

A Comparative Study of Feature Selection in Tsetlin Machines

Vojtech Halenka¹, Ole-Christoffer Granmo¹, Lei Jiao¹, and Per-Arne Andersen¹

Department of ICT, University of Agder, Grimstad, Norway
 {vojtech.halenka,ole.granmo,lei.jiao,per.andersen}@uia.no

AUC area under curve
FS Feature Selection
NN Neural Network
DNN Deep Neural Network
TA Tsetlin automaton
TAT Tsetlin automaton team
TM Tsetlin machine
ROAR Remove and Retrain
ROAD Remove and Debias

Abstract. Feature Selection (**FS**) is crucial for improving model interpretability, reducing complexity, and sometimes for enhancing accuracy. The recently introduced Tsetlin machine (**TM**) offers interpretable clause-based learning, but lacks established tools for estimating feature importance. In this paper, we adapt and evaluate a range of **FS** techniques for **TMs**, including classical filter and embedded methods as well as post-hoc explanation methods originally developed for neural networks (e.g., SHAP and LIME) and a novel family of embedded scorers derived from **TM** clause weights and Tsetlin automaton (**TA**) states. We benchmark all methods across 12 datasets, using evaluation protocols, like Remove and Retrain (**ROAR**) strategy and Remove and Debias (**ROAD**), to assess causal impact. Our results show that **TM**-internal scorers not only perform competitively but also exploit the interpretability of clauses to reveal interacting feature patterns. Simpler **TM**-specific scorers achieve similar accuracy retention at a fraction of the computational cost. This study establishes the first comprehensive baseline for **FS** in **TM** and paves the way for developing specialized **TM**-specific interpretability techniques.

1 Introduction

Machine learning models often benefit from **FS** to identify the most informative input features, which can improve generalization and interpretability. In Deep Neural Network (**DNN**), a rich array of methods exists to assess feature importance (e.g. SHAP, LIME) [1, 2]. However, the **TM** — a logic-based learning algorithm using Tsetlin automaton team (**TAT**) — has not yet seen a comparable

development of such interpretability-assisting tools. The **TM** achieves competitive accuracy on various tasks while producing human-readable rules [3], but questions remain on how to best discern which input features drive its decisions.

FS is particularly relevant for **TMs** to examine large learned rule sets and focus interpretation on key factors. While classical **FS** methods and Neural Network (**NN**)-inspired explainers could potentially be applied to **TMs**, their efficacy in this new context is largely unexplored. Prior work has hinted that **TMs** may handle high-dimensional inputs without explicit **FS** [4, 5], yet a systematic comparison of **FS** methods for **TMs** is missing.

In this paper, we bridge that gap by conducting a comparative study of **FS** methods applied to **TM** across a variety of benchmark datasets. We evaluate methods from all three major **FS** categories—filter, embedded, and wrapper—adapting methods from both traditional statistics and modern **NN** explainability. Our contributions are as follows:

- We adapt and apply four post-hoc explainers to Tsetlin Machines’ binary inputs.
- We propose a family of **TM**-internal scorers derived from clause weights and **TAs** states.
- We benchmark 23 methods on 12 UCI datasets using four evaluation protocols, tracking area under curve (**AUC**) and ranking time.
- We compare speed–quality trade-offs, ranking correlations, and outline **TM**-specific **FS** strategies.
- Finally we select a few representatives from clusters of the **FS** and plot them for explainability comparison.

2 Related Work

FS methods are generally categorised into:

- **Filter methods** evaluate feature relevance using data characteristics independent of any specific learning algorithm (e.g., Mutual Information [6], Chi-square [7], Variance thresholds). They are computationally efficient but ignore model interactions.
- **Wrapper methods** search for an optimal subset by repeatedly evaluating a predictive model. Modern wrappers include permutation importance [8] and other post-hoc probes, capturing interactions at high compute cost.
- **Embedded methods** perform **FS** as part of model training (e.g., tree-based split importance [9]). They balance filter speed with wrapper specificity but are model-specific.

Gradient-based and other wrapper explainers have been widely adopted in **NNs**, driven by the challenge of interpreting these black-box models:

- **LIME** fits local linear surrogates [2].
- **SHAP** uses Shapley values for global attributions [1].
- **Integrated Gradients (IG)** accumulates gradients along the inputs [10].

- **SmoothGrad** and its variants average noisy-gradient attributions [11].

The **TM** learns human-readable propositional clauses via **TAT** [3]. Prior work indicates **TMs** scale well with high-dimensional inputs (e.g., text n-grams [4], hypervector encodings [5]), but no systematic study of **FS** or feature importance has been conducted for **TMs**. This work fills that gap.

3 Overview of Tsetlin Machines

4 A Brief Overview of Tsetlin Machines

In short, the **TM** [3] is a logic-based pattern recognizer that learns human-readable propositional rules from binary inputs. It is built upon two key components: *clauses*, which are conjunctions of input *literals*, and **TAs**, which decide whether to include or exclude each literal in each clause. Learning proceeds via stochastic reinforcement on each automaton, using two types of feedback (detailed below). Figure 1 illustrates the overall architecture.

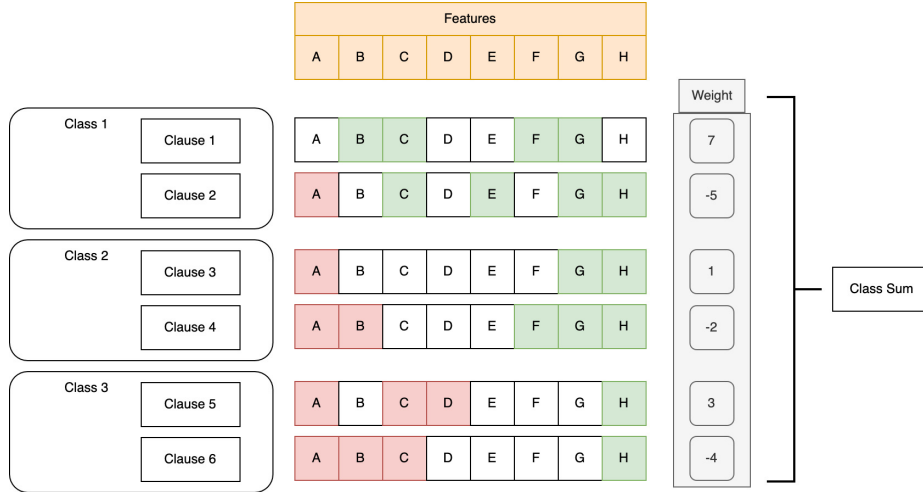


Fig. 1. Tsetlin Machine architecture: features (columns), clauses (rows), and class assignments. Green literals are included, red literals are excluded, white are irrelevant.

Given a binary input

$$\mathbf{x} = (x_1, \dots, x_d) \in \{0, 1\}^d,$$

each feature x_i is a literal $l_i \equiv x_i$ and has a negation $\neg l_i \equiv (1 - x_i)$. A *clause*

$$C_\ell(\mathbf{x}) = \bigwedge_{i \in S_\ell} (l_i)^{\delta_i} \in \{0, 1\}$$

is a conjunction over some subset $S_\ell \subseteq \{1, \dots, d\}$ of literals. For example,

$$C_\ell(\mathbf{x}) = l_1 \wedge \neg l_3 \wedge l_7 \iff (x_1 = 1) \wedge (x_3 = 0) \wedge (x_7 = 1).$$

Each clause $\ell = 1, \dots, M$ is assigned to a class $c_\ell \in \{1, \dots, C\}$ and carries a polarity $p_\ell \in \{+1, -1\}$, so that when the clause *fires* ($C_\ell(\mathbf{x}) = 1$), it casts a vote with a weight w_ℓ .

$$w_\ell C_\ell(\mathbf{x}), \quad \text{where } w_\ell \in R.$$

The *class sum* for class c aggregates these votes:

$$s_c(\mathbf{x}) = \sum_{\ell=1}^M p_\ell C_\ell(\mathbf{x}) \mathbf{1}\{c_\ell = c\}.$$

The final prediction is $\hat{c} = \arg \max_{c=1, \dots, C} s_c(\mathbf{x})$.

Each literal–clause pair is governed by a Tsetlin Automaton, a two-action finite state machine that learns via stochastic Type I and Type II feedback [12]:

- **Type I feedback** rewards clauses that correctly fire on positive examples, reinforcing inclusion of literals that keep the clause true.
- **Type II feedback** penalizes clauses that erroneously fire on negative examples, reinforcing inclusion of literals that make the clause false.

Over many epochs, each automaton converges to either *include* or *exclude* its literal, thereby shaping each clause to capture relevant patterns.

Key hyperparameters are:

- the total number of clauses M (often split equally per class and polarity),
- the *threshold* T , which controls the rate of feedback,
- the *granularity* s , which tunes how narrowly clauses match input patterns.

Despite their logical form, TMs scale to high-dimensional data (e.g., text n -grams, thermometer-encoded continuous features) by distributing learning across many simple automata.

Because each clause is a human-readable conjunction of literals, TMs offer innate rule-based explanations. However, large datasets typically yield hundreds or thousands of clauses, making manual inspection intractable. **FS**—identifying which original features appear most frequently or carry the greatest weight across clauses—thus provides a natural way to distill and visualize the key signals driving TM decisions.

5 Feature Selection Methods

Table 1 summarizes all **FS** methods evaluated: classical filters, post-hoc/wrapper explainers, and TM-internal embedded scorers.

We directly apply SHAP, LIME, IG, SmoothGrad, and Permutation Importance to the trained TM. We further define indicator functions

$$I_{\ell, f} = \mathbf{1}\{\text{clause } \ell \text{ contains the literal } l_f\},$$

Table 1. Overview of adapted and proposed feature-scoring methods, their origins, and key references.

Method	Description / References
Filter Methods	
MutualInfo	Measures mutual dependence between feature f and label. Based on Battiti (1994) [6] and Brown et al. [13].
Chi2	χ^2 test of independence between f and label. Based on Forman (2003) [7].
Variance	Removes features with low variance (baseline filter).
Random	Assigns random scores (sanity check).
Embedded (TM-Internal) Methods	
Relevance	Class-weighted TA engagement frequency: $\sum_{\ell} \alpha_c \mathbf{1}\{\ell \text{ selects } f\}$. Granmo (2020) [3].
TM-Weight	$s_f = \max_{c,\ell} w_{\ell} $ for those literals ℓ involving feature f . Proposed here.
CW-Sum	$s_f = \sum_{c=1}^C \alpha_c W_{c,f} $, where α_c is class-weight. Proposed here.
Support-CW-Sum	As CW-Sum but $\epsilon_c = 1 - \alpha_c$ replaces α_c . Proposed here.
CW-Feat	Let $\beta_{c,f} = W_{c,f} / \sum_{c'} W_{c',f} $. Then $s_f = \sum_{c=1}^C \beta_{c,f} W_{c,f} $. Proposed here.
Margin	Sort $\{ W_{c,f} \}_c$ as $a_{(1)} \geq a_{(2)} \geq \dots$; set $s_f = a_{(1)} - a_{(2)}$. Proposed here.
Entropy	Inverted Shannon entropy of $\{ W_{c,f} \}_c$. Based on Shannon (1948).
Gini	Gini index on $\{ W_{c,f} \}_c$. Based on Gini (1912).
Stability	Ratio of mean to std. of clause-weight history across epochs. Nogueira & Brown (2017) [14].
TaylorCrit	First-order change in true-class sum when flipping a literal. Inspired by Montavon et al. (2018) [15].
VarDropout	Sensitivity under random binary feature masks. Srivastava et al. (2014) [16].
AblationImpact	Impact of permanently ablating each feature. Zeiler & Fergus (2014) [17].
SmoothStabil	Sensitivity to small random input perturbations. Smilkov et al. (2017) [11].
Wrapper Methods	
Dropout	Mask-one-feature sensitivity (leave-one-out). Koh & Liang (2017) [18].
PermImportance	Permutation importance of each f . Fisher et al. (2018) [8].
Attribution / Ensemble NN-Inspired Methods	
SHAP	Shapley-value based contributions ϕ_f . Lundberg & Lee (2017) [1].
LIME	Local surrogate linear approximation. Ribeiro et al. (2016) [2].
IG	Integrated Gradients attribution. Sundararajan et al. (2017) [10].
SmoothGradSq	Squared-average of gradients over noisy inputs. Smilkov et al. (2017) [11].
VarGrad	Variance of gradients over noisy inputs. Smilkov et al. (2017) [11].

Then for each class c and feature f we accumulate

$$W_{c,f}^+ = \sum_{\ell=1}^M w_{\ell} I_{\ell,f}^+ \mathbf{1}\{c_{\ell} = c\}, \quad W_{c,f}^- = \sum_{\ell=1}^M w_{\ell} I_{\ell,f}^- \mathbf{1}\{c_{\ell} = c\}.$$

We set

$$\text{netW}_{c,f} = W_{c,f}^+ - W_{c,f}^-, \quad \text{absW}_{c,f} = |\text{netW}_{c,f}|.$$

For the PosNeg variants we define

$$\text{sumW}_{c,f} = W_{c,f}^+ + W_{c,f}^-, \quad \text{absSumW}_{c,f} = |\text{sumW}_{c,f}|.$$

6 Experimental Setup

Table 2 summarizes dataset properties, and baseline TM accuracy/F1. Continuous features are binned into 10 thermometer levels. Splits were set to 60/20/20 if no standard split is provided.

Table 2. Summary of datasets, TM parameters, and performance with all features at 30 epochs, 500 clauses, 10 thermometer bins per feature; parameters s and T found with 100 Optuna trials (s [0.9 to 20.0], T {50,200,300,500,800})

Dataset	#Sam.	#Feat.	#Cls	Source	s	T	Acc (%)	F ₁ (%)
Hierar. Bool.	500	20	2	generated	3.00	600	68.00	64.84
Parity	500	20	2	generated	3.00	600	44.00	42.53
Feature Interact.	500	20	2	generated	3.00	600	51.00	50.76
Balance Scale	625	4	3	UCI	3.00	600	88.00	61.12
Banknote	1 372	4	2	UCI	3.00	600	97.09	97.07
Breast Cancer	569	30	2	sklearn	18.33	500	94.74	94.35
Digits	1 797	64	10	sklearn	6.92	50	96.67	96.64
Ecoli	336	7	8	UCI	3.00	600	85.29	70.77
Glass	214	9	6	UCI	13.10	50	62.79	43.00
Heart Disease	303	13	2	OpenML	3.32	300	81.48	81.38
Ionosphere	351	34	2	OpenML	19.82	200	91.55	90.90
Iris	150	4	3	sklearn	14.80	300	90.00	90.15
Pima Diabetes	768	8	2	OpenML	8.63	50	72.73	67.45
Sonar	208	60	2	OpenML	3.30	50	88.10	87.92
Spambase	4 601	57	2	UCI	3.00	600	89.90	89.31
Steel Plates Faults	1 941	27	7	OpenML	8.34	50	78.92	74.16
Transfusion	748	4	2	UCI	3.00	600	77.33	56.35
Vehicle (Statlog)	846	18	4	OpenML	1.17	50	75.88	75.13
Wine	178	13	3	OpenML	9.70	800	94.44	94.53

TM models were implemented in Python 3.12.4 on Mac M1 Pro 16 GB RAM. We fix 500 clauses, 10 bins per feature, and tune (s, T) via Optuna (100

trials). Epochs=30 for score history. See Table 2 for found hyperparameters. Each Accuracy measure taken is an average over 10 trials.

We assess FS rankings by these evaluation protocols:

- **Insertion/Deletion**: mask/unmask top- k features
- **ROAR**: remove top- k features
- **ROAD**: replace top- k with marginal samples

ROAD should be the most advanced evaluation protocol, however, others are included for thorough comparison. Each protocol creates several subsets of features given by the chosen **FS** technique. These subsets are fed into 10 new **TMs**, with the same parameters as the original where we took the ranking from. Their accuracy is then averaged and plotted against the number of features K .

Figure 2 and 3 illustrates an example of these curves. Comparison of those two figures shows how the spread of assigned importance by a different FS correlates with the natural intuition. As in, Digits is has many redundant features, and medical datasets tend to have a single telling pattern, that once is found, defines the label. Moreover, Digits contains continuous features, and therefore thermometer encoding was used, this increased the redundancy.

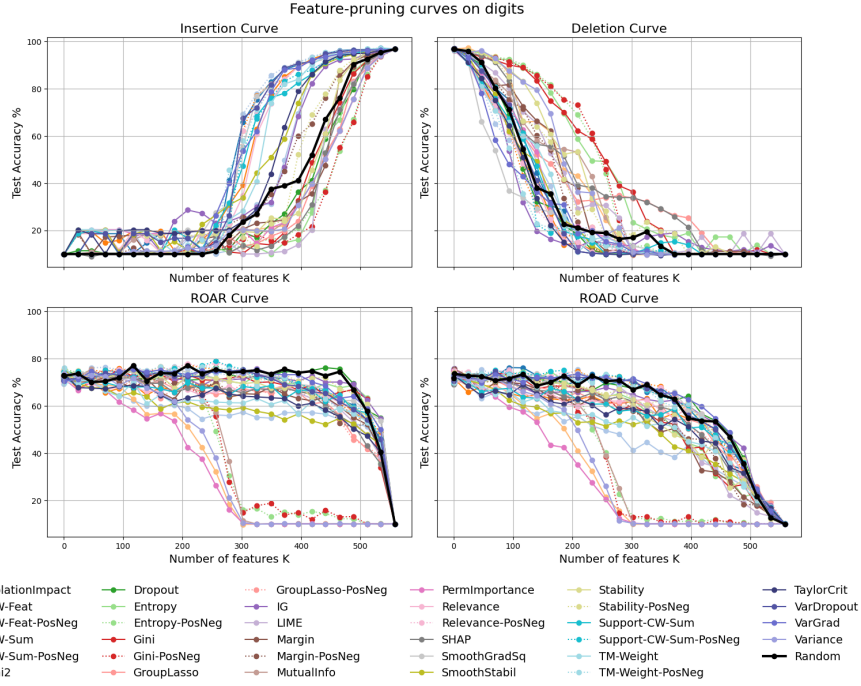


Fig. 2. Example of the pruning curves of the used Evaluation Protocols on the Digits dataset. Black is random as a reference.

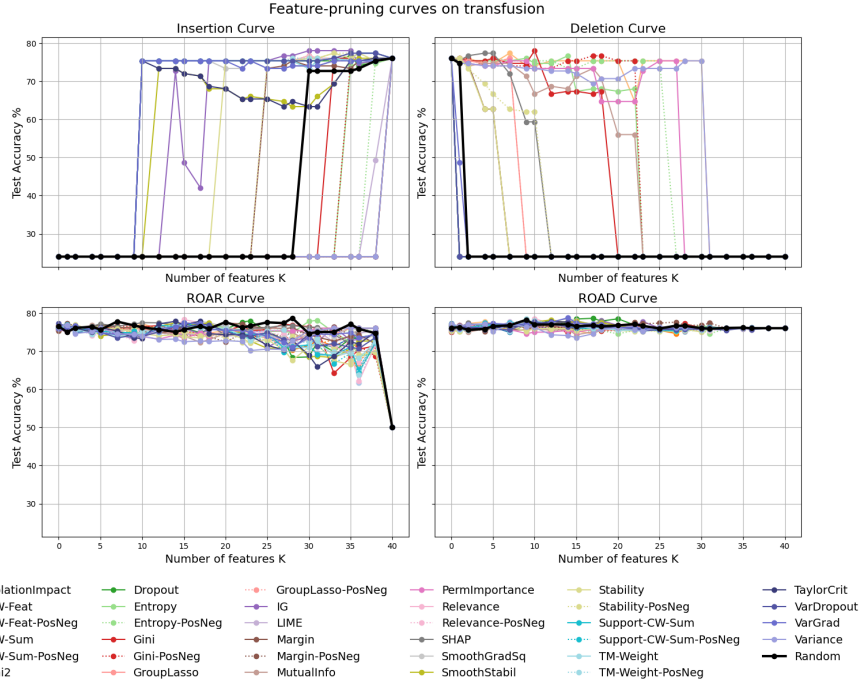


Fig. 3. Example of the pruning curves of the used Evaluation Protocols on the Transfusion dataset. Black is random as a reference.

ROAD curve of the Transfusion datasets seems nearly constant for all **FS** techniques, which could be caused by a much lower achieved F1 score by the used model, than in the Digits case. Achieved F1 score was 56 for Transfusion and 96 for Digits.

7 Results

Table in Figure 4 shows the count of each method appearing in the top-5 per protocol.

Chi2, Variance and MutualInfo filter methods, seem to dominate right after PermImportance wrapper. This dominance of simple filters can be attributed to our thermometer-encoding, which expands each original variable into several binary features; marginal-based statistics (variance, X^2 , MI) then readily identify those with the greatest spread or class-association, without needing to model inter-feature interactions. Best among its category are Chi2, PermImportance and Stability-PosNeg.

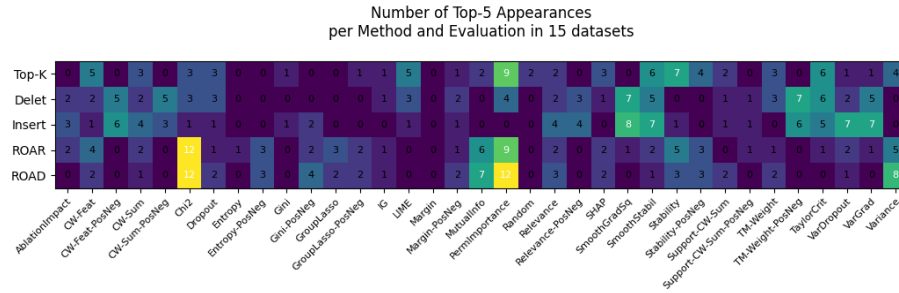


Fig. 4. Appearances in top-5 **AUC** scores across datasets and protocols. We look for consistency among all, or the most advanced Evaluation protocol only - the **ROAD**.

Figure 5 plots per-dataset top-3 methods' **AUC** vs. normalized ranking time. It also shows efficiency of each category vs its **AUC** with **ROAD** evaluation. Here we see that filters (red) are extremely fast but plateau at moderate **AUC**, whereas wrappers (blue) achieve slightly higher **AUC** at much greater cost—reflecting their retraining or perturbation loops—and embedded methods (green) strike a middle ground by leveraging TM internals at modest overhead.

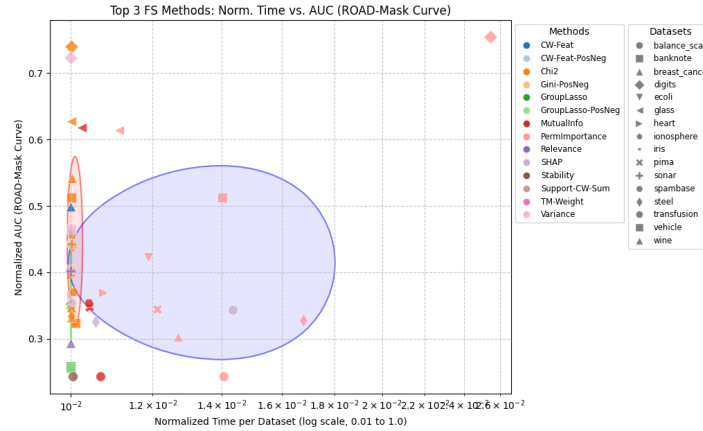


Fig. 5. Normalized **AUC** vs. ranking time for top-3 methods per dataset. Grouped average of each category by colour Red is Filter, Green is Embedded, Blue is Wrapper

We used pairwise Euclidean distances between each method’s average ROAD-mask AUC profiles across all datasets, and average-linkage clustering to produce the dendrogram in Figure 6. A few noteworthy patterns emerge:

- Filter methods (MutualInfo, Chi2, Variance) form a cleanly separated branch, reflecting their shared reliance on simple, data-centric statistics.
- LIME and SHAP do not cluster together; instead each associates with a different subset of embedded scorers, hinting at their distinct attribution biases even though both are “post-hoc” wrappers, which suggests that SHAP’s global Shapley-based attributions and LIME’s local surrogate fits emphasize different aspects of TM clause outputs.
- The two TM-internal scorers CW-Feat and CW-Sum sit in their own tight subcluster, underscoring their similar nature.
- Proposed embedded methods occupy a middle ground—structurally closer to each other than to filters or neural wrappers—thereby highlighting the unique, clause-based nature of TM feature selection.

Together, these clusters confirm that methods group by their reliance on marginal versus interaction-aware versus clause-based signals.

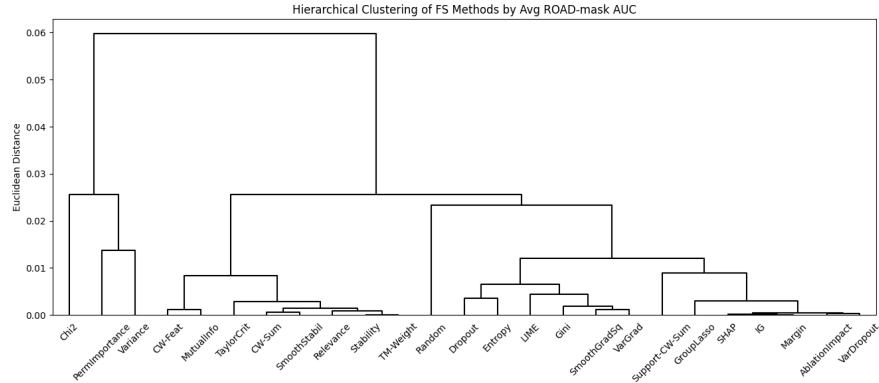


Fig. 6. Dendrogram of the used FS techniques, averaged among all datasets

Lastly, we validate alignment of the feature rankings with human intuition on the Digits dataset (Fig. 7). *Filters* Chi² and MutualInfo light up similar high-variance pixels at the digit boundaries, Variance focuses on the outline. *Wrappers* PermanentImportance, IG LIME and VarGrad reveal pixels whose shuffling most degrades class sums, but are unreadable. *Embedded* CW-Feat, TayCrit, show numbers; AlbImpa and VarDrop are heavily selective with their importance, and show individual pixels. SmoothStability and Margin seem to produce similar results to Wrappers. Entropy, Gini, Stability, show the numbers almost readably.

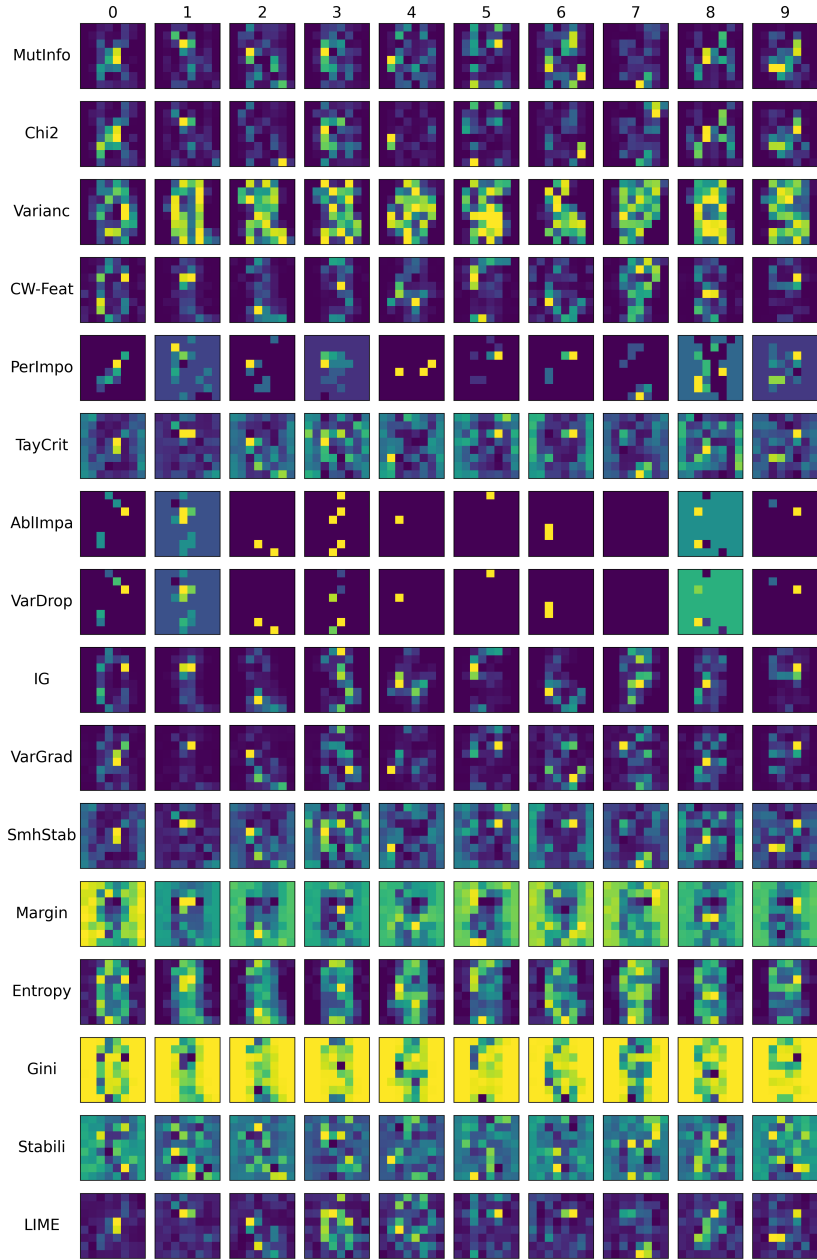


Fig. 7. Plots of the important features according to different FS Techniques - Proposed TM Embedded FS techniques are the most interpretable

8 Discussion and Future Work

Chi2 filter method either dominates, or is in the top 5 rankings, according to ROAD. It is on par with Permanent Importance wrapper, and Variance filter is lagging only one place behind them. Such high ranking of all the filters suggests that filters would be a great addition to the TM. This is most likely due to the thermometer encoding of all input features.

Future work could leverage feature scores to highlight not only important inputs but also which clauses and rules they influence, enabling symbolic-level interpretability.

Our experiments show that simple TM-embedded scorers (e.g. CW-Sum, TM-Weight) often match or exceed the quality of NN-inspired explainers at far lower cost. Classical filters are fast but miss higher level interactions. Key insights:

- Embedded TM scorers effectively capture clause-based interactions.
- Wrappers remain costly on Boolean inputs.
- Hybrid two-stage FS (filter and embedded) could further reduce cost.
- Online adaptive FS during TM training seems like a promising direction.
- FS ranking may be used to help with explainability, once the number of clauses grows too large

We presented the first systematic benchmark of filter, wrapper/post-hoc, and embedded FS methods on TMs, evaluated via Insertion/Deletion, ROAR, and ROAD. TM-internal scorers offer a middle ground, speed-quality trade-off. Moreover they also carry on the TM’s natural explainability, whilst other FS method fail at showing a human understandable scores. Future work will explore hybrid and online FS strategies, and compare directly against the same techniques on neural networks.

References

1. S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” *NeurIPS*, vol. 30, 2017.
2. M. T. Ribeiro, S. Singh, and C. Guestrin, “” why should i trust you?” explaining the predictions of any classifier,” in *SIGKDD*, 2016, pp. 1135–1144.
3. O.-C. Granmo, “The tsetlin machine—a game theoretic bandit driven approach to optimal pattern recognition with propositional logic. arxiv 2018,” *arXiv preprint arXiv:1804.01508*, 2018.
4. G. T. Berge, O.-C. Granmo, T. O. Tveit, M. Goodwin, L. Jiao, and B. V. Matheussen, “Using the Tsetlin Machine to Learn Human-Interpretable Rules for High-Accuracy Text Categorization,” *IEEE Access*, vol. 7, pp. 115 134–115 146, 2019.
5. V. Halenka, A. K. Kadhimi, P. F. A. Clarke, B. Bhattarai, R. Saha, O.-C. Granmo, L. Jiao, and P.-A. Andersen, “ Exploring Effects of Hyperdimensional Vectors for Tsetlin Machines ,” in *2024 International Symposium on the Tsetlin Machine (ISTM)*. IEEE Computer Society, Aug. 2024, pp. 1–8.
6. R. Battiti, “Using mutual information for selecting features in supervised neural net learning,” *IEEE Trans. on NN*, vol. 5, no. 4, pp. 537–550, 1994.
7. W.-Y. Loh, “Classification and regression trees,” *IEEE Trans. NN*, vol. 1, no. 1, pp. 14–23, 2011.
8. A. Fisher, C. Rudin, and F. Dominici, “All models are wrong, but many are useful,” 2019. [Online]. Available: <https://arxiv.org/abs/1801.01489>
9. L. Ceriani and P. Verme, “The origins of the gini index,” *J. Econ. Inequality*, vol. 10, pp. 421–443, 2012.
10. M. Sundararajan, A. Taly, and Q. Yan, “Axiomatic attribution for deep networks,” in *Proceedings of the 34th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 70. PMLR, 06–11 Aug 2017, pp. 3319–3328.
11. D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, “Smoothgrad: removing noise by adding noise,” *arXiv:1706.03825*, 2017.
12. M. L. Tsetlin, “Automaton theory and modeling of biological systems,” *MATH. SCI. ENGNG*, 1973.
13. G. Brown, A. Pocock, M.-J. Zhao, and M. Luján, “Conditional likelihood maximisation: a unifying framework for information theoretic feature selection,” *J. Mach. Learn. Res.*, vol. 13, no. null, p. 27–66, Jan. 2012.
14. S. Nogueira, K. Sechidis, and G. Brown, “On the stability of feature selection algorithms,” *JMLR*, vol. 18, no. 174, pp. 1–54, 2018.
15. G. Montavon, W. Samek, and K.-R. Müller, “Methods for interpreting and understanding dnn,” *Digit. Signal Processing*, vol. 73, pp. 1–15, 2018.
16. N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *JMLR*, vol. 15, no. 56, pp. 1929–1958, 2014.
17. M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks,” in *ECCV 2014*. Springer, 2014, pp. 818–833.
18. P. W. Koh and P. Liang, “Understanding black-box predictions via influence functions,” in *ICML*. PMLR, 2017, pp. 1885–1894.