

# AIガバナンス構築の肝となる モデルの独立検証の重要性

既存のモデル管理の枠組みにおける  
検証の視点を有効に活用せよ

AI（人工知能）の活用が広がると同時に、AIの管理（AIガバナンス）の重要性が高まっている。海外の金融機関では、AIガバナンスはモデル管理（Model Risk Management = MRM）の枠組みに組み込むことが一般的である。本稿では、まずMRMの枠組みを生かしたAIガバナンスの考え方とその利点について触れつつ、モデルの独立検証の重要性を述べる。その上で、独立検証では、どのような点に留意してAI（含む生成AI）の検証を行うべきか、個別のAIを具体例に挙げながら考え方を示す。さらに、今後のAIガバナンスのポイントにも簡単に触れたい。

## AIガバナンスの 土台としてのモデル管理

金融機関においてAIの活用が広がっている。その範囲も定型的な業務から、効率化・高度

化を狙った業務まで拡大し、さまざまな用途で活用されている。

金融庁も今年3月に公表した「AIディスカッションペーパー」（第1・0版）において、  
(AI活用で)「チャレンジしないリスクも十分に認識される

シニアマネジャー

竹川 正浩



あづさ監査法人  
金融統轄事業部  
アドバイザリー統轄事業部  
ディレクター

田中 康浩



# AIガバナンスの下でのモデル検証

あり、モデル管理（MRM）の枠組みを活用すべきという考え方が一般的である（注1）。MRMでは、AIの開発者・使用者である1線に対し、1線から独立した2線（モデル検証者）が検証を行う。

米国のMRMに係る監督は厳格である。FRB（連邦準備制度理事会）とOCC（通貨監督庁）のガイダンス（注2）が公表されて以降、金融機関は相当のリソースをかけてMRMの態勢を整備・高度化してきた。MRMでは、モデル開発からモdel特定期付与・検証・承認・使用中のモニタリング・再検証、使用停止に至るライフサイクルに応じて、1線と2線の役割と責任を定め、モデルの管理を求めている。

MRMが経営管理に深く浸透している米国では、AIをMRMで管理するという発想は合理的である。一方、本邦金融機関では、独立検証を経てAIを使用するという発想は乏しく、AI（モデル）をMRMで管理するという考え方があまり浸透していない。

## モデル管理の利点は三つの「仕組み」

MRM活用の利点は、以下の三つの「仕組み」がすでにMRMに備わっていることにある。

一つ目は、ライフサイクル管理である。AIにもモデルのライフサイクルと同じ概念が存在する。

具体的には、構想から開発の際にテストが行われ、承認を受けて使用される。使用中も、挙動などに関するモニタリングが行われる。こうしたライフサイクル管理は、MRMが得意とする分野である。すなわち、MRMを活用すれば、新たなAI管理の仕組みは不要である。

二つ目は、リスク格付けとインベントリー（一覧表）である。モデルのリスク格付けは、リスクベース・アプローチの基礎となる重要な概念である。米国G-SIBs（グローバルな金融システム上重要な銀行）ではモデルは数千以上あり、すべてについて厳格な管理を行うことは不可能である。管理にめりりりを付けるためにも、AIにリスク格付けを付与することは重要である。

また、格付けと共にAIをインベントリーに登録しておけば、使用中のAIの「見える化」につながる。これは、AI管理のインフラを提供するものだといえる。

三つ目は、検証・承認である。あらゆるモデルは2線による独立検証・承認を経て使用される。

モデル検証では、開発者にとって都合の良いモデルになつていいなか、リスクや弱点・制約、使用上の留意点は何かなど、多面的な評価がなされる。米国ではモデル検証は厳格であり、リスクの高いモデルであれば10ページを超える評価書（検証報告書）を作成し、承認を行う（否認や条件付き承認もあり得る）。

MRMで最も重要な仕組みが、独立検証である。これによって安心・安全なAI活用を担保できる。

MRMの独立検証では、あらゆるモデルに対して、データ、機械学習アルゴリズムのランダムフォレストは、複雑な構造を持つため、伝統的な統計モデル

## 従来型AIで論点となる説明可能性や公平性の担保

本邦ではMRMでAIの検証を行う事例はまだ少数である。そこで次に、どのような点に留意してAIの検証を実施すべきか、従来型AIと生成AIのそれぞれについて、イメージしやすいよう架空の例を交えて考え方を示す（注5）。

まず図表1のような従来型AIを取り上げる。

このタイプのAIが採用する機械学習アルゴリズムのランダムフォレストは、複雑な構造を持つため、伝統的な統計モデル

のように特定の出力に至った理由を直接的に説明することはできない。しかしSHAP (SHapley Additive exPlanations) 値やLIME (Local Interpretable Model-agnostic Explanations) 値から出力に対する特徴量の貢献度を知ることで、ランダムフォレストの説明可能性を一定程度、評価することができる。

また、伝統的な統計モデルやほかのAIを用いて、モデル構造の理解のしやすさと精度の両面から比較評価を行うことも有用である（ベンチマーкиング）。モデル精度が大きく変わらなければ、説明可能性を優先して伝統的な統計モデルの利用を促す結論に至ることもある（この場合、伝統的な統計モデルがチャンピオンモデルになり、AIがチャレンジャーモデルになる）。

公平性の観点では、人種・国籍・性別等の情報を特徴量として採用すると、不当な差別や不利益が生じる場合がある。こうした点を防ぐために、属性グループ別にスコアの分布を確認し、特定のグループが不当に偏った評価になつていなかを分析す

ることが必要になる。AIの出力についても、公平性が保たれているかを検証することが重要である。

〔図表1〕 従来型AI（例）の概要

|        |  |
|--------|--|
| 名称     | 取引データに基づく個人の信用スコアリングAI   |
| リスク格付け | 高（AIの手法を用いており複雑であり、審査という重要な用途で活用）  |
| 使用部署   | 個人ローン審査部門  |
| 入力データ  | 預金口座の動きや属性情報、給与振り込み、公共料金支払い等   |
| 出力     | 信用力スコア   |
| AIの特徴  | 決定木を発展させたランダムフォレスト等のAIを採用するケースが多い。ランダムフォレストは多数の決定木から構成され、個々の決定木のスコアを集約することで最終的な出力を得る。特徴量と出力の非線形な関係を表現することができる。また、オーバーフィッティングしにくく、外れ値にも強い |

（出所） KPMGジャパン作成（図表2も同じ）

〔図表2〕 生成AI（例）の概要

|        |   |
|--------|---|
| 名称     | 融資業務における財務分析AI  |
| リスク格付け | 中（AIの手法を用いており複雑であるが、分析の参考程度の活用）   |
| 使用部署   | 融資部門  |
| 入力データ  | 分析対象企業名、決算年度、過去の財務分析結果、同業他社の財務分析結果等   |
| 出力     | 当該年度における財務分析結果  |
| AIの特徴  | ドラフトエージェントが入力データに基づき財務分析を実施し、ドラフトを作成。評価エージェントがレビューを行い、修正を指示。修正不要と判定されると、ドラフトエージェントが全体を統合した財務分析リポートを作成 |

（注）本稿ではエージェント機能を持つ生成AIをイメージして例示している。

## 一層高い生成AIの検証

一方の生成AIは、従来型AIよりさらに複雑かつブラックボックス性が高く、確立されたここで挙げた検証の視点はロジック等の手法やテストの一部であるが、リスク格付けがA Iシステム）で階層的に検証するアプローチを紹介する。

従来型AIと同じく架空の例として、図表2のような生成AIをイメージする。

まずモデルレベルの検証では、ドラフトエージェントおよび評価エージェントの中核となる生成AIモデル単体の性能を評価

する。性能評価指標は、さまざま

だらう。

まなベンチマークが存在する。

例えば、幅広い学問分野や一般的な知識の理解の評価、数学の問題を解く能力の評価、コード

生成能力の評価、金融データに

対する回答能力の評価等である

(注6)。今回は財務分析での利用が想定されるため、数学の問題を解く能力や、金融データに対する回答能力の評価に係る指標に重点を置くべきだらう。

また、リスクが高くないモデルについては、検証の負荷を減らす目的で、第三者が公開しているモデルの順位表を参照して、簡易的に検証を行うことも一案である。なお、基盤モデルの多くがベンダーモデルであるとい

う現状を踏まえると、「使用目的に適合したベンダーモデルが選定されているか」もポイントになる。ベンダーが公表している検証報告書の確認が中心になるが、公表内容が限定的である点に留意が必要である。

以上のようなモデル自体の性能のほか、契約や利用条件(入力したデータが再学習に利用されない等)に関する評価も重要

## リスクに備えて 人による判断が不可欠

続いてタスクレベルの検証で

は、タスクを担うエージェント単位で評価を行う。モデルに加えRAG(検索拡張生成)やガードレール等の要素が適切に組み合わされて、自律的にタスクを遂行できるかを検証する。適合率、再現率、F1スコア(注7)といった従来型A.I.でも使用されている指標に加えて、タスク固有の指標も加味して評価を行う。

今回は、過去に人が行つた財務分析結果を正として、どれだけ近い出力が得られたかを意味的類似度指標で確認することが案だらう。また、ハルシネーションのほか、タスクの内容次第では著作権やプライバシーの侵害等についても評価が必要となる。

システムレベルの検証では、ユースケース単位で評価を行う。財務分析レポートの作成というユースケースの場合、財務分析

かを検証する。

「期待どおりの財務分析リポートが生成されたか」といった、入力データから出力まで(End-to-End)の成功率の評価も重要な要素である。財務分析リポートが期待外れであった場合、その原因究明のためシステムの動きを詳細に把握する必要があり、監視やログ収集機能等を有するツールを活用することは有用だらう。こうしたツールは、使用中のA.I.のモニタリングにも有効である。

以上のような階層的な検証によってリスクの軽減を図ることができる。なかでも生成A.I.は、ユースケース単位での評価がボイントになる。だが、A.I.の潜在的なリスクを完全になくすことは非現実的である。財務分析結果の確定や結果を踏まえた融資判断といった重要なステップには、人による判断(Human in the Loop)が欠かせない。

やドラフト作成、ドラフトのレビューといった個々のタスクを適切に定義できているか、それらのタスクを実行するエージェントを適切に組み合わせているかを検証する。

また、ベンダーモデルはすべての情報が開示されるわけではないため、モデルレベルの検証が限定的にならざるを得ない。モデルが頻繁にアップデートされたり、サイバー攻撃にさらされたりする可能性があることも踏まえると、使用中のA.I.モニタリングの重要性は高くなる。

さらに生成A.I.では、ITシステム、コンプライアンス、リガル、データ等の専門知識がより広く深くなる。MRM部署が中心となり、これらの関連部署と連携を取りながら検証を進めることが重要である。

## 過度なブレーキを かけない発想を

米国では、トランプ政権の発足を受けて、規制・監督の強度が弱まる方向である。一方でMRMについては、頑健なMRMや独立検証の仕組みが、安心・安全なモデルやA.I.の活用に寄与しているとの声が聞かれている。

F.R.B./O.C.C.のガイドンスを土台とするMRM自体は、A

I（特に生成AI）ガバナンスの最適解になるわけではない。

理想は、本稿で示したような検証の視点等をMRMに取り込んでいくかたちである。MRMの重要性は、AI時代においてむしろ高まると考えられる（注8）。

半面、AIをMRMに組み込むことのデメリットにも留意が必要である。MRMは厳格な管理を要求するために、AIの活用にブレーキがかかる可能性がある。

活用とそれに伴うリスク管理のバランスを取るために、リスクの高いAIは厳格な独立検証を行う一方、リスクの低いAIは管理の強度を緩める、または管理から外すという考え方もある。生成AIについては、個人業務や分析のアシスタント等に使用が限定されるのであればなおさらである。AIの活用に過度なブレーキを掛けない範囲でMRMに組み込む発想が求められる。

\* \* \*

AIの進化はとどまるところを知らない。AIガバナンスの

観点では、伝統的なMRMだけではAI（特に生成AI）に対応し切れないことも事実である。

関連部署との連携に加え、AIの進化に合わせてMRMも進化が必要だろう（将来、MRMエンジニアののような概念が出現するかも知れない）。

AIは金融以外でも、あらゆる業態・ビジネスで活用されている。金融やMRMを超えて、さまざまな業態・ビジネスのAI活用事例や、それを受けたリスク管理の理解も重要な。

本稿で紹介したMRMや独立検証の視点が、金融機関のみならずさまざまな業態におけるAIガバナンスの進化および安心・安全なAI活用の一助になれば幸いである。

4 AI特有のリスク等に関する詳細は、「KPMG Trust ed AI」フレームワークを参照。

5 従来型AIと生成AIについては、金融庁「AIディスクッションペーパー（第1・0版）」の6～7ページの記載が参考になる。

6 生成AIの進化のスピードはさまざま、それに合わせて性能評価指標も日々さまざまなものが出現している。AI開発者・検証者等の実務者は、評価に際して最適な性能指標等の最新の動向をフォローすることが重要になる。

7 適合率と再現率のバランスを示す指標。

8 最新のFSBのペーパー

MANAGEMENT」（2021年、SR11-7）

3 Conceptual Soundness→ Outcome Analysis→モデル検証の中で中核を成す概念であり、前掲注2（SR11-7）で詳しく述べられている。車に例えると、前者は車種に見合った「エンジンの設計」について評価し、後者は「車が適切に動くか」の評価といえる。

4 AI特有のリスク等に関する詳細は、「KPMG Trust ed AI」フレームワークを参照。

5 従来型AIと生成AIについては、金融庁「AIディスクッションペーパー（第1・0版）」の6～7ページの記載が参考になる。

6 生成AIの進化のスピードはさまざま、それに合わせて性能評価指標も日々さまざまなものが出現している。AI開発者・検証者等の実務者は、評価に際して最適な性能指標等の最新の動向をフォローすることが重要になる。

7 適合率と再現率のバランスを示す指標。

8 最新のFSBのペーパー

（Monitoring Adoption of Artificial Intelligence and Related Vulnerabilities in the Financial Sector））もMRMの重要性について触れている。

たなか やすひる  
京都大学経済学部卒。日本銀行やボストンコンサルティンググループを経て、17年KPMG/あずさ監査法人。モデル管理・検証やAIガバナンスに関するアドバイザリー業務を統括。金融監督局に出向し、「モデル・リスク管理に関する原則」を策定。著書に『詳解金融機関のためのモデル・リスク管理』（編著）や『マネロン対策等の新論点』（共著）。

たけかわ まさひる  
大阪大学大学院工学研究科修了。信用リスク関連のデータベース機関等を経て、19年KPMG/あずさ監査法人。信用リスク管理領域を中心にAI・機械学習等を踏まえたアドバイザリー業務を提供。本邦大手金融機関に向け、リスク管理業務に従事。

1. (注)1 KPMGジャパン「本邦金融機関のAIガバナンスの在り方」を参照。なお、MRM全般については、田中康浩・曾我部淳編著『詳解金融機関のためのモデル・リスク管理』（中央経済社、2024年）が詳しい。

2 「SUPPLEMENTARY GUIDANCE ON MODEL RISK