

グラフ属性の非類似度に着目した Random Forest からの解釈可能決定集合の抽出

Extraction of Interpretable Decision Sets from Random Forest focusing on Dissimilarity of Graph Attributes

松山 航太¹ 尾崎 知伸^{1*}
Kota Matsuyama¹ Tomonobu Ozaki¹

¹ 日本大学 文理学部

¹ College of Humanities and Sciences, Nihon University

Abstract: Interpretable decision sets (IDS) is a representative framework for building accurate predictive models with high interpretability directly from scratch. In this paper, by extending the core idea in IDS, we propose a two stage framework for deriving interpretable models on graph classification. In the first stage, a set of rules are derived through discovery of association rules and learning tree ensembles on subgraph patterns. The final models are constructed in the second stage based on our developed optimization functions for selecting a small number of simple and accurate rules that are dissimilar to each other.

1 はじめに

社会ネットワーク分析や、化学、生物科学分野の発展に伴い、近年、グラフ構造データ、すなわちノードとそれらを結ぶエッジから構成される構造データを対象とした知識発見技術の重要性が高まっている。グラフ構造データを対象とした分類モデル構築に関しては、部分構造を利用した決定木モデル [1, 2, 3] やカーネル法 [4, 5] 等に加え、深層学習技術に基づく手法 [6, 7, 8] も数多く提案されている。一般に、モデルの精度と解釈性・理解容易性にはトレードオフの傾向が見られ、決定木モデルに代表されるルールベースモデルは理解容易ではあるが精度が十分ではなく、またカーネル法や深層学習技術に基づく手法は、高精度ではあるがモデルの解釈が困難であるという弱点が指摘されている。本研究では、これらの問題を軽減し、グラフ構造データを対象とした高精度かつ解釈容易なモデルの構築を目的に、機械学習の解釈性研究 [9, 10] の代表的な成果の一つである解釈容易決定集合 (Interpretable Decision Sets; IDS) [11] をグラフ構造データへと拡張する。

機械学習の解釈性研究は、(1) 大域的な説明 (学習済みモデルを可読性の高い解釈可能なモデルで表現する)、(2) 局所的な説明 (特定の入力に対するモデルの予測根拠を提示する)、(3) 説明可能なモデルの設計 (最初から可読性の高い解釈可能なモデルを構築する) 等に大

別することが可能である [9]。IDS は、説明可能なモデルの設計に関する代表的な手法であり、LeGo アプローチ [12] に基づく二段階処理、すなわち関連ルール列挙 [13] に基づく局所ルール群の生成と、劣モジュラ最適化 [14] に基づくルール群選択を通じ、解釈容易なモデルを構築する。本研究では、IDS のアイデアを踏襲し、各段階における処理をグラフ構造データを対象としたものに置き換えることを考える。詳しくは後述するが、第一段階におけるルール生成に関しては、2通りの方法を検討する。第一の方法は、IDS のその素直な拡張であり、部分グラフパターン発見技術 [15] を用いてグラフ構造データをアイテム集合化することで、部分グラフをアイテムとする関連ルールを導出する。第二の方法は、分類モデル構築を通じたルール生成であり、部分グラフを属性とするランダムフォレストモデル [16] を構築し、各木の根から葉までのパスをルールとして抽出する。一方、第二段階におけるルール選択に関しては、ルールに含まれる部分グラフの大きさや類似度を考慮した新たな評価関数を導入する。これらを通じ、グラフ構造データからの解釈容易な分類ルール集合の導出を目指す。

本論文の構成は以下の通りである。2章で関連研究について言及する。3章でIDSの要点を整理する。4章で、グラフ構造データからのクラス関連ルールの導出とランダムフォレストモデルの構築に加え、グラフ構造を考慮したルール集合選択基準を導入し、IDSをグラフ構造データへと拡張する。5章で評価実験を行い、

*連絡先：日本大学文理学部情報科学科
〒156-8550 東京都世田谷区桜上水 3-25-40
E-mail: tozaki@chs.nihon-u.ac.jp

最後に 6 章でまとめと今後の課題を述べる。

2 関連研究

本研究では、部分グラフを属性とするランダムフォレストモデルからのルール抽出を行っている。これに関連し文献 [17] では、inTrees[18] と defragTrees[19] と呼ばれる技術を基に、部分グラフを属性とする分類木及び回帰木アンサンブルからのルール抽出が行われている。inTrees では不要な条件の削除やルールの集約を伴う集合被覆アルゴリズムを用い、また defragTrees では EM アルゴリズムやベイズ推定を用いた確率的生成モデルの最適化を通じ、それぞれ少数ルールからなる分類モデルが構築される。これに対し本研究では、グラフ属性の大きさや数、類似性に基づく最適化基準を採用している。

グラフ構造データを対象とした機械学習モデルの解釈性研究に関しては、局所的な説明、すなわち与えられた事例の分類結果に寄与する部分構造と属性を提示する GNNExplainer[20] や GraphLIME[21] と呼ばれる手法が提案されている。これに対し本研究は、事例に対する説明ではなく、モデルそのものの導出を目的としている。また文献 [22] では、本研究と同じく、IDS のグラフ構造データへの拡張が提案されている。本研究では、第一段階におけるルール列挙の多様化と、第二段階における最適化基準の見直しを行っており、[22] をより発展させたものであると言える。

3 解釈可能決定集合

本章では、文献 [11] に従い IDS の形式的な枠組みを示す。

属性ベクトル \mathbf{x} とクラスラベル y の組 (\mathbf{x}, y) を事例とするデータベース $D = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ から分類モデルを構築することを考える。〈属性, 演算子, 値〉の 3 項組をアイテムと呼び、その連言 (アイテム集合) を s と表記する。アイテム集合の全体集合を \mathcal{S} 、クラスラベルの集合を \mathcal{C} とし、アイテム集合 $s \in \mathcal{S}$ を前件、クラスラベル $c \in \mathcal{C}$ を後件とするルール $s \rightarrow c$ を (s, c) と表記する。またルールの集合 $\mathcal{R} = \{(s_1, c_1), \dots, (s_k, c_k)\} \subseteq \mathcal{S} \times \mathcal{C}$ を決定集合と呼ぶ。IDS の目的は、(後述する) 解釈性及び精度に関する評価関数 $f_1 \sim f_7$ の重み付き和を最大化するルール集合

$$R^* = \arg \max_{\mathcal{R} \subseteq \mathcal{S} \times \mathcal{C}} F(D, R) \text{ where } F(D, R) = \sum_{i=1}^7 \lambda_i f_i(\mathcal{R})$$

を求めることであり、以下に示す 2 段階処理の基づきこの目的を達成する。

$$\begin{aligned} f_1(D, \mathcal{R}) &= |\mathcal{S}| - |\mathcal{R}| \\ f_2(D, \mathcal{R}) &= |\mathcal{S}| \cdot \max\{|s'| \mid (s', c') \in \mathcal{S}\} - \sum_{(s, c) \in \mathcal{R}} |s| \\ f_3(D, \mathcal{R}) &= |D| \cdot |\mathcal{S}|^2 - \sum_{r_i, r_j \in \mathcal{R}, i < j, c_i = c_j} |\text{cov}(r_i) \cap \text{cov}(r_j)| \\ f_4(D, \mathcal{R}) &= |D| \cdot |\mathcal{S}|^2 - \sum_{r_i, r_j \in \mathcal{R}, i < j, c_i \neq c_j} |\text{cov}(r_i) \cap \text{cov}(r_j)| \\ f_5(D, \mathcal{R}) &= |\{c \mid (s, c) \in \mathcal{R}\}| \\ f_6(D, \mathcal{R}) &= |D| \cdot |\mathcal{S}| \\ &\quad - \sum_{(s, c) \in \mathcal{R}} |\{(\mathbf{x}, y) \in D \mid (\mathbf{x}, y) \in \text{cov}((s, c)), y \neq c\}| \\ f_7(D, \mathcal{R}) &= \left| \left\{ (\mathbf{x}, y) \in D \mid \begin{array}{l} (s, c) \in \mathcal{R}, y = c, \\ (\mathbf{x}, y) \in \text{cov}((s, c)) \end{array} \right\} \right| \end{aligned}$$

図 1: IDS における評価関数

1. アイテム集合列挙技術 [13] を用いて、集合 \mathcal{S} を導出する。
2. 各評価関数 f_i が劣モジュラ関数であることに着目し、平滑化局所探索 [14] を用いて R^* (の近似解) を獲得する。

以上示した通り、IDS では、解釈性および精度に関する複数の評価関数を同時に考慮することで、解釈性に優れかつ精度の高いルール群の抽出を試みる。また、各評価関数を劣モジュラ関数へと限定することで、効率的なルール集合選択を実現している。

IDS で採用されている評価関数を図 1 にまとめる。図中において、 $\text{cov}(r)$ は、ルール $r = (s, c) \in \mathcal{R}$ の支持集合を表す。関数 f_1 および f_2 は決定集合 \mathcal{R} の簡潔さ評価する関数であり、条件の少ない少数ルールで構成される決定集合に対し高い評価値を与える。また、関数 f_3 および f_4 は \mathcal{R} の有意義性を評価する関数であり、各ルールを支持する事例集合の重複が小さい場合、すなわち各ルールがそれぞれバラバラの事例を説明する場合に高い評価値を与える。関数 f_5 は、 \mathcal{R} によって予測可能なクラス種数であり、クラスに関する網羅性の観点から \mathcal{R} を評価する。関数 f_6 および f_7 は精度に関する関数であり、 f_6 は適合性、 f_7 は再現性の概念に相当する。

4 IDS のグラフ構造データへの拡張

本章では、グラフ構造データを対象とした IDS の拡張として、クラスラベルが付与されたグラフ集合を入力

とし、(1) ルール集合の抽出と (2) 最適化ルール集合の選択の二段階処理を用いたモデル構築手法を提案する。

4.1 相関ルール発見によるルール集合の抽出

\mathcal{L} をラベルの全体集合とし、ラベル付きグラフ $g = (V_g, E_g, l_g)$ を頂点集合 V_g と辺集合 $E_g \subseteq V_g \times V_g$ 、ラベル関数 $l_g : V_g \cup E_g \rightarrow \mathcal{L}$ の 3 項組で表現する。ラベル付きグラフ g とそのクラスラベル y の対 $t = \langle g, y \rangle$ を事例とするグラフデータベースを $G = \{(g_1, y_1), \dots, (g_n, y_n)\}$ と表記する。

グラフ g_i がグラフ g_j の部分グラフであることを $g_i \sqsubseteq g_j$ と表記する。この表記を用い、データベース G における部分グラフパターン p の支持度を

$$\text{sup}(p, G) = |\{(g, y) \in G \mid p \sqsubseteq g\}| / |G|$$

と定義する。支持度に関する閾値 σ ($0 < \sigma \leq 1$) に対し、条件 $\text{sup}(p, G) \geq \sigma \wedge \forall q \sqsubseteq p (\text{sup}(q, G) > \text{sup}(p, G))$ を満たす部分グラフを頻出フリー部分グラフ [23] と呼び、その全体集合を $\mathcal{F}_{G, \sigma}$ と表記する。

$\mathcal{F}_{G, \sigma}$ を用いてグラフデータベース G を頻出フリー部分グラフをアイテムとするトランザクションデータベース

$$D_{G, \sigma}^{AR} = \{(\{p \in \mathcal{F}_{G, \sigma} \mid p \sqsubseteq g_i\}, y_i) \mid \langle g_i, y_i \rangle \in G\}$$

へと変換する。また変換後の $D_{G, \sigma}^{AR}$ に対して、制約 $\forall g_i, g_j \in s [g_i \sqsubseteq g_j]$ 、すなわち包含関係にある部分グラフ同士を含まないという制約のもとで一般的な頻出パターン発見技術を適用し、部分グラフを構成要素とする相関ルール $r = (s, c) \in S_{G, \sigma}^{AR} \subseteq 2^{\mathcal{F}_{G, \sigma}} \times C$ の全体集合 $S_{G, \sigma}^{AR}$ を導出する。導出されたルール集合 $S_{G, \sigma}^{AR}$ を対象に、第二段階目の処理として評価関数を最適化するルール群の選択を行う。

4.2 ランダムフォレストモデルを用いたルール集合の抽出

相関ルール発見によるルール集合 $S_{G, \sigma}^{AR}$ の抽出に加え、グラフ属性を利用したランダムフォレストモデルからのルール集合抽出を考える。

グラフ $\langle g, y \rangle \in G$ と部分グラフパターン $p \in \mathcal{F}_{G, \sigma}$ に対し、 g が p を含むかを表す述語 $\text{att}(\langle g, y \rangle, p)$ を以下のように定義する。

$$\text{att}(\langle g, y \rangle, p) = \begin{cases} \langle p, =, 1 \rangle & (p \sqsubseteq g) \\ \langle p, =, 0 \rangle & (p \not\sqsubseteq g) \end{cases}$$

この述語を用い、グラフデータベース G を別のデータベース

$$D_{G, \sigma}^{RF} = \{(\{\text{att}(g_i, p) \mid p \in \mathcal{F}_{G, \sigma}\}, y_i) \mid \langle g_i, y_i \rangle \in G\}$$

$$f_2(D, \mathcal{R}) = |\mathcal{S}| - \frac{\sum_{(s, c) \in \mathcal{R}} |s|}{\max\{|s'| \mid (s', c') \in \mathcal{S}\}}$$

$$f_3(D, \mathcal{R}) = |\mathcal{S}| - \frac{2}{|D|(|\mathcal{S}| - 1)} \sum_{r_i, r_j \in \mathcal{R}, i < j, c_i = c_j} |\text{cov}(r_i) \cap \text{cov}(r_j)|$$

$$f_4(D, \mathcal{R}) = |\mathcal{S}| - \frac{2}{|D|(|\mathcal{S}| - 1)} \sum_{r_i, r_j \in \mathcal{R}, i < j, c_i \neq c_j} |\text{cov}(r_i) \cap \text{cov}(r_j)|$$

$$f_6(D, \mathcal{R}) = |\mathcal{S}| - \frac{\sum_{(s, c) \in \mathcal{R}} |\{(x, y) \in D \mid (x, y) \in \text{cov}((s, c)), y \neq c\}|}{|D|}$$

図 2: 変形後の評価関数

へと変換し、そこからランダムフォレストモデルを構築する。ランダムフォレストモデルにおける根から葉までの各パスをルールと捉え、その全体集合を $S_{G, \sigma}^{RF}$ と表記する。ここでルール $r = (s, c) \in S_{G, \sigma}^{RF}$ には、述語の否定、すなわち「ある部分グラフパターンを持たない」という否定条件が含まれる可能性がある。これにより、より表現力の高いルール群の獲得が期待できる。

4.3 グラフ構造に着目した評価関数の導入

IDS では、第二段階目の処理として、劣モジュラ最適化に基づくルール選択を採用している。これに関連し、本研究では、部分グラフの大きさや数、類似性に着目し、「互いに非類似でシンプルかつ少数の高精度なルール群」に対して高い評価値を与えることを目的とした 3 種の評価関数を提案する。また合わせて、既存 7 つの評価関数に対して値域をある程度揃えるための簡単な変形を導入する。これらの計 10 の評価関数に対し、その重み付き和

$$F_G(D, \mathcal{R}) = \sum_{i=1}^{10} \lambda_i f_i(D, \mathcal{R})$$

を最大化するルール集合

$$\mathcal{R}^* = \arg \max_{\mathcal{R} \subseteq \mathcal{S} \times C} F_G(D, \mathcal{R})$$

を最終的なモデルとして獲得することを提案する。

図 2 に、変形後の既存評価関数を示す。各式から明らかなように、これらの評価関数の値域は $0 \sim |\mathcal{S}|$ となる。また以下に、今回提案するグラフ構造を考慮した 3 つの評価関数を示す。

グラフサイズ： 関数 f_8 は、ルールに含まれるグラフサイズ（頂点数+辺数+ラベル種数）の総計に基づき値を算出するものであり、サイズが小さいルールに対して高い評価値を与える。

$$f_8(D, \mathcal{R}) = |\mathcal{S}| - \frac{1}{\max_{(s',c') \in \mathcal{S}} \text{size}(s')} \sum_{(s,c) \in \mathcal{R}} \text{size}(s)$$

where $\text{size}(s) = \sum_{g \in s} (w_v |V_g| + w_e |E_g| + w_l |\{l_g(o) \mid o \in V_g \cup E_g\}|)$,
and $w_v + w_e + w_l = 1$

ルール内非類似度： 関数 f_9 は、各ルール内のグラフ間非類似度に着目した評価関数であり、各ルール条件部が互いに非類似な部分グラフで構成される場合に高い評価を与える。また、2つのグラフ g_i, g_j 間の非類似度には、グラフ編集距離 $ged(g_i, g_j)$ を採用する。

$$f_9(D, \mathcal{R}) = |\mathcal{S}| - \frac{1}{\max_{(s',c') \in \mathcal{S}} \text{dist}(s')} \sum_{(s,c) \in \mathcal{R}} \text{SIM}(s)$$

where $\text{SIM}(s) = \max_{(s',c') \in \mathcal{S}} \text{dist}(s') - \text{dist}(s)$
and $\text{dist}(s) = \frac{1}{|s| C_2} \sum_{g_i, g_j \in s, i < j} ged(g_i, g_j)$

ルール間非類似度： 関数 f_{10} はルール集合内の各ルール間のグラフの非類似度に着目した評価関数であり、ルール集合が非類似なルール同士で構成されている場合に高い評価を与える。

$$f_{10}(D, \mathcal{R}) = |\mathcal{S}| - \frac{2 \sum_{(s_i, c_i), (s_j, c_j) \in \mathcal{R}, i < j} \text{SIM}(s_i, s_j)}{(|\mathcal{S}| - 1) \max_{(s'_i, c'_i), (s'_j, c'_j) \in \mathcal{S}} \text{dist}(s'_i, s'_j)}$$

where $\text{SIM}(s_i, s_j) = \max_{(s'_i, c'_i), (s'_j, c'_j) \in \mathcal{S}} \text{dist}(s'_i, s'_j) - \text{dist}(s_i, s_j)$
and $\text{dist}(s_i, s_j) = \frac{1}{|s_i| |s_j|} \sum_{g_i^x \in s_i} \sum_{g_j^y \in s_j} ged(g_i^x, g_j^y)$

5 評価実験

5.1 データセットと実験設定

提案手法の有効性を確認するため、株式会社 LIFULL が国立情報学研究所の協力により研究目的で提供している「LIFULL HOME'S データセット¹」を基に構築されたグラフデータセット [17] を用いた実験を行った。データセットは、東京 23 区内の賃貸マンション 480 件の各間取り図をグラフ化したものであり、各グラフ（間

表 1: データセットに関する統計量
クラス Y (240 件) クラス N (240 件)

	V	E	\lambda	V	E	\lambda
平均	25.49	26.63	9.84	24.31	25.57	9.34
分散	3.92	4.56	0.01	6.78	7.55	0.02
最小値	14.00	13.00	8.00	2.00	1.00	1.00
Q1	21.00	21.00	9.00	20.00	20.00	9.00
中央値	24.00	26.00	10.00	23.00	24.00	9.00
Q3	28.50	31.00	10.00	27.00	29.50	10.00
最大値	49.00	52.00	11.00	56.00	59.00	11.00

|V| は頂点数, |E| は辺数, |\lambda| はラベル種数を表す

表 2: 第一段階で抽出されたルール数

	クラス Y	クラス N
$S_{G,\sigma}^{AR}$	67,706	901
$S_{G,\sigma}^{RF}$	1,523	1588

取り図) に対し、推定賃料と実賃料との差に基づくクラス (Y または N) が付与されている。また、頂点ラベル種数は 11、辺ラベルは存在しない。表 1 にデータセットの概要を示す。

ルール集合の抽出は、以下の手順で行った。まず、頻出部分グラフマイナー $gSpan[24]$ と適切な後処理を用いて 7,616 件からなる頻出フリー部分グラフの集合 $\mathcal{F}_{G,\sigma}$ を導出した。次に、 $\mathcal{F}_{G,\sigma}$ を用いて構築される $D_{G,\sigma}^{AR}$ に対し、最小支持度 0.3、最小確信度 0.5 を用いてルール集合 $S_{G,\sigma}^{AR}$ を抽出した。また同様に、 $\mathcal{F}_{G,\sigma}$ を用いて構築される $D_{G,\sigma}^{RF}$ からランダムフォレストモデルを構築し、ルール集合 $S_{G,\sigma}^{RF}$ を抽出した。なお、各ルール集合に含まれるルール数は、表 2 に示す通りである。

実験では、各クラス 500 件ずつ、計 1000 件のルールをランダムに選択することでルール集合 \mathcal{S} を構築し、そこから最適化基準に従い \mathcal{R}^* を抽出する操作を 10 回繰り返した。また得られた \mathcal{R}^* に対し、図 3 に示す指標を用いて評価を行った。

#_Rules と Avg_Length は、それぞれ f_1 と f_2 に対応するモデルの簡潔性に関する評価指標であり、値が小さい方が優れていると判断する。一方、モデルの有意義性に関しては $Frac_Overlap$ および $Frac_Classes$ に基づき評価を行う。 $Frac_Overlap$ は f_3 と f_4 に対応する。値域は 0 ~ 1 であり、値が小さい方が優れていると判断する。 $Frac_Classes$ は f_5 に対応するクラスに関する網羅性を表す。今回は 2 値分類を対象としているので値域は {0.5, 1.0} となる。モデルの精度に関しては、 f_7 に対応する $Frac_Uncovered$ を用いて評価を行う。この基準は、説明されない事例の割合であり、値が小さい方が優れていると考える。また、グラフ構造に焦点を当てた評価指標として、それぞれ $f_8 \sim f_{10}$ に対応する Graph_Size、

¹<https://www.nii.ac.jp/dsc/idr/lifull/homes.html>

$\#_Rules(\mathcal{R}) = \mathcal{R} $
$Avg_Length(\mathcal{R}) = \frac{1}{ \mathcal{R} } \sum_{(s,c) \in \mathcal{R}} s $
$Frac_Overlap(R) = \frac{ D }{ s \binom{C}{2}} \sum_{r_i, r_j \in \mathcal{R}, i < j} cov(r_i) \cap cov(r_j) $
$Frac_Classes(\mathcal{R}) = \{c \mid (s, c) \in \mathcal{R}\} / C $
$Frac_Uncovered(R) = 1 - \bigcup_{r \in \mathcal{R}} cov(r) / D $
$Graph_Size(R) = f_8(D, R)$
$In_Graph_Dist(R) = f_9(D, R)$
$Bet_Graph_Dist(R) = f_{10}(D, R)$

図 3: \mathcal{R}^* に対する評価指標

In_Rule_Graph_Dist, Bet_Rule_Graph_Dist を用いる。

5.2 結果と考察

実験結果を表 3 に示す。なお表中の F_G^M は、 R^* の選択基準として評価関数

$$F_G^M(D, R) = \sum_{i \in \{1, \dots, 10\} \setminus M} \lambda_i f_i(D, R)$$

すなわち M に含まれる評価関数を取り除いた場合を採用したことを表す。また従って、 $F_G^{\{8,9,10\}}(D, R) = F(D, R)$ は IDS オリジナルの評価関数に相当する。

まず、提案した評価関数 F_G を用いた場合と IDS オリジナルの評価関数 $F_G^{8,9,10}$ を用いた場合との結果の違いを比較する。相関ルールを経由した手法 ($S_{G,\sigma}^{AR}$ からのルール選択) では、精度 (AUC) と簡潔性 (Ave.Length), 有意義性 (Frac.Overlap) の指標に関し、提案手法が既存手法より優れていることが確認できた。加えて、提案手法の方が劣っている指標はなく、これらの結果より、新たに導入した評価関数 $f_8 \sim f_{10}$ が有効に働いていると結論付けることができる。

一方、ランダムフォレストを経由した手法 ($S_{G,\sigma}^{RF}$ からのルール選択) では、ほとんどの指標において $F_G - F_G^{8,9,10}$ 間で有意な差が見られず、導入した評価関数 $f_8 \sim f_{10}$ が必ずしも有効に働いているわけではないことが分かる。また相関ルールを経由した手法と比べ、ルール長が長く、また被覆されない事例が多いことが分かる。これらの違いは、第一段階で抽出されるルールの数および質に起因するものと考えられる。相関ルールを用いた場合、少数の事例を説明する局所的なルールが多数導出される。一方、ランダムフォレストモデルを用いた場合、否定情報が含まれるなど表現能力は向上するが、樹状モデルのため、他のルールとの関係が考慮された大域的なルールが生成されることになる。この違いが、最終的な結果に大きな影響を与えていると考えられる。

次に、グラフ構造に関する評価基準 $f_8 \sim f_{10}$ について考察する。相関ルールを用いた手法では F_G^{10} と $F_G^{8,9,10}$ の結果が類似していること、および F_G^9, F_G^{10} と F_G の結果が類似していることがそれぞれ確認できる。 F_G^{10} は、ルール間グラフ非類似度に関する評価関数 f_{10} を排除した場合であり、簡潔性 (Ave.Length) 以外の F_G^{10} の結果が $F_G^{8,9,10}$ と類似するということから、ルール間グラフ非類似度が精度や有意義性の向上に強く寄与していることが示唆される。また、 F_G^9 は簡潔性 (Ave.Length) のみ $F_G^{8,9,10}$ に比べて結果が劣っていることから、ルール内非類似度は簡潔性の向上に寄与していることが示唆される。一方で、 F_G^8 の結果が F_G の結果と類似するということから、 f_8 (グラフサイズ) は大きな影響を持たないことが伺える。

6 まとめ

本研究では、グラフ構造データに対する高精度かつ解釈容易な分類モデルの構築を目的に、IDS の拡張を行った。第一段階におけるルール抽出に関しては、相関ルール発見を用いる手法とランダムフォレストモデルを用いる手法を提案した。また第二段階におけるルール選択に関しては、グラフ構造に着目し、互いに非類似でシンプルかつ少数の高精度なルール群に対して高い評価値を与える新たな評価関数を導入した。

今後の課題としては、ルール抽出手法の精緻化と更なる評価関数の開発があげられる。加えて、多様な実データを対象とした定量評価および人間による主観評価の実施が必要であると考えている。

謝辞： 本研究では、株式会社 LIFULL が国立情報学研究所の協力により研究目的で提供している LIFULL HOME'S データセットを利用した。

参考文献

- [1] B. Bringmann and A. Zimmermann : Tree² - Decision Trees for Tree Structured Data, *Proc. of the 9th European Conference on Principles and Practice of Knowledge Discovery in Databases*, pp.46–58 (2005)
- [2] P. C. Nguyen, K. Ohara, A. Mogi, H. Motoda, and T. Washio : Constructing Decision Trees for Graph-Structured Data by Chunkingless Graph-Based Induction, *Proc. of the 10th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining*, pp.390–399 (2006)
- [3] 尾崎 知伸, 渡沼 智己, 大川 剛直 : 多次元構造データからの分類知識の獲得, *人工知能学会論文誌*, Vol.22, No.2, pp.173–182 (2007)

表 3: 実験結果

	F_G	$F_G^{\{8\}}$	$F_G^{\{9\}}$	$F_G^{\{10\}}$	$F_G^{\{8,9,10\}}$	F_G	$F_G^{\{8\}}$	$F_G^{\{9\}}$	$F_G^{\{10\}}$	$F_G^{\{8,9,10\}}$
	AR					RF				
AUC	0.61	0.61	0.62	0.57	0.58	0.48	0.48	0.48	0.48	0.48
#_Rules	2.00	2.00	2.00	2.00	2.00	2.10	2.10	2.10	2.10	2.10
Ave_Length	2.65	2.60	2.85	2.70	2.80	2.97	2.97	3.00	2.97	3.07
Frac_Overlap	0.51	0.51	0.51	0.68	0.67	0.11	0.11	0.11	0.11	0.10
Frac_Classes	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
Frac_Uncovered	0.02	0.02	0.02	0.02	0.02	0.71	0.71	0.71	0.71	0.70
Graph_Size	499.24	499.23	499.26	499.45	499.48	497.98	497.98	497.95	497.98	497.87
In_Graph_Dist	499.83	499.84	499.39	499.79	499.32	499.50	499.50	499.48	499.50	499.48
Bet_Graph_Dist	499.996	499.996	499.996	499.996	499.996	499.996	499.996	499.996	499.996	499.996

- [4] S. V. N. Vishwanathan, N. N. Schraudolph, R. Kondor, and K. M. Borgwardt : Graph Kernels, *Journal of Machine Learning Research*, vol.11 pp.1201–1242 (2010)
- [5] N. M. Kriege, F. D. Johansson, and C. Morris : A survey on graph kernels, *Applied Network Science*, Vol.5, Article Number: 6 (2020)
- [6] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu : A Comprehensive Survey on Graph Neural Networks, arXiv:1901.00596v4 (2019)
- [7] F. Errica, M. Podda, D. Bacciu, and A. Micheli : A Fair Comparison of Graph Neural Networks for Graph Classification, *Proc. of the 8th International Conference on Learning Representations* (2020)
- [8] Z. Ying, J. You, C. Morris, X. Ren, W. Hamilton, and J. Leskovec : Hierarchical Graph Representation Learning with Differentiable Pooling, *Advances in Neural Information Processing Systems 31*, pp.4800–4810 (2018)
- [9] 原 聡 : 私のブックマーク「機械学習における解釈性」, *人工知能*, Vol.33, No.3, pp.366–369 (2018)
- [10] A. B. Arrieta et al. : Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI, *Information Fusion*, Vol.58, pp.82–115 (2020)
- [11] H. Lakkaraju, S. H. Bach, and J. Leskovec : Interpretable Decision Sets: A Joint Framework for Description and Prediction, *Proc. of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp.1675–1684 (2016)
- [12] A. Knobbe, B. Crémilleux, J. Fürnkranz, and M. Scholz : From Local Patterns to Global Models: The LeGo Approach to Data Mining, *Proc. of the ECML/PKDD-08 Workshop on From Local Patterns to Global Models (LeGo-08)*, pp 1–16 (2008)
- [13] R. Agrawal and R. Srikant : Fast Algorithms for Mining Association Rules in Large Databases, *Proc. of the 20th International Conference on Very Large Data Bases*, pp.487–499 (1994)
- [14] U. Feige, V. S. Mirrokni, and J. Vondrák : Maximizing non-monotone submodular functions, *SIAM Journal on Computing*, Vol.40, No.4, pp.1133–1153 (2011)
- [15] H. Cheng, X. Yan, and J. Han : Mining Graph Patterns, In C. Aggarwal, and H. Wang (eds) *Managing and Mining Graph Data*, pp.365–392, Springer (2010)
- [16] L. Breiman : Random Forests, *Machine Learning*, Vol.45, Issue 1, pp.5-32 (2001)
- [17] T. Ozaki : Extraction of Characteristic Subgraph Patterns with Support Threshold from Databases of Floor Plans, *Proc. of the 2019 Seventh International Symposium on Computing and Networking*, pp.197–203 (2019)
- [18] H. Deng : Interpreting Tree Ensembles with inTrees, *International Journal of Data Science and Analytics*, Issue 4/2019 (2018)
- [19] S. Hara and K. Hayashi : Making Tree Ensembles Interpretable: A Bayesian Model Selection Approach, *Proc. of the 21st International Conference on Artificial Intelligence and Statistics*, pp.77-85 (2018)
- [20] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec : GNNExplainer: Generating explanations for graph neural networks, *Advances in Neural Information Processing Systems 32 (NIPS 2019)*, pp.9244–9255 (2019)
- [21] Q. Huang, M. Yamada, Y. Tian, D. Singh, D. Yin, and Y. Chang : GraphLIME: Local Interpretable Model Explanations for Graph Neural Networks, arXiv:2001.06216 (2020)
- [22] 松山 航太, 尾崎 知伸 : グラフ構造データを対象とした解釈可能決定集合の拡張, 第 24 回 人工知能学会 インタラクティブ情報アクセスと可視化マイニング研究会, SIG-AM-24-05, pp.25–29 (2020)
- [23] Z. Zeng, J. Wang, J. Zhang, and L. Zhou : FOGGER : An Algorithm for Graph Generator Discovery, *Proc. of 12th International Conference on Extending Database Technology*, pp.517–528 (2009)
- [24] X. Yan and J. Han : gSpan : Graph-based Substructure Pattern Mining, *Proc. of the 2002 IEEE International Conference on Data Mining*, pp.721-724 (2002)