

# 数値計算ライブラリの自動チューニングにおける XAI適用の試み

An Adaptation of XAI to Auto-tuning for Numerical Calculation Library

青木将太<sup>1)</sup>, 片桐孝洋<sup>2)</sup>, 大島聡史<sup>3)</sup>, 永井亨<sup>4)</sup>, 星野哲也<sup>5)</sup>

Takahiro Katagiri, Shota Aoki, Satoshi Ohshima, Toru Nagai

1) 学(情報) 名古屋大学大学院 情報学研究科 (〒464-8601 愛知県名古屋市千種区不老町 情報学研究科,

E-mail: aoki@hpc.itc.nagoya-u.ac.jp)

2) 博(理) 名古屋大学情報基盤センター 教授 (〒464-8601 愛知県名古屋市千種区不老町 名古屋大学情報基盤センター, E-mail: katagiri@cc.nagoya-u.ac.jp)

3) 博(工) 九州大学情報基盤研究開発センター 准教授 (〒819-0395 福岡市西区元岡744,

E-mail: ohshima@cc.kyushu-u.ac.jp)

4) 博(理) 名古屋大学情報基盤センター 助教 (〒464-8601 愛知県名古屋市千種区不老町 名古屋大学情報基盤センター, E-mail: nagai@cc.nagoya-u.ac.jp)

4) 博(理) 名古屋大学情報基盤センター 准教授 (〒464-8601 愛知県名古屋市千種区不老町 名古屋大学情報基盤センター, E-mail: hoshino@cc.nagoya-u.ac.jp)

We explain an adaptation of Explainable AI (XAI) to auto-tuning problem for parameter tuning on a PICCG solver, which is one of preconditioned sparse iterative solvers. Result from SHAP, which is an XAI tool, indicates that predicted execution time for the solver from AI output can be well-explained in our experiment.

**Key Words** : Explainable AI, Auto-tuning, PICCG Solver

## 1. はじめに

人工知能 (AI) が出力する答えを検証が不十分なまま利用することで, 社会的な問題を引き起こすことが懸念されている. またAI出力結果については, WEB等の公知で入手可能な文書から生成されたものがあるため, 自動生成された文章等は, 真実が出力される保証はない.

以上のAIの現状では, AIの出力に対して, 人間がなんらかの検証を事前に行わなくてはならない. しかしAIモデルからの莫大な出力を, すべて人間が確認することは現実的ではない. つまるところ, このAI出力の検証工程の自動化や, 検証工数の削減手法が必須となる. そのため, AIモデルを構築する際のコストを減らす研究開発がされている.

一方, 我々は数値計算分野での性能チューニング工数の削減を目的として, ソフトウェア自動チューニング (Software Auto-tuning, AT) 技術[1]へのAI適用を行ってきた. ここでのAT技術は, 単にモデル上のパラメタ調整をする技術ではない. AT技術は, 高性能コード生成も行う性能チューニングの総合的技術であり, コンパイラ最適化やアルゴリズム選択などの上位の最適化対象も取り扱う. ここでは, 性能パラメタ調整の観点にのみ着目する.

本研究では, 疎行列反復解法上に現れる性能パラメタチューニング事例を例題として, AIを適用したAT機能を実現した場合における, AIが出力する予測時間に関する

説明性について検証を行う.

## 2. XAIツール

### (1) 概要

AI出力の説明性を高めることで「信頼されるAI」を実現するためのツール (Explainable AI (XAI)ツール) がいくつか提案されている. さらにこのXAIツールは, いくつかはオープンソース化がなされており, 容易に利用可能である. 以下に代表的なXAIツールを説明する.

### (2) LIME

幅広い学習済みモデルに適用できるXAIツールの1つとして, LIME (Local Interpretable Model-agnostic explanations)[2] が知られている.

LIMEはLocal surrogate model であり, ブラックボックスモデルの個々の予測を説明するために用いられる解釈可能なモデルである. 分類器が特定の予測を行った理由を, 人間が理解できるように提示する機能を持っている. それぞれの特徴がどの程度, 分類に貢献しているか調べることにより, 分類器の予測結果を説明することができる. 分類器の予測結果を用いるため, 任意の分類器に適用できるのが特徴である.

LIMEは個別の事例を説明するのに使われる「局所説明ツール」である.

## (3) SHAP

LIMEに加えて良く利用されるXAIツールとして、SHAP (SHapley Additive exPlanations) [3] が知られている。SHAPは、協力ゲーム理論のシャープレイ値 (Shapley Value) を機械学習に応用したものである。そのため、採用されている評価基準に、理論的な妥当性がある。シャープレイ値の計算は、厳密にすると計算量が高い。そのため、近似的にシャープレイ値を算出する手法が研究されており、SHAPで利用されている。

SHAPは全体的な傾向を説明するのに使われる、大域説明ツールである。

## 3. 不完全コレスキー分解前処理付CG法 (PICCG法)

## (1) 概要

不完全コレスキー分解 (Incomplete Cholesky Decomposition, IC) は、疎な対称行列  $A$  を係数とする連立一次方程式  $Ax = b$  ( $A \in \mathbb{R}^{n \times n}$ ,  $x, b \in \mathbb{R}^n$ ) のような連立1次方程式を解く反復法の1つである共役勾配法 (Conjugate Gradient (CG) Method) の前処理として活用される。IC分解前処理付きのCG法は、PICCG法と呼ばれる。

ICを用いて行列  $A$  を、以下の式(1)の分解をすることを考える。

$$A = U^t D U + R \quad (U, D \in \mathbb{R}^{n \times n}) \quad \dots (1)$$

ここで、 $U$ : 上三角行列、 $D$ : 対角行列、 $R$ :  $A$  とIC分解後の  $U^t D U$  との差の行列、とする。このとき、行列の要素値が0であった  $A$  の要素が、分解行列  $U$  において、非ゼロ要素となる場合がある。これを、**fill-in** と呼ぶ。

基本的なIC分解は、fill-inをすべて棄却する実装になっている。つまり、全て要素を0要素として扱う。このことで、分解後の行列の非ゼロ要素を少なく保ち、メモリ量と計算量を減らすことを狙う。しかし反面、 $A$  と  $U^t D U$  の差が大きく異なる場合は、うまく機能しないかもしれない。すなわち、前処理行列として機能しなくなることがある。

## (2) 閾値付きIC分解前処理とアルゴリズム

ここでは、閾値付きIC分解について説明する。

式(1)の行列分解過程で、分解行列  $U$  の非ゼロ要素について、**fill-inレベル**を持つようにする。Fill-inレベルとは、非零要素をどこまで許容するかレベルである。通常、行当たりや対角要素当たりを対象として、行列や扱う問題対象の特性を考慮し、このレベルを定義することが多い。

本研究では、東京大学情報基盤センターの河合による実装[4]を利用する。この実装では、**最大fill-inレベル** ( $m$ )、および**閾値** ( $t$ ) を設定する。

最大fill-inレベルでは、指定以下のレベルのfill-inに対して、閾値より小さいfill-in対象の要素値を0とみなして、閾値以上のfill-inを許容する。このことで、IC分解前処理と同程度以上の収束性と、より少ない非ゼロ要素数の行列に分解することで、メモリ量を低く抑えつつ高速化を狙うものである。つまり、最大fill-inレベルを増やせば、一般的に、反復回数は減る。しかし、行列が密に近づき演算量が増えるため、反復回数の削減と演算量の増加のトレードオフがある。つまりは、これらは性能チューニングパラメタになる。著者が扱っている問題では経験的に、最大fill-inレベル2が最適であるが、当然、扱う問題の性質に依存する。

以上のアルゴリズムの概略を、図-1に記載する。ここで、 $a_{i,j}$ :  $A$  の  $i, j$  要素、 $d_{i,i}$ :  $D$  の  $i, i$  要素、 $u_{i,j}$ :  $U$  の  $i, j$  要素、 $f_{i,j}$ :  $u_{i,j}$  の fill-in レベル、 $t$ : 0 とみなす閾値、 $m$ : 最大 fill-in レベルである。

$$\begin{aligned} d_{i,i} &= a_{i,i} - \sum_{k=1}^{i-1} u_{i,k} d_{i,k} u_{k,i} \\ f_{i,j} &= \begin{cases} 0, & a_{i,j} \neq 0 \\ f_{i,k} + f_{k,i} + 1, & \text{else} \end{cases} \\ u_{i,j} &= \begin{cases} d_{i,j}^{-1} (a_{i,j} - \sum_{k=1}^{i-1} u_{i,k} d_{i,k} u_{k,j}), & f_{i,j} \leq m \wedge |u_{i,j}| \geq t \\ 0, & \text{else} \end{cases} \end{aligned}$$

図-1 閾値付き IC 前処理の概要

## 4. 適用事例と分析

## (1) 計算機環境

AIの教師データを取得するため、名古屋大学情報基盤センター設置のスーパーコンピュータ「不老」Type I サブシステム[5]を利用した。また、機械学習は、「不老」Type II サブシステム[5]のGPUを利用した。

Pythonは ver. 3.6.13, Tensorflowは ver. 2.4.1, SHAPは ver. 0.39.0を利用している。

## (2) PICCG法の実行時間を予想するモデル

Tensorflowを使用し、以下の入力と出力を設定する。

- 入力：
  - 係数行列の特徴画像 (山田ら[6]の手法により作成)
  - 最大 fill-in レベル、閾値
- 出力：閾値付きICCG法の計算時間

畳み込み層、プーリング層: 2層、全結合層: 3層のCNNを作成した。モデルの概要を図-2に示す。

学習設定として、エポック数は200、バッチサイズは256、活性化関数はReLU、最適化はAdam法、損失関数は平均二乗誤差とした。

## (3) データセットの作成

差格子子によってメッシュ分割された三次元領域において定常熱伝導方程式を解くアプリケーション (P3D) [6]

を対象にし、有限体積法に基づき非構造格子型のデータとして考慮したデータを利用する[7]。ここで、熱伝導率 $\lambda$ の分布( $\lambda_2 \leq \lambda_1$ )の問題空間であり、熱伝導率が $\lambda_1$ の層の中に、伝導率 $\lambda_2$ の層を入れ込んだものである。この $\lambda_2$ の値が小さくなると、生成される係数行列の条件数が大きくなっていく。つまり悪性問題に近づく、反復解法では解きにくい問題になるので、問題の性質を $\lambda_2$ で制御できる特徴を持つ。

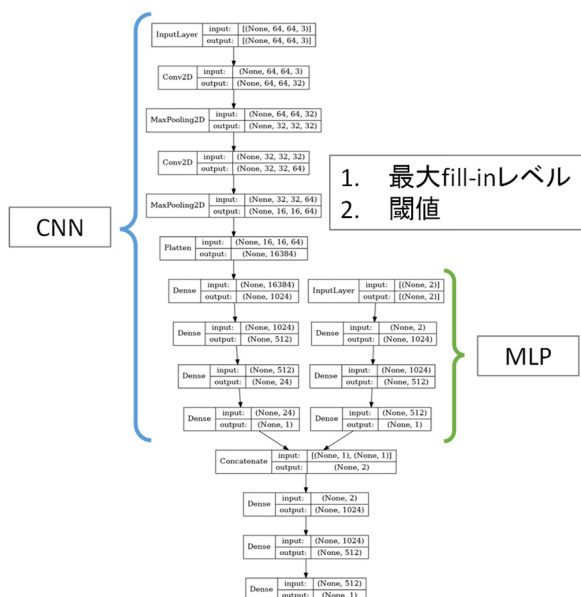


図-2 モデルの概要

教師データは、「不老」TypeIサブシステムで以下の条件で計算時間を計測して取得する。

- 係数行列
  - サイズ: 4096×4096, 32768×32768, 262144×262144 【3種類】
  - 各サイズ条件数を変えた (伝導率 $\lambda_2$ を変化) 【90種類】
- 最大fill-inレベル: 0, 1, 2 【3種類】
- 閾値: 0.001~0.02 (0.001間隔) 【199種類】

以上から教師データ総数は、各サイズで【41,073個】となる。また、テストデータ総数は、各サイズで【10,269個】となる。

各サイズ、に閾値付きICCG法の実行時間を予測する回帰モデルを作成する。

#### (4) SHAPによる説明結果

まず、図-3によるテストデータ (実測実行時間) と、AIモデルによる予測結果 (予測実行時間) の平均絶対誤差は0.000526であり、5%以下の誤差であったため、モデル生成は妥当である。かつ、本テストデータの範囲において、AIによるATが、誤差5%で実現できたことも意味している。

次に、この生成したAIモデルが、妥当な回答をしているか、SHAPの説明により検証する。図-3に、SHAPによる問題

サイズ32768×32768の説明結果を示す。

図-3の説明から、以下の解釈ができる。

- 閾値が一定以下になると、閾値が小さいほど実行時間が短くなる傾向がある
- Fill-inレベル2の場合、閾値が0.075以下において、実行時間が短くなっていく傾向がある

以上のSHAPにより説明された傾向は、対象アルゴリズムの知見を使ったものではないことに注意する。しかしこの傾向は、閾値付きIC前処理のICCG法のアルゴリズム上、および、今回のテストデータの実測実行時間の分布傾向から妥当な解釈であった。したがって、生成したAIモデルが妥当であることを示している。

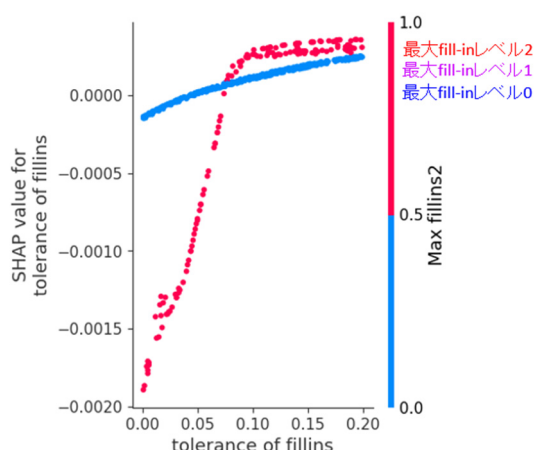


図-3 問題サイズ 32768×32768 の説明結果

#### 5. おわりに

本研究では、説明可能なAI (XAI) のツールとして、SHAPを利用して、疎行列反復解法に現れる性能パラメータチューニングにAIを活用する自動チューニング (AT) 機能を開発する場合の説明可能性について、事例をもとに結果を検証した。

検証結果から、SHAPを利用した本事例の説明については、妥当に説明されていることを確認した。また、本稿では説明していないが、SHAPにより妥当な解説がされていないことを発見し、原因を解析した。その結果、説明変数を定式化する際に、あるテクニックを適用したところ、解の精度と、SHAPによる妥当な説明が可能となる事例も発見した。本件については、当日報告する予定である。

今後の課題として、さらに多くの事例解析を行い、XAIツールが数値計算処理におけるAT機能にも有効であることを示していく必要がある。加えて、XAIツールからの説明に基づき、有効となる説明変数の追加や削減を自動で行い、AT機能の性能向上に寄与する仕組みの研究開発も必要である。

謝辞: 本研究は、科学技術研究費補助金、基盤研究 (S), 「(計算+データ+学習) 融合によるエクサスケール時代

の革新的シミュレーション手法」(課題番号:19H05662), および, 学際大規模情報基盤共同利用・共同研究拠点、および、革新的ハイパフォーマンス・コンピューティング・インフラ (課題番号: jh220022) の支援による。また, PICCGソルバーの知識提供について, 東京大学情報基盤センター河合直聡 特任助教に感謝の意を表する。

#### 参考文献

- [1] T. Katagiri, and D. Takahashi: Japanese Auto-tuning Research: Auto-tuning Languages and FFT, *Proc. of the IEEE*, 106 (11), pp.2056-2067, 2018.
- [2] M. T. Ribeiro, S. Singh, and C. Guestrin: Why should I trust you?: Explaining the predictions of any classifier, *Proc. of 22nd ACM SIGKDD*, pp.1135-1144, 2016.
- [3] S. M. Lundberg, and S. Lee: A unified approach to interpreting model predictions, *Proc. of the 31st International Conference on Neural Information Processing Systems (NIPS' 17)*, pp. 4768-4777, 2017.
- [4] 河合直聡, 中島研吾, 前処理付きクリロフ部分空間法への低/任意精度の適用, 京都大学数理解析研究所研究集会, 2022
- [5] 名古屋大学情報基盤センター スーパーコンピュータ「不老」ホームページ  
<https://icts.nagoya-u.ac.jp/ja/sc/overview.html#type2>  
[2023年3月21日閲覧]
- [6] K. Yamada, T. Katagiri, H. Takizawa, K. Minami, M. Yokokawa, T. Nagai, and M. Ogino: Preconditioner Auto-Tuning Using Deep Learning for Sparse Iterative Algorithms, *Proc. of 2018 Sixth International Symposium on Comput. and Networking Workshops (CANDARW)*, pp. 257-262, 2018.
- [7] 中島研吾, 大島聡史, 埴敏博, 星野哲也, 伊田明弘, ICCG法ソルバーのIntel Xeon Phi向け最適化, 情報処理学会研究報告 (2016-HPC-157-16) , 2016.