IMLI: An Incremental Framework for MaxSAT-Based Learning of Interpretable Classification Rules

Bishwamittra Ghosh and Kuldeep S. Meel School of Computing, National University of Singapore

Interpretable Machine Learning

- The wide adoption of machine learning in the critical domains has propelled the need for interpretable techniques
- Interpretable machine learning model provides end users the reasoning behind decision making
- We propose an incremental approach to MaxSAT based interpretable rule learning framework

Experimental Results

Dataset	Size	Features	NN	SVC	RIPPER	MLIC	IMLI
PIMA	768	134	77.92 % 0.32 s	75.32 % 0.37 s	75.32 % 2.58 s	75.97 % Timeout	73.38 % 0.74 s
Credit- default	30000	334	79.61 % 872.97 s	80.69 % 847.93 s	80.97 % 20.37 s	80.72 % Timeout	79.41 % 32.58 s
Twitter	49999	1050	Timeout	Timeout	95.56 % 98.21 s	94.78 % Timeout	94.69 % 59.67 s



Existing Approach

- Reduces the learning problem as a MaxSAT query
- Generates interpretable rules expressed in CNF
- To generate a k clause CNF rule for a dataset of n samples over m boolean features, the number of clause of MaxSAT query is O(n * m * k)

Proposed Approach

- We attribute large formula size of the MaxSAT query for the poor scalability of the existing approach
- We propose a partition-based incremental learning framework

- Each cell in last 5 columns: test accuracy (%) and training time (s)
- IMLI exhibits better training time by costing a little bit of accuracy



- CNF(1) denotes the result for CNF rule with 1 clause
- Rule size decreases as the number of partitions (p) increases

Dataset	RIPPER	MLIC	IMLI
Parkinsons	2.6	2	8
Ionosphere	9.6	13	5
WDBC	7.6	14.5	2
Blood	1	3	3.5
Adult	107.55	44.5	28
PIMA	8.25	16	3.5
Tom's HW	30.33	2	2.5
Twitter	21.6	20.5	6
Credit	14.25	6	3



- Divide the training data into a fixed number of partitions
- The MaxSAT query constructed for partition i is based on the training data for partition i and the rule learned until partition i 1

Key Contribution

• IMLI makes p queries to MaxSAT solvers with each query of the size $O\left(\frac{n}{p} * m * k\right)$

- Each cell denotes the average length of the generated rule
- IMLI generates shorter rules compared to other models

Interpretable Rules

- Credit-default Dataset
 - A client will default if
 - Education type = other OR
 - repayment status in September: payment delay > 1 month OR repayment status in August: payment delay > 2 months OR repayment status in June: payment delay > 2 months

Pima Indians Diabetes Dataset

A person is tested positive for diabetes if





Paper

Source Code

https://bishwamittra.github.io https://www.comp.nus.edu.sg/~meel/ Plasma glucose concentration > 125 AND Triceps skin fold thickness ≤ 35 mm AND Diabetes pedigree function > 0.259 AND Age > 25 years

Conclusion

- IMLI achieves up to three orders of magnitude improvement in training time
- The generated rules appear to be reasonable, intuitive, and more interpretable