれば中位格付であるRank3,4からの寄与度が顕著であり、これも債務者数割合で説明が可能と考えられる。表1は2008年度における業種×格付別の寄与度であるが、製造業はRank1の寄与度が大きく、Rank2～4も同様の傾向であるのに対し、不動産業はRank3と4の寄与度が大きいことがわかった。図3は時系列分散の寄与度分解であり、業種ではその他、格付では低格付であるRank6,7が全体の動きとは逆相関となった。

6. おわりに

本稿では、Integrated Gradientsを用いて格付AR値の変化につき寄与度分解する手法を提案し、実際の債務者統合データに適用した。今回採用した手法により、相対デフォルト率から通常のデフォルト率に変更したことで結果解釈がわかりやすくなり、また、業種セグメントに関する寄与度分解も合わせて可能となった。

¹データは所属機関が独自に金融機関から許可を得て収集。

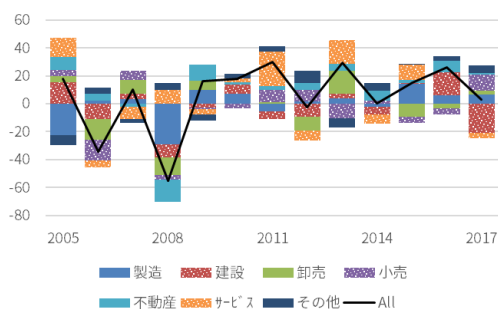


図1: 格付AR値変化の業種別寄与度分解

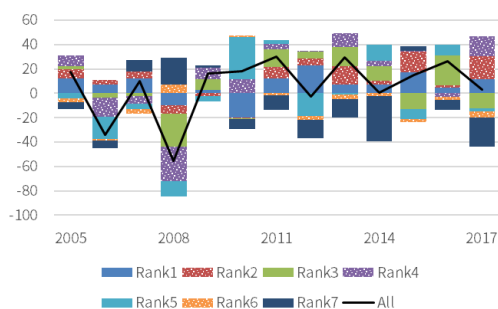


図2: 格付AR値変化の格付別寄与度分解

表1: 2008年度の業種・格付別寄与度分解

	製造	建設	卸売	小売	不動	サー	その他
Rank1	-8.5	1.5	-2.4	-0.3	-0.2	0.1	0.0
Rank2	-6.7	-0.2	-2.0	-1.0	-1.0	1.6	2.6
Rank3	-7.5	-2.3	-3.1	-4.4	-7.8	-3.2	1.2
Rank4	-6.6	-9.1	-4.0	0.1	-9.0	2.0	-1.6
Rank5	-3.1	-3.6	-4.3	1.0	1.0	-4.7	1.0
Rank6	0.3	1.7	1.2	1.5	0.5	1.8	0.3
Rank7	2.8	2.9	1.9	0.0	0.3	12.5	1.3

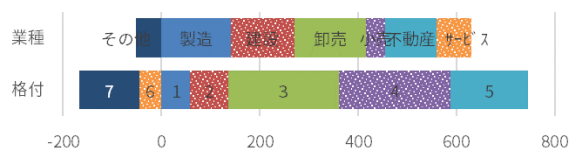


図3: 時系列分散の業種/格付別寄与度分解

参考文献

- [1] 山下智志, 三浦翔 (2011). 信用リスクモデルの予測精度- AR値と評価指標-. 朝倉書店.
- [2] 木村和央, 宋明子, 友添峻希 (2019). 格付AR値変化の寄与度分解. 日本オペレーションズ・リサーチ学会 2019年秋季研究発表会予稿集.
- [3] M Sundararajan, A Taly, Q Yan (2017). Axiomatic attribution for deep networks. International conference on machine learning, 3319-3328.