

研究課題 臨床や地域の課題分析への応用を目指す解釈性の高い  
高精度なデータマイニング手法の確立



八戸工業大学・工学部システム情報工学科

研究者名 島内 宏和

データから規則を見出す機械学習の技術は著しく発展し、その手法の一つである深層学習を採用した人工知能 (Artificial Intelligence, AI) が、医療や法律など様々な分野において目覚ましい成果を挙げている。精神医学分野や地域における課題に取り組む際には稀な症例や事象を扱うことがあるが、そのような場合には必然的に限られたデータしか利用できない。他方、健常者のデータはより多く得られることが多いため、データセットはインバランスなものとなる。稀な症例・事象は、データセットの中では外れ値と見なせることがあり、そのような場合には外れ値検出と呼ばれる手法が有効になりうる。また、近年 AI をブラックボックスとして扱うだけでなく、その予測の根拠を与えることを目指した「説明可能な人工知能 (Explainable Artificial Intelligence, XAI)」の研究が活発に進められている。

このような背景の下、本研究では高精度かつ特徴量の重要度の評価が可能な教師なし外れ値検出アルゴリズムを、教師なし表現学習の手法と深層学習の手法の一種である敵対的生成ネットワーク (Generative Adversarial Network, GAN) により構築した。アルゴリズムは、複数の教師なし異常検知アルゴリズムによる予測スコア (Transformed Outlier Scores, TOS) を特徴と見なし元のデータに追加することで特徴空間を拡張し、GAN により拡張された外れ値のデータを拡張した上で、アンサンブル法の一種である勾配ブースティングによる分類を行うという構造をとっており、Zhao と Hryniewicki による Extreme Gradient Boosting Outlier Detection を拡張したものとなっている。構築したアルゴリズムは、Zhao と Hryniewicki による先行研究と比較して、複数の外れ値検出のためのベンチマークデータセット上でより高い性能を示した (ROC と Precision@n で評価、表を参照)。さらに、構築したアルゴリズムは、教師なし学習により TOS を得る際にその特徴量重要度を Permutation Importance により獲得しておき、最終的な出力を行う勾配ブースティングの特徴量重要度と統合することで、どの特徴が分類に寄与したか分析を行うことも可能である。今後は、TOS の生成に利用する教師なし学習や GAN のネットワークに修正を加えることで、さらなる外れ値検出性能向上を目指す。

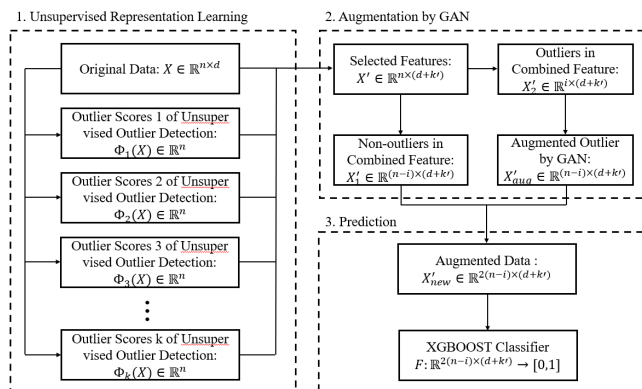


図. 提案した外れ値検出アルゴリズム

Dataset		Scores of Classification	
		TOS+XGB	Presented Algorithm
Arrhythmia	ROC	0.8783 (0.0366)	0.8875 (0.0296)
	PREC@N	0.6084 (0.0715)	0.6415 (0.0537)
Cardio	ROC	0.9973 (0.0017)	0.9979 (0.0009)
	PREC@N	0.9331 (0.0211)	0.9329 (0.0227)
Letter	ROC	0.9675 (0.0149)	0.9713 (0.0134)
	PREC@N	0.7259 (0.0582)	0.7490 (0.0536)
Mammography	ROC	0.9315 (0.0143)	0.9396 (0.0167)
	PREC@N	0.6667 (0.0312)	0.6703 (0.0290)
Speech	ROC	0.8965 (0.0384)	0.9175 (0.0421)
	PREC@N	0.3646 (0.0733)	0.3694 (0.0854)

表. 提案アルゴリズムの性能  
(30回の独立試行における平均と標準偏差)