

Discovering Association Rules in High-Dimensional Small Tabular Data

Erkan Karabulut¹ ✉ (e.karabulut@uva.nl), Daniel Daza², Paul Groth¹, Victoria Degeler¹
¹University of Amsterdam, ²Amsterdam University Medical Center



1 Learning Rules?

Knowledge discovery: Reveal associations between data features, e.g., columns of a given table (ARM).

Interpretable inference: Draw conclusions using learned rules instead of black box models, such as classification rules.

Formalization: Table with k features $F = \{f_1, \dots, f_k\}$, each with classes $f_i^1, \dots, f_i^{c_i}$. Define the item universe $I = \{f_i^j \mid 1 \leq i \leq k, 1 \leq j \leq c_i\}$.

Each row (transaction) $T \subset I$ satisfies $\forall i \in \{1, \dots, k\}, \exists! j \in \{1, \dots, c_i\}, f_i^j \in T$. An association rule is $X \rightarrow Y$ with $X, Y \subset I, X \cap Y = \emptyset, |Y| = 1$.

Logical form: $X \rightarrow Y \equiv (\neg \wedge_{x \in X} x) \vee y$ (Horn clause in CNF).

Gene expression datasets:

gene1	gene2	gene3	...	gene18107
normal	normal	normal		normal
normal	normal	high	...	high
normal	normal	normal		low

Gao et al. 2015.

$Gene2(high) \wedge Gene29(high) \rightarrow Gene14(low)$

Example: 20 features, 4 values each, $a = 4 \Rightarrow$ **5, 186, 240** rules!

Rule Explosion:

Even small datasets can generate millions of redundant or trivial rules --- slow, hard to interpret, and unscalable.

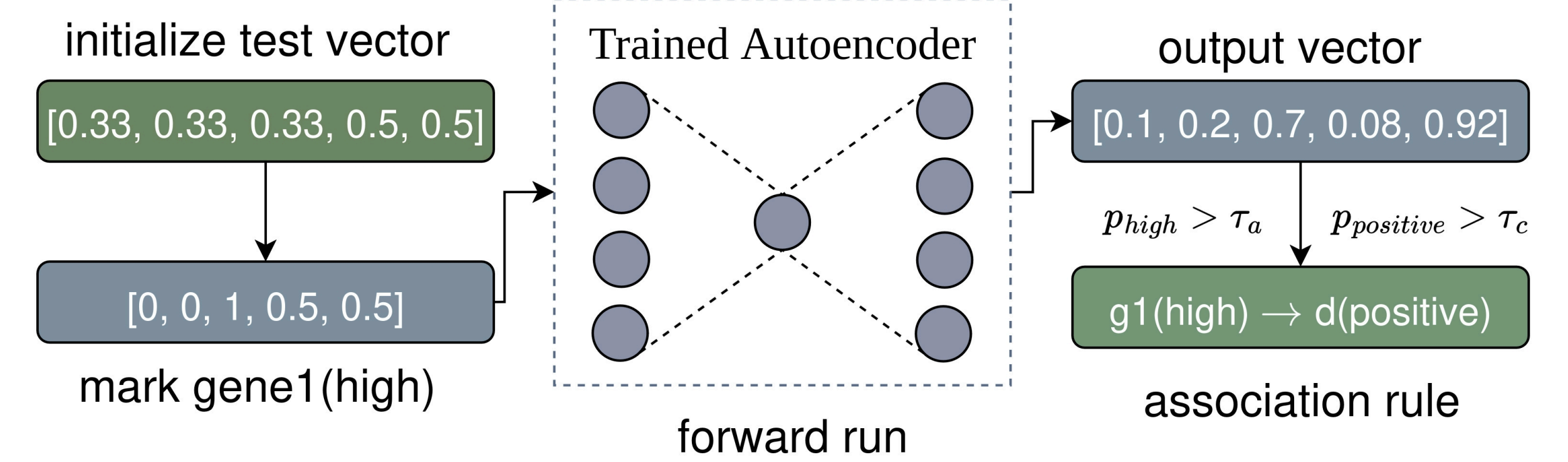
$X \rightarrow Y, |X| + |Y| \leq a$

$$\# \text{ Rules} = \sum_{p=1}^a c_p! \left(\prod_{i \neq p} (c_i + 1) - 1 \right)$$

2 Neurosymbolic ARM

Intuition: Autoencoders capture feature associations via reconstruction. If, after training, a forward pass with marked categories A reconstructs categories C with high probability, then $A \rightarrow C \setminus A$ (no self-implication).

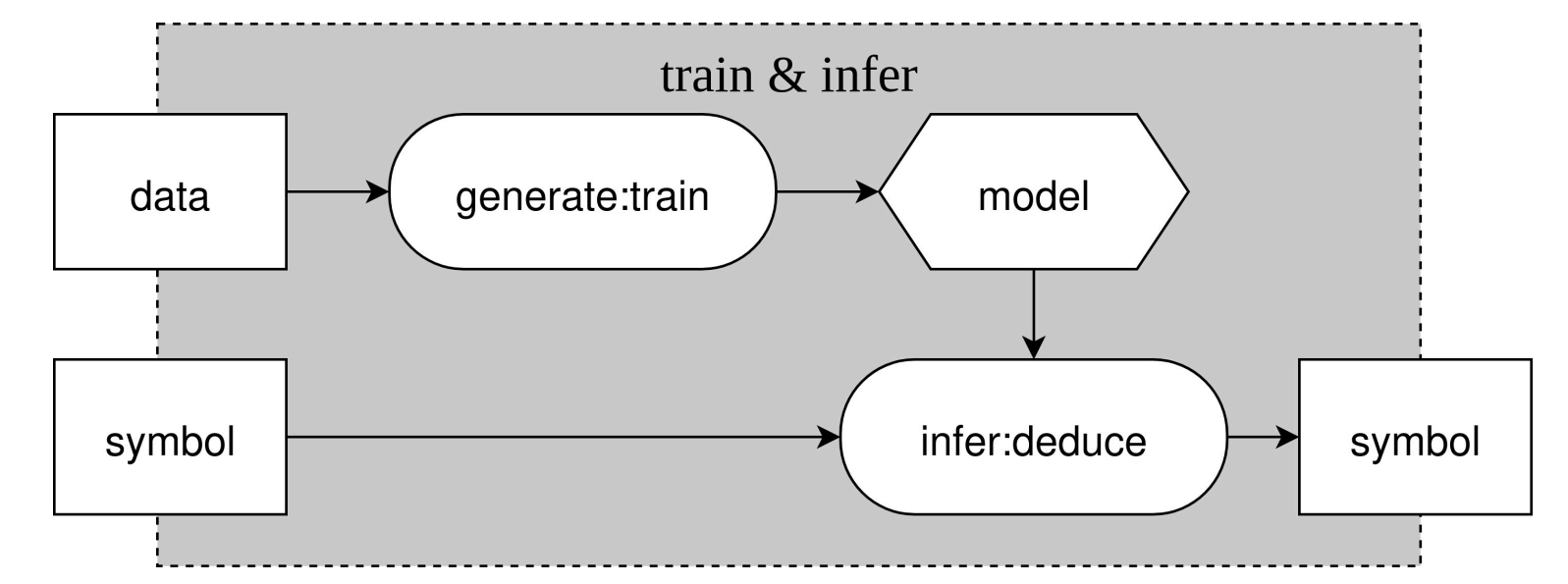
gene1 = {low, normal, high}, disease = {negative, positive} $\tau_a = 0.5, \tau_c = 0.8$



(Aerial+), Karabulut et al. 2025

Observation:

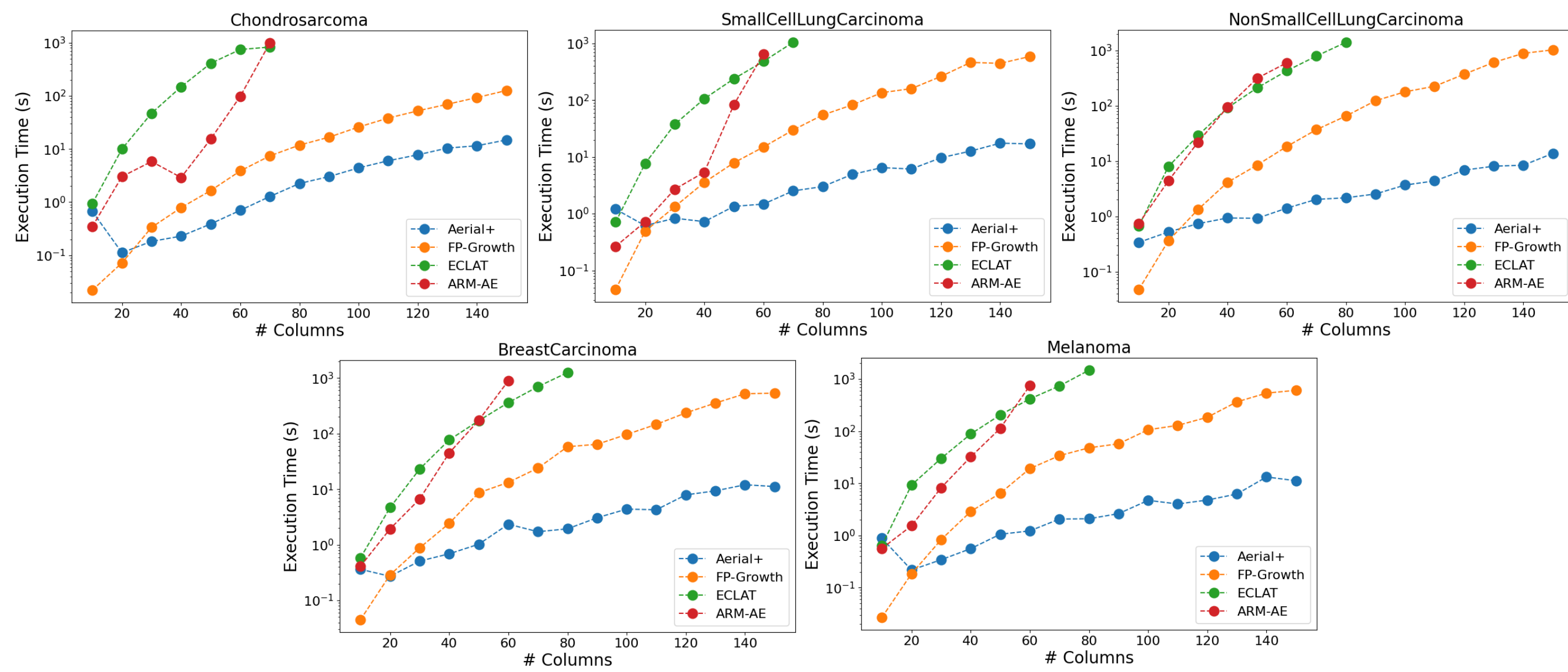
The reconstruction objective of Aerial+ prioritizes significant patterns rather than redundant patterns!



(Boxology), van Bekkum et al. 2021

3 Empirical Analysis Neurosymbolic ARM is Scalable!

Empirical result: NeSy method scales 1-2 orders of magnitude faster, despite low data!

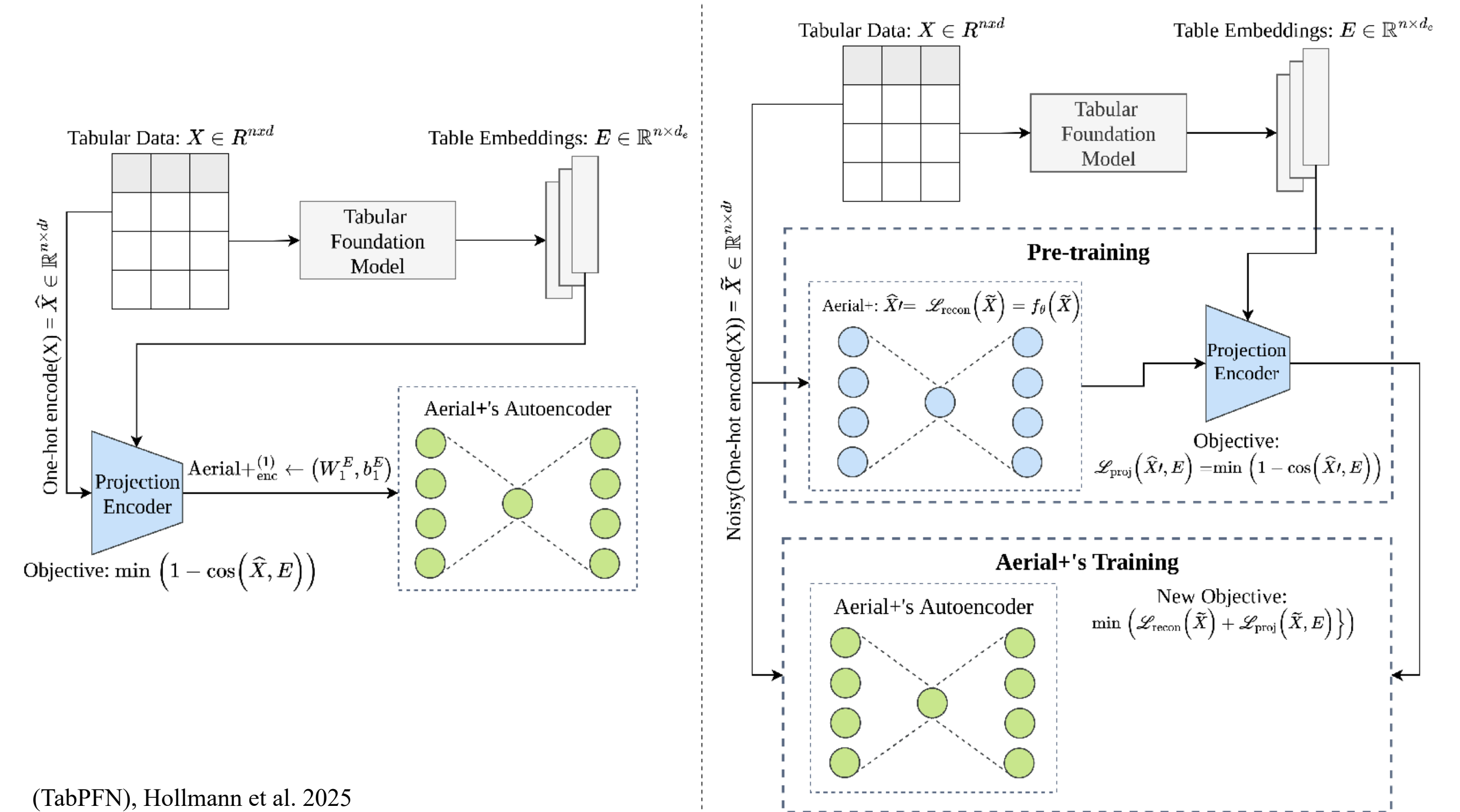


Dataset	# Columns	# Rows	Algorithm	Type	Parameters
Chondrosarcoma	18006	6	Aerial+	Neurosymbolic	$a = 2, \tau_a = 0.5, \tau_c = 0.8$
SmallCellLungCarcinoma	18237	60	ARM-AE	Neurosymbolic	$M=2, N= R / C , L=0.5$
NonSmallCellLungCarcinoma	18108	86	FP-Growth	Algorithmic	(both) antecedents = 2, min_conf=0.8,
BreastCarcinoma	18061	51	ECLAT	Algorithmic	min_support=0.5 * $\mathbb{E}[\text{support}(R)]$
Melanoma	17902	55			

4 How to Boost Knowledge Discovery in a Low-Data Regime?

Weight initialization (Aerial+WI)

Double loss (Aerial+DL)



(TabPFN), Hollmann et al. 2025

Weight initialization from tabular foundation models (Aerial+WI) and semantic alignment of table embeddings with code layers using a dual loss function (Aerial+DL)

5 Validation, Interpretation and Future Research

Empirical result: A concise set of higher-quality association rules!

Approach	# Rules	~Rule Coverage	~Support	~Confidence	Data Coverage	~Zhang's Metric	Exec. Time (s)
Chondrosarcoma							
Aerial+	200	0.23	0.21	0.921	0.533	0.784	2.25
Aerial+WI	75	0.217	0.206	0.945	0.524	0.813	5.80
Aerial+DL	75	0.235	0.219	0.947	0.536	0.828	5.36
SmallCellLungCarcinoma							
Aerial+	1576	0.068	0.041	0.579	0.835	0.476	10.58
Aerial+WI	664	0.076	0.052	0.633	0.715	0.577	13.48
Aerial+DL	1338	0.070	0.044	0.597	0.816	0.513	18.23
NonSmallCellLungCarcinoma							
Aerial+	1620	0.059	0.035	0.584	0.823	0.554	18.03
Aerial+WI	978	0.078	0.057	0.663	0.698	0.639	28.67
Aerial+DL	1453	0.053	0.028	0.547	0.849	0.501	24.27
BreastCarcinoma							
Aerial+	1017	0.072	0.046	0.641	0.816	0.575	9.64
Aerial+WI	590	0.077	0.052	0.686	0.686	0.644	12.09
Aerial+DL	535	0.078	0.050	0.652	0.761	0.590	15.31
Melanoma							
Aerial+	1220	0.067	0.035	0.545	0.888	0.440	13.09
Aerial+WI	773	0.070	0.038	0.575	0.772	0.496	13.19
Aerial+DL	859	0.071	0.038	0.566	0.860	0.461	16.49

Metrics:

$\text{Support}(X \rightarrow Y) = |\{T : X \cup Y \subseteq T\}| / n$

$\text{Confidence}(X \rightarrow Y) = |\{T : X \cup Y \subseteq T\}| / |\{T : X \subseteq T\}|$

$\text{Coverage}(X \rightarrow Y) = |\{T : X \subseteq T\}| / n$

$$\text{Zhang}(X \rightarrow Y) = \frac{\text{conf}(X \rightarrow Y) - \text{conf}(X' \rightarrow Y)}{\max(\text{conf}(X \rightarrow Y), \text{conf}(X' \rightarrow Y))}$$

• **Dataset dependency:** Algorithmic methods' execution time increases with data density (many high-support itemsets), while Aerial+ maintains constant polynomial-time extraction regardless of density.

• **Validation scope:** Testing on more diverse domains and datasets with higher instance-to-feature ratios ($n \gg d$) needed to assess convergence and rule quality

• **Foundation model constraint:** Current approach limited to TabPFN (only available tabular foundation model with table embedding interface); designed for classification/regression rather than column associations

• **Background knowledge in knowledge discovery:** What other types of knowledge, e.g., structured or bayesian, can be utilized in knowledge discovery?

Code



Paper

