

勾配ブースティング決定木の可視化に基づく インタラクティブな開発・運用支援手法

柏山 美結¹ 廣川 暢一² 伊藤 貴之¹

概要：

機械学習を用いた業務の効率化は、需要予測の自動化などの形で活発に実用化されている。一方で、運用されている予測モデルは年々複雑化している。複雑なモデルを解釈するためには、予測に使用される訓練データの理解に加えて、モデル自体の構造やモデルに特定のデータを入れた時の振る舞いを把握することも重要である。そこで本研究では、実際に運用段階で頻繁に使用されているアンサンブル学習モデルの学習メカニズムを3次元空間で可視化しモデル構造の理解を支援することで精度劣化原因の特定を促すインタラクション手法を提案する。本報告では Housing Dataset への適用事例から提案手法の有効性を議論する。

1. はじめに

機械学習技術を用いて開発された予測モデルは、幅広い産業分野にて頻繁に運用されている。開発されたモデルは年を追うごとに巨大・複雑になり、開発者でも全貌を把握するのが困難になりつつある。そこで、機械学習モデルの構造の分析や評価を支援する手法として情報可視化技術が注目されている。Linhao[1] および Zijie ら [2] は、学習データの性質やモデルの性能等から複数のモデルが比較できるような可視化システムを提案した。また、Nagasaka ら [3] は、没入型環境における深層学習モデルの可視化システムの提案をしており、3次元空間での可視化の有効性を示している。

本研究では、運用段階において広く利用されているアンサンブル学習の勾配ブースティング決定木を対象とし、モデルの構造理解を支援するための3次元空間での可視化の一手法を提案する。また、予測モデルの解釈性を向上させるためのインタラクション手法についても議論する。

2. 関連研究

機械学習 (ML) モデルの解釈可能性と予測精度を高めることを目的とした ML モデル可視化に関する研究は、既に多く発表されている。Kovalerchuk ら [4] は、機械学習のための新たな2つの決定木可視化手法を提案した。属性間の関係や決定木構造、決定木内のデータフロー等を観察し分析することで、モデルの過度の一般化やオーバーフィッ

ティングを防ぐことに有効であることが報告されている。Chatzimparmpas ら [5] は、多様なアンサンブル学習ベースの ML モデルからの意思決定支援を目的として、可視化分析システムを提案した。各モデルの性能や重要な特徴量、訓練データとテストデータの比較結果を表示し、モデルについての意思決定を行う ML エンジニアとドメイン専門家の共通理解や共同フローにより、モデルの分析、調整を可能にしている。

3. 提案手法

3.1 アンサンブル学習モデル可視化

本論文では勾配ブースティング決定木 (GBDT) を対象とする。GBDT は、最初の決定木を基に、弱学習器が持つ誤差を順に学習していくことで全体の予測精度を高める手法であり、様々な問題に適用できることから幅広く用いられている。しかしながら、従来の決定木の可視化手法では個々の弱学習器の構造やデータフローの分析は容易であるが、GBDT のように弱学習器間に直列的な関係性がある場合、その関係性の可視化は想定していない。また、アンサンブルモデルを構成する弱学習器の数が膨大になると、全ての弱学習器を可視化することは現実的ではない。そこで本研究では、予測に使用されるデータの特徴量やモデルの性能結果だけでなく、弱学習器間関係性やデータセットの違いによるモデルの学習過程の変化も同時に観察可能な新しい3次元可視化手法を提案する。

ここでは、1) 個々の弱学習器の構造・性能に関する情報、2) 弱学習器間関係性に関する情報、および3) 異なるデータセット間でのデータフローの差異に関する情報を

¹ お茶の水女子大学

² 日本電気株式会社

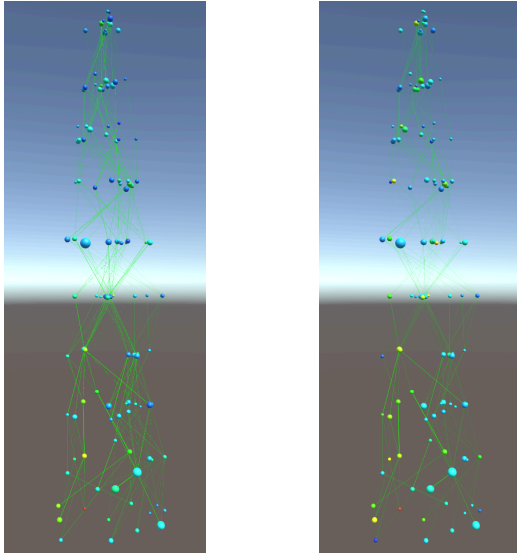


図 1 訓練データによる可視化 図 2 テストデータによる可視化

可視化することを試みる。1) は弱学習器である決定木の葉の数、葉ごとのサンプル数と残差、および葉同士の特徴量類似度とし、また 2) は隣接する決定木の葉に含まれるサンプルの重複度、3) は異なるデータセット間での 2) の差分とした。

3.2 インタラクション手法

次に、Boston Housing DataSet を用いて構築した GBDT モデル (LightGBM) を、提案手法により 3 次元空間上に可視化した結果を図 1 に示す。ここでは、ある住宅の価格を部屋数や近隣の犯罪率といった 13 属性の量質混合データに基づき、20 個の弱学習器からなる GBDT モデルにより推定した結果を示す。描画には Unity2021.3.8 を用いた。

3.2.1 モデルの視覚表現

図 1 において、シーン内の球体オブジェクトは、弱学習器の各々の葉を示しており、球体の直径および色は葉に含まれるサンプル数およびそれらの平均予測誤差をそれぞれ示す。ここでは、同一の x-y 平面に描画されている球体は一つの弱学習器に含まれる葉を表し、z 軸に沿って複数の弱学習器の葉の情報が描画されている。

球体オブジェクトを接続するリンクは弱学習器間の関係性を示しており、隣接する層に属する二つの葉に対し、それぞれに含まれるサンプルの重複度に応じてリンクが太くなるように設定した。

また、x-y 平面上における球体の座標は、1) 同じ弱学習器の葉同士の特徴量類似度、および 2) 隣接する弱学習器間でのサンプルの重複度に基づき決定した。

3.2.2 モデルの開発・運用支援手法

着目する葉をクリックすることで、葉内に含まれるの全てのデータ経路をハイライトする。リンクの色を赤く変化させることで経路を示す。また同時に、葉内のデータの詳細をグラフ表示する。これにより、図 1, 2 に示した全体

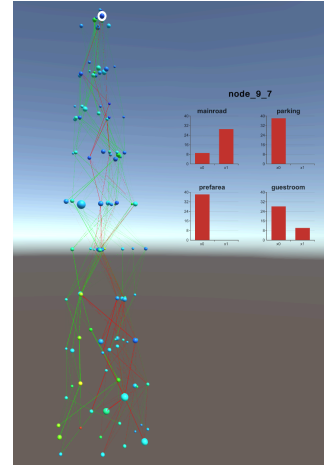


図 3 データ経路をハイライト、葉内のデータ詳細を出力

図からモデルの学習結果を俯瞰することで、特異な振る舞いをしている葉ノードの同定やデータセット間の差分についての理解を視覚的に支援するとともに、注目する弱学習器やデータサンプルについての詳細な情報を提示することで、モデルの微調整や予測誤差要因の分析作業に貢献する。

4. 今後の課題

本手法を他のデータセットにも適用させることで、よりモデルの解釈性を向上させるためのインタラクティブについて議論を進める。また、予測モデルの精度監視技術など、MLOps と呼ばれる運用中の課題を解決するための技術開発が進められている [6]。本研究においても、運用中の精度劣化原因を探索するような MLOps における技術開発にも適用させたい。

参考文献

- [1] Linhao Meng, et al.: *ModelWise: Interactive Model Comparison for Model Diagnosis, Improvement and Selection*(2022).
- [2] Zijie J. Wang, et al.: *TIMBERTREK: Exploring and Curating Sparse Decision Trees with Interactive Visualization* (2022).
- [3] Hikaru Nagasaka, Motoya Izuhara : *Interactive Visualization of Deep Learning Models in an Immersive Environment* (2021).
- [4] Dunn, Boris Kovalerchuk Andrew, Alex Worland, and Sridevi Wagle. *Interactive Decision Tree Creation and Enhancement with Complete Visualization for Explainable Modeling*.arXiv preprint arXiv:2305.18432 (2023).
- [5] Angelos Chatzimparmpas , et al.: *VisRuler: Visual analytics for extracting decision rules from bagged and boosted decision trees* (2023).
- [6] Sculley, D., Holt, G., Golovin, D., Davydov, E., Phillips, T., Ebner, D., Chaudhary, V., Young, M., Crespo, J.- F., and Dennison, D.: *Hidden Technical Debt in Machine Learning Systems*, in *Advances in Neural Information Processing Systems*, Vol. 28 (2015).