# INFERRING THE INVISIBLE: NEURO-SYMBOLIC RULE DISCOVERY FOR MISSING VALUE IMPUTATION

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

025

026

028 029

031

033

034

035

036

037

040

041

042

043

044

046

047

048

049

050 051

052

Paper under double-blind review

#### **ABSTRACT**

One of the central challenges in artificial intelligence is reasoning under partial observability, where key values are missing but essential for understanding and modeling the system. This paper presents a neuro-symbolic framework for latent rule discovery and missing value imputation. In contrast to traditional latent variable models, our approach treats missing grounded values as latent predicates to be inferred through logical reasoning. By interleaving neural representation learning with symbolic rule induction, the model iteratively discovers—both conjunctive and disjunctive rules—that explain observed patterns and recover missing entries. Our framework seamlessly handles heterogeneous data, reasoning over both discrete and continuous features by learning soft predicates from continuous values. Crucially, the inferred values not only fill in gaps in the data but also serve as supporting evidence for further rule induction and inference—creating a feedback loop in which imputation and rule mining reinforce one another. Using coordinate gradient descent, the system learns these rules end-to-end, enabling interpretable reasoning over incomplete data. Experiments on both synthetic and real-world datasets demonstrate that our method effectively imputes missing values while uncovering meaningful, human-interpretable rules that govern system dynamics.

#### 1 Introduction

Neural-symbolic reasoning combines the pattern recognition power of neural networks with the precision and interpretability of symbolic reasoning (Hitzler & Sarker, 2022; Yang et al., 2024). This hybrid paradigm enables AI systems to detect complex patterns in unstructured data while reasoning about them in a structured and explainable manner.

Traditional rule induction methods extract explicit patterns from observed data but often fail when *some observations are missing or incomplete* (Campero et al., 2018; Claire Glanois, 2022). These approaches can effectively learn surface-level rules, yet their ability to fully explain the underlying system is limited when essential data points are absent. For example, in healthcare diagnostics, critical measurements may be missing or noisy, making accurate imputation necessary for reliable reasoning.

Probabilistic models such as Markov Logic Networks (MLNs) (Richardson & Domingos, 2006) handle missing data by treating unobserved facts as latent predicates. However, they typically rely on a *fixed* rule base and *expensive joint inference*, limiting scalability and adaptability in large or heterogeneous datasets (Oltramari et al., 2020). In contrast, we propose a neuro-symbolic system that *co-learns rules and imputations* in a single differentiable loop, enabling fast forward-chaining inference and end-to-end learning.

Our core idea is a closed loop between imputation and rule discovery. Given partially observed tables with discrete and continuous attributes, we treat each missing, entity-specific entry as an unknown fact and apply learned rules in a forward-chaining pass to predict it. These predictions are compared to the observed entries via a supervised loss, and backpropagation updates the rule parameters and soft predicates. Crucially, improved imputations provide additional evidence for discovering and refining rules in subsequent passes. This self-reinforcing loop leads to better imputations, improved rule induction, and stronger downstream inference.

To enable multi-hop reasoning at scale, many targets require compositional explanations in the form of chains and disjunctions. We optimize rule embeddings using *asynchronous coordinate gradient descent*, updating one rule or clause at a time while holding others fixed. This mirrors step-wise reasoning and ensures monotone loss progress on a smooth surrogate. For disjunctive heads, we

adopt a sequential covering strategy to harvest diverse clauses, followed by joint fine-tuning using a soft-OR aggregator (LogSumExp) to reconcile interactions. This staged procedure reliably recovers long chains and disjunctive theories under high missingness while keeping computation tractable.

Our framework handles heterogeneous data by learning *soft predicates* for continuous features (using sigmoid thresholds and slopes) and combining them with discrete predicates through differentiable logical operators. Specifically, we use soft-min to approximate logical AND and soft-max to approximate logical OR. This approach enables uniform forward chaining over mixed data types without requiring pre-discretization.

**Contributions.** We summarize our contributions as follows: (*i*) We introduce a closed-loop neurosymbolic framework in which imputation and rule discovery mutually reinforce each other, rather than treating imputation as a preprocessing step. (*ii*) We develop a scalable coordinate gradient descent scheme, combined with sequential covering and joint fine-tuning, that enables multi-hop and disjunctive rule learning even under high missingness. (*iii*) We design a unified differentiable forward-chaining engine that handles both discrete and continuous attributes through soft predicates and smooth logical operators. (*iv*) We empirically validate our approach on synthetic chain and disjunction tasks, as well as real-world datasets (Birds, Heart, SPECT), demonstrating that it recovers human-interpretable rules while achieving strong imputation accuracy and downstream prediction performance.

#### 2 RELATED WORK

Our work is at the intersection of neuro-symbolic Inductive Logic Programming (ILP) and missing value imputation.

**Neural Embedding-based ILP.** Embedding-based models are widely used for Knowledge Base (KB) completion like TransE (Bordes et al., 2013), TransH (Wang et al., 2014), and TransR (Lin et al., 2015). Complex (Trouillon et al., 2016) introduces complex-valued embeddings for asymmetric relations, while multi-hop reasoning methods like Guu et al. (2015) leverage path-based embeddings for traversing knowledge graphs. However, these approaches often face limitations in reasoning power.

Recent advances in ILP integrate symbolic logic with neural networks. Rocktäschel & Riedel (2017) propose *Neural Theorem Proving (NTP)*, which uses a differentiable backward-chaining method. Then, Campero et al. (2018) introduces a neural forward-chaining differentiable rule induction network. However, both rely on hand-designed templates. Claire Glanois (2022) advances these models by incorporating a hierarchical structure, enabling more flexible rule induction. Nevertheless, these methods are primarily designed for fully-observed data and struggle to handle missing values.

**Rule-Based Missing Value Imputation.** Traditional missing data imputation methods, ranging from statistical techniques like SOFT-IMPUTE (Mazumder et al., 2010) to deep learning models like MMDL (Li et al., 2020), typically rely on statistical patterns and do not leverage explicit logical rules to govern inter-variable relationships (see Appendix A for a detailed overview).

Recent works have started to bridge rule-based reasoning and missing value imputation. For instance, Chen et al. (2023) employ various interpretable machine learning techniques to address the missing value problem, but their methods are not explicitly rule-based. Closer to our approach, MINTY (Stempfle & Johansson, 2024) utilizes a rule-based model to handle missing data; however, it does not leverage neuro-symbolic reasoning to learn the intricate relationships between observed and missing values as we do. Other non-neural approaches, such as the work by Wang et al. (2017) on synthesizing data completion, also tackle the problem but lack of the representation learning capabilities of neural networks. Our work is distinct in its tight integration of neural learning for representation and symbolic reasoning for both rule discovery and imputation, forming a feedback loop where each component enhances the other.

#### 3 BACKGROUND

**Predicate.** In the context of logic-based AI systems, a predicate is a fundamental Boolean logic variable used to describe properties of or relationships between entities. Predicate variables are grounded by data, being True or False, and serve as the basic building blocks for logical expressions. For instance, a predicate like  $Has\_Fever(Patient)$  denotes whether a patient has a fever, while  $Use\_Drug(Patient)$  specifies whether a drug treats a particular patient. These predicates capture essential aspects of the system's state and relationships.

#### Logic Rules and Forward Chaining. We represent knowledge with Horn clauses

$$f: Q \leftarrow P_1 \wedge P_2 \wedge \dots \wedge P_h,$$
 (1)

where  $P_1, \ldots, P_h$  (the body) are conditions and Q (the head) is the conclusion. Given observed facts (the evidence set  $\mathcal{E}$ ), we perform forward chaining: whenever all body predicates of a rule are (approximately) satisfied by facts in  $\mathcal{E}$ , the rule fires and adds Q to  $\mathcal{E}$ . Importantly, newly inferred facts are immediately recycled as evidence, enabling multi-hop reasoning—cascades of rule applications that derive conclusions not reachable in a single step.

**Latent Predicates and Rule Learning.** We use the term *latent predicate* to denote an unobserved fact tied to concrete entities (and, when relevant, timestamps) within the same relational schema as observed predicates. Latent predicates may be Boolean or soft-valued (degrees of truth); they represent missing-but-specific facts we wish to infer. Our goal is to learn *Horn rules* of the form Eq. (1) that capture regularities among observed predicates and support inference about latent ones—i.e., rules whose heads or intermediate conclusions may involve latent predicates, enabling principled completion of missing facts.

**Expressive Rule Forms.** We consider rules that capture rich logical structure, including conjunctions (AND), disjunctions (OR via multiple clauses), and *chained dependencies*. For example, a latent predicate  $Q_k$  may be characterized by

$$Q_k = (P_1 \wedge P_2) \vee (P_3 \wedge P_4),$$

or by multi-hop compositions such as

$$Q_1 = P_1 \wedge P_2, \quad Q_2 = P_3 \wedge P_4, \quad Q_3 = (Q_1 \wedge P_5) \vee (Q_2 \wedge P_6).$$

This view accommodates both single-step and multi-step (multi-hop) reasoning patterns within a unified Horn-rule framework. We also allow *predicate invention*: introducing unlabeled latent predicates that are not predefined in the schema but are useful intermediates for explaining the data. These invented predicates participate in rules just like observed ones. After rules are discovered, their roles can be *post-hoc interpreted* by inspecting the clauses in which they appear and their relationships to observed predicates.

#### 4 MODEL: NEURO-SYMBOLIC FORWARD CHAINING NETWORK

Consider problems where some information or features are incomplete. Our goal is to learn a set of logical rules that explain how each predicate with information can be imputed based on evidence from feature space **X**.

These missing variables are inferred through a rule-learning process, allowing the model to uncover hidden relationships in the data. For clarity, we identify the predicates with missing information as U, also named as "latent predi-

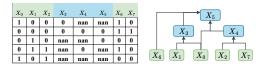


Figure 1: Example of missing variables imputation with rule discovery.  $X_i$  with nan is the predicates with missing information, which can be inferred by the logic rules from X.

cate" in our setting. Though in our experiments, we do not strictly distinguish between feature predicates, as any of them can be incomplete and serve as latent predicates. In more general settings with a predictive label Y, we can view Y as one of the latent predicates, making the rule learning and prediction for Y equivalent to inferring latent predicates U with rules.

To summarize, our model learns logical rules to infer latent predicates  $\mathbf{U}$  by discovering hidden structures within data, as an example illustrated in Figure 1. This rule induction process identifies logical relationships among observable predicates  $\mathbf{X}$  and other inferred latent predicates. By explicitly learning these structures, our approach enhances both inference capability and interpretability, offering clear insights into complex, otherwise hidden dependencies. The key idea is summarized in Figure 2, with details presented in the following sections.

#### 4.1 Model Preparation: Pretrained Predicate Embeddings

We begin by defining two sets of predicates:  $\mathbf{X} = \{X_1, \dots, X_n\}$  represents the set of observable predicate variables, and  $\mathbf{U} = \{U_1, \dots, U_m\}$  denotes the set of predicate variables with missing information that the model aims to discover and define. Our framework is designed to handle **both binary (categorical) and continuous features within a unified logical structure**. Binary features are treated as standard logical predicates. For continuous features, we introduce a mechanism to derive a "soft" truth value, effectively creating learnable predicates from them. This allows the model to reason over heterogeneous data types, as detailed in Section 4.2.

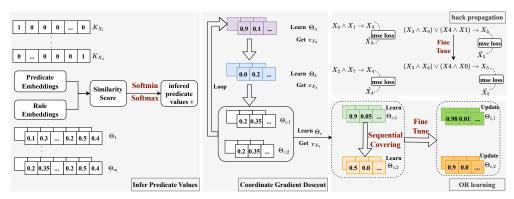


Figure 2: Model framework. Rule embeddings  $\Theta$  are optimized using coordinate gradient descent. In each learning step, predicate values are inferred via the Softmin-Softmax operation (Eqs. (3, 4,5)). For disjunctive (OR) rule learning, sequential hard covering is applied, followed by fine-tuning of the learned rule embeddings (Section 5.2). Errors are back-propagated using MSE loss between inferred predicate values and the small portion of observed latent predicate samples, constituting a weak-supervision setting.

As mentioned before, we do not distinguish X and U in the experiment, as any predicates can be the predicate with missing information. We just use separate notations for model description. We initialize a fixed, unique embedding for each predicate, whether observable or missing. For example, these embeddings can be instantiated as one-hot vectors within an embedding space of dimension d. We denote the collection of embeddings for observable predicates as  $K_X$  and incomplete predicates as  $K_U$ . These predicate embeddings remain frozen throughout the rule learning phase and serve as a foundational dictionary, enabling the interpretation of the composition of learned rules by relating rule components back to specific predicates.

With the predicate representations defined, we next describe the core of our model: the representation of logical rules and the mechanism by which inferences are drawn.

#### 4.2 MODEL BACKBONE: RULE REPRESENTATION AND INFERENCE

In our NS-FCN framework, logical rules are materialized as learnable rule embeddings, which are the primary trainable parameters. Our model employs an asynchronous coordinate descent learning process. This learning scheme is particularly well-suited for discovering complex logical structures such as chained dependencies (where one latent predicate forms part of the definition of another) and disjunctive rules (where a latent predicate can be satisfied by one of several distinct conditions).

#### 4.2.1 Specification of Rule Embeddings $\Theta$

The parameter  $\Theta = [\Theta_f]_{f \in \mathcal{F}}$  represents the set of logical rules we aim to learn, denoted as  $\mathcal{F} := \{f\}$ . Each matrix  $\Theta_f$  encodes a logical rule that can infer latent predicates.

Conjunctive Rule Embedding. For a latent predicate  $U_j$  that is defined by a single conjunctive rule (e.g.,  $U_j = X_a \wedge X_b$ ), its corresponding rule embedding  $\Theta_f = [\theta_1, \dots, \theta_h] \in \mathbb{R}^{d \times h}$ . Here, h represents the number of predicates forming the body of the conjunctive rule (the arity of the conjunction, e.g., h = 2 for  $X_a \wedge X_b$ ), and d is the dimensionality of the predicate embeddings. Each of the h rows in this matrix is learned to align with the embedding of one of the constituent predicates in the rule's body.

Disjunctive Rule Embeddings. If a latent predicate  $U_k$  is defined by a disjunction of  $R_k$  distinct conjunctive clauses (e.g.,  $U_k = \bigvee_{r=1}^{R_k} (\text{clause}_r)$ ), it will be associated with a set of  $R_k$  distinct rule embeddings, denoted  $\{\Theta_{k,1}, \ldots, \Theta_{k,R_k}\}$ . Each individual rule embedding  $\Theta_{k,r}$  is itself an  $h_r \times d$  matrix, representing the r-th conjunctive clause, where  $h_r$  is the arity of that specific clause.

All rule embeddings are initialized randomly prior to training and are subsequently optimized as described in Section 5. Given these rule embeddings, the model infers the truth values (or continuous approximations thereof) of latent predicates through a carefully defined inference mechanism.

**Parameters for Continuous Predicates.** For each continuous feature  $f \in \mathcal{F}_C$ , where  $\mathcal{F}_C$  is the set of continuous features, the model learns two additional scalar parameters: a threshold  $\theta_f$  and a slope  $\beta_f$ . These parameters are used to define a learnable soft predicate function that maps the continuous feature value to a probabilistic truth value, as explained next.

219

220

221

222

224 225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240 241

242

243

244

245

246

247

249

250

251

253

254

256 257

258 259

260 261

262

264

265

266 267

268

#### 4.2.2 INFERRING PREDICATE VALUES

The latent predicates is inferred based on the current state of observable predicates, any previously inferred latent predicate values, and the learned rule embeddings  $\Theta$ .

**Predicate Matching.** Each column  $\theta_i (j = 1, ..., h)$  in the rule embedding  $\Theta_f$  is matched with a corresponding predicate embedding. This matching is achieved by finding the predicate embedding most similar to  $\theta_i$  using cosine similarity:

$$K_j^* = \operatorname*{argmax} \cos(K, \theta_j), \quad j = 1, \dots, h \tag{2}$$

 $K_j^* = \operatorname*{argmax} \cos\left(K, \theta_j\right), \quad j = 1, \ldots, h$  (2) where  $\mathbf{K} = \mathbf{K}_X \cup \mathbf{K}_U$  represents the set of all available predicate embeddings. The inverse mapping I(K) maps a predicate embedding  $K \in \mathbb{R}^d$  back to its corresponding index. Thus, indices  $1, \ldots, (n+m)$  correspond to n+m predicate embeddings.

**Predicate Truth Values.** Once the best matching predicate  $K_i^*$  is identified for a rule component  $\theta_i$ , we determine its truth value, denoted as  $\tau_i$ . The calculation depends on whether the corresponding feature is binary or continuous:

- 1) For a binary feature (e.g., from one-hot encoding), its truth value is its current value in the data tensor:  $\tau_j = \boldsymbol{v}^t \left( I\left( K_i^* \right) \right)$ .
- 2) For a *continuous feature*, its truth value is computed using a learnable **soft predicate** function (a sigmoid):  $\tau_j = \sigma(\beta_{f_j} \cdot (x_{f_j} - \theta_{f_j}))$  where  $x_{f_j}$  is the value of the feature corresponding to  $K_j^*$ ,  $\theta_{f_j}$  and  $\beta_{f_j}$  are its learned parameters, and  $\sigma(\cdot)$  is the sigmoid function. This allows learning soft boundaries like " $x_{f_i} > \theta_{f_i}$ ".

Conjunctive Clause Inference (Soft-AND). The value for a conjunctive clause is then computed by aggregating the contributions of all its components, modeling a Soft-AND operation. The contribution of each component j is the product of its similarity score and its truth value. The aggregated value is:

$$v = \prod_{j=1,\dots,h} \cos\left(K_j^*, \theta_j\right) \cdot \tau_j,\tag{3}$$

where  $v^t$  is the current value for observable predicates or any previously imputed values. At the beginning,  $v^t$  is all from observable predicates. With the optimization steps of coordinate descent,  $v^t$  is updated based on the refined  $\Theta$ .

To address the potential issue of diminishing values, we can use the min function instead: v = $\min_{i=1,\ldots,h} \left\{ \cos \left( K_i^*, \theta_i \right), \tau_i \right\}.$ 

However, to make this function differentiable, we approximate the min function using the softmin function. For each component j, there are two terms: the similarity score  $\cos(K_i^*, \theta_i)$  and the truth value  $\tau_j$ . The softmin is applied to the set of all 2h such terms:

softmin 
$$(x_1, \dots, x_{2h}; \Theta) = -\frac{1}{\tau} \log \left( \frac{1}{2h} \sum_{i=1}^{2h} e^{-x_i/\tau} \right)$$
 (4)

where each  $x_i$  represents one of the 2h terms (all similarity scores and all truth values), and  $\tau$  is a temperature parameter controlling the smoothness of the approximation. As  $\tau$  approaches 0, the softmin function approximates the behavior of the hard min function.

**Disjunctive Rule Inference (Soft-OR).** When a latent predicate  $U_k$  is defined by a disjunction of multiple conjunctive clauses,  $U_k = \bigvee_{r=1}^{R_k} \text{clause}_{k,r}$ , its final inferred value  $v_{U_k}$  is determined by aggregating the values of its individual clauses  $\{v_{\text{clause}_{k,1}}, \dots, v_{\text{clause}_{k,R_k}}\}$ . This aggregation is performed using the LogSumExp (LSE) function, which serves as a differentiable soft-OR operator:

$$v_{U_k} = \frac{1}{\beta} \log \sum_{r=1}^{R_k} \exp(\beta \cdot v_{\text{clause}_{k,r}}), \tag{5}$$

where  $\beta$  is a temperature parameter. As  $\beta \to \infty$ , the LSE function increasingly approximates the true max operator, thereby hardening the OR logic. Conversely, smaller values of  $\beta$  yield a softer aggregation. The model's ability to discover meaningful rules and infer latent predicate states accurately hinges on an effective learning procedure. We now outline the training methodology employed to optimize the rule embeddings  $\Theta$ .

#### 5 MODEL LEARNING

The core of our model learning process involves training the rule embeddings  $\Theta$  by minimizing a loss function that quantifies the discrepancy between the inferred values of latent predicates and

their partially observed truth values. Our approach leverages a sequential and staged optimization strategy, drawing parallels with coordinate descent and incorporating elements of rule covering, particularly for disjunctive rules. This is typically followed by a joint fine-tuning phase for rules involving disjunctions.

#### 5.1 COORDINATE GRADIENT DESCENT FOR RULE OPTIMIZATION

We employ a coordinate gradient descent approach, iteratively optimizing the embedding  $\Theta_j$  for each predicate  $U_j$  while holding the embeddings of other predicates fixed. The order in which predicates  $U_j$  are selected for optimization is randomized in each complete pass (cycle) through all learnable latent predicates. Such optimization progress is similar to human thinking strategy, as we humans usually draw conclusions step by step.

During the optimization step for a specific predicate  $U_j$  within a cycle, the inferred value  $v_{U_j}$  is obtained by Eq. 3 or Eq. 4 as mentioned in the previous Section. The Mean Squared Error (MSE) loss is computed between the inferred value  $v_{U_j}$  and its observed value  $U_{j,\text{obs}}$ , exclusively for instances where  $U_j$  is observed, which can be viewed as a **weakly supervision** setting:

$$\mathcal{L}_{U_j} = \text{mean}((v_{U_j}[\text{mask}_j] - U_{j,\text{obs}}[\text{mask}_j])^2), \tag{6}$$

where  $\operatorname{mask}_j$  is a binary vector indicating observed instances of  $U_j$ . The rule embedding  $\Theta_j$  is then updated using gradients from this loss.

After its training epochs within a cycle, if  $\Theta_j$  meets the criteria for a "perfect rule" (i.e., the imputation accuracy of missing variables is larger than 0.99 and a marginal loss drop is less than  $10^{-3}$ ), the parameters of  $\Theta_j$  will be frozen for efficient computing in subsequent cycles.

#### 5.2 SEQUENTIAL COVERING AND FINE-TUNING OF DISJUNCTIVE RULES

**Sequential Covering.** When a latent predicate  $U_k$  is hypothesized to be formed by a disjunction of multiple clauses (e.g.,  $U_k = \text{clause}_{k,1} \vee \text{clause}_{k,2} \vee \cdots \vee \text{clause}_{k,R_k}$ ), its constituent rule embeddings  $(\Theta_{k,1},\Theta_{k,2},\ldots,\Theta_{k,R_k})$  are learned in a sequential manner. This iterative procedure—training a rule embedding for a clause and then conceptually "covering" the samples it explains—is repeated for all  $R_k$  rule clauses intended for the disjunctive predicate  $U_k$ .

The process begins by training the first rule embedding,  $\Theta_{k,1}$ , to capture one set of conditions that satisfy  $U_k$ . The inferred value  $v_{\text{clause}_{k,1}}$  is computed, and the loss  $\mathcal{L}_{\text{clause}_{k,1}}$  (as per Eq. 6) is minimized against the partially observed  $U_{k,\text{obs}}$ .

The learning of multiple rule clauses for a predicate  $U_k$  proceeds sequentially. After an initial clause,  $\Theta_{k,1}$ , is trained to a point where it effectively explains a subset of positive instances for  $U_k$ , a hard covering step is employed. Specifically, training instances are considered "well-explained" if the output of  $\Theta_{k,1}$  (i.e.,  $v_{\text{clause}_{k,1}}$ ) for these instances exceeds a high confidence threshold (e.g., 0.99). These "well-explained" instances are then removed from the active training set. The training of  $\Theta_{k,1}$  concludes at this stage, and the subsequent rule clause  $\Theta_{k,2}$  is then trained on the remaining, unexplained instances of  $U_k$ . This iterative hard covering approach encourages further clause discovery of distinct rules that satisfy  $U_k$ .

Joint Fine-tuning of Disjunctive Rules. After the individual rule clauses for a disjunctive predicate  $U_k$  have been initialized through the sequential training and covering strategy, a joint fine-tuning phase is employed to refine these rules collectively. In this phase, the optimizer simultaneously updates all associated rule embeddings  $\{\Theta_{k,1},\ldots,\Theta_{k,R_k}\}$  for  $U_k$ . The MSE loss is computed between the combined soft-OR output  $v_{U_k}$  (obtained using Eq. 5, which aggregates the evidence from all  $R_k$  clauses) and the observed values  $U_{k,\text{obs}}$ . Given that latent predicates are, by definition, not always directly measurable, this MSE is calculated based on the small fraction of instances where the true state of the hidden predicate  $U_k$  is actually observed in the training data, which is a weakly supervised scenario:  $\mathcal{L}_{U_k,\text{finetune}} = \text{mean}((v_{U_k}[\text{mask}_k] - U_{k,\text{obs}}[\text{mask}_k])^2)$ .

The optimization details, including Adam optimizer parameters and rule embedding normalizations, are illustrated in the Appendix B.1.

#### 6 EXPERIMENTS

#### 6.1 SYNTHETIC DATA EXPERIMENTS

We use synthetic datasets to evaluate our model's ability to learn chained and disjunctive rules under partial observability (Figure 3). Each dataset is built from observable Bernoulli variables, with missing predicates defined by ground truth rules and made partially available (10%-30% observability)

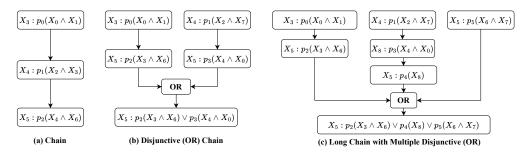


Figure 3: Example rule structures of synthetic experiments.

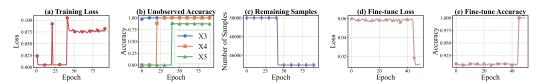


Figure 4: An example of loss and imputation accuracy during coordinate optimization (Obs. Ratio = 0.2, seed = 42). We assume the training order is  $X_3$ ,  $X_4$ ,  $X_5$ . Epochs 0–19 correspond to rule learning for  $X_3$ ; epochs 20–39 for  $X_4$ ; and epochs 40-end for  $X_5$ . Remaining samples identified how many samples are "well-explained" during the hard covering phase.

under an MCAR setting. The task is to learn rule embeddings that capture the ground truth logic, evaluated by *Rule Discovery Accuracy* and *Imputation Accuracy*. Our method is robust to MAR and MNAR mechanisms as well (Appendix D.1).

Table 1: Results for synthetic data example Figre 3(b) with an observation ratio of 0.2. Metrics are averaged over 20 random seeds on a dataset of 50,000 samples. Ground truth rules are underlined.

	Imp. Acc. (Before FT)	Imp. Acc. (After FT)	Train Loss (Before FT)	Train Loss (After FT)	Learned Rules	Rule Acc.
$X_3$	$1.00\pm0.000$	/	$0.005 \pm 0.000$	1	$X_0 \wedge X_1$	1.00
$X_4$	$0.95 \pm 0.010$	/	$0.041 \pm 0.005$	1	$X_2 \wedge X_7, X_0 \wedge X_7, X_2$	0.80
$X_5$	$0.93 \pm 0.003$	$0.96 \pm 0.002$	$0.063 \pm 0.003$	$0.067 \pm 0.001$	$ \begin{array}{c} (X_0 \wedge X_4) \vee (X_3 \wedge X_6) \\ \hline (X_3 \wedge X_4) \vee (X_3 \wedge X_6) \\ (X_0 \wedge X_1) \vee (X_0 \wedge X_4) \end{array} $	0.40

**Results and Analysis.** We analyze example (b) from Figure 3 (full results of observation ratio at 0.1 and 0.3 are in Appendix D.3). Table 1 shows that our model achieves near-perfect recovery for simple conjunctive rules  $(X_3, X_4)$  and high imputation accuracy for the complex disjunctive rule  $(X_5)$ . Figure 4(a)-(b) illustrates stable training dynamics. For  $X_5$ , the model uses sequential covering (Figure 4(c)), with "well-explained" examples reducing the remaining set. The fine-tuning (FT) phase is followed, which corrects the rule structure and boosts accuracy (Figure 4(d)-(e)). The corresponding ablation study (Table 2) confirms that fine-tuning is critical for disjunctive rules, increasing unobserved imputation accuracy for  $X_5$  from 0.87 to 1.00.

Our asynchronous coordinate descent is robust to different rule optimization orders (Table 3, Appendix Figures 7-9) and is data-efficient, recovering complex rules with as few as 4,000 samples (Appendix Figure 6). While coordinate descent requires different cycle numbers, Appendix Table 8 demonstrate minimal time and memory costs.

Convergence Analysis of Asynchronous Coordinate Descent. Exact rule-set induction reduces to the minimum-set-cover problem, which is NP-hard, so like any practical rule learner, we do not claim global optimality. Instead, we frame search as asynchronous block-coordinate descent on a smooth surrogate loss: at each step, we update a single rule embedding in closed form, which guarantees the loss never increases yet keeps each move computationally cheap. To guard against poor local minima, we (i) freeze a rule only after this rule is perfectly learned, and (ii) launch diverse initializations. Across 20 runs on synthetic datasets (Tables 11-16), this strategy delivers < 1.3%

378 379

Table 2: Ablation Study: Effect of Fine-tuning on  $X_5$  (Disjunctive Rule) Learning

Metric for $X_5$	Before Fine-tuning	After Fine-tuning
Recovered Rule Structure	$(X_0 \wedge X_2) \vee (X_0 \wedge X_4)$	$(X_3 \wedge X_6) \vee (X_0 \wedge X_4)$
$ \hline \textbf{Imputation Accuracy for } X_5 \text{ (Unobserved)} $	0.8729	1.0

386 387

384

Table 3: Impact of **rule optimization order** on learning progress. Use the example (a) of Figure 3. Note: ✓ denotes successful learning for the respective predicate.

Run 2

Run 3

393 397

399 400

401 402 403

404 405 406

407

408

413

420

421

422

428 429

430

431

Cycle Optimization Order  $[X_5, X_4, X_3]$  $[X_3, X_5, X_4]$  $[X_3, X_4, X_5]$  $X_3 \checkmark, X_4, X_5$  $X_3 \checkmark, X_4 \checkmark, X_5$  $X_3 \checkmark, X_4 \checkmark, X_5 \checkmark$ Cycle 1 Rule Accu. Imputation Accu.,  $X_3:1.00,0.005$  $X_3:1.00,0.005$  $X_3:1.00,0.005$ Train Loss  $X_4:0.87,0.074$  $X_4: 1.00, 0.004$  $X_4: 1.00, 0.004$  $X_5:0.94,0.053$  $X_5:0.94,0.035$  $X_5: 1.00, 0.003$ Optimization Order  $[X_3, X_5, X_4]$  $[X_5, X_3, X_4]$ Rule Accu. Cycle 2  $X_3 \checkmark, X_4 \checkmark, X_5$  $X_3 \checkmark, X_4 \checkmark, X_5 \checkmark$ Imputation Accu.,  $X_3:1.00,0.005$  $X_3:1.00,0.005$ Train Loss  $X_4:1.00,0.004$  $X_4: 1.00, 0.004$  $X_5:0.94,0.035$  $X_5: 1.00, 0.003$ Optimization Order  $[X_3, X_4, X_5]$ Cycle 3 Rule Accu.  $X_3 \checkmark, X_4 \checkmark, X_5 \checkmark$  $X_3:1.00,0.005$ Imputation Accu., Train Loss  $X_4:1.00,0.004$  $X_5: 1.00, 0.003$ 

Run 1

imputation performance variance, and the top-ranked learned rules consistently match ground truth rules. Such empirical evidence suggests that while global optimality is infeasible, the optimizer reliably converges to the same high-quality local solution suitable for deployment.

#### 6.2 REAL-WORLD DATA EXPERIMENTS

Metric

We validate our approach on three diverse real-world datasets, comparing it with four rule-based advanced models and a black-box MLP. For each dataset, we randomly miss some features. We then evaluated the models on their ability to impute these missing values as well as their performance on a downstream target classification task. Preprocessing details and baselines are provided in Appendices C.3 and C.2.

Our evaluation highlights NS-FCN's unique ability to handle heterogeneous data types. A key distinction is that NS-FCN directly models continuous features, whereas baseline methods are restricted to binary inputs, forcing discretization (e.g., for trestbps, binning values into < 120, 120 - 140, > 140 mmHg as low, normal, and high).

For *logical reasoning*, we used the Birds dataset (Tafjord et al., 2021) with a 90% missing ratio for two key predicates. As shown in Table 4, NS-FCN achieves perfect imputation accuracy (1.00) and, crucially, **perfectly recovers the ground truth logical rules**, highlighting its superior capability in deciphering underlying logical structures.

In medical diagnosis, we used the UCI Heart Disease (Detrano et al., 1989) and SPECT Heart (Kurgan et al., 2001) datasets, introducing 30% missingness. We also vary the observation ratio from 0.3 to 0.9 and results in Table 18 shows comparable performance with only 30% of the data observed. On the Heart Disease dataset, with its mix of continuous and categorical features, NS-FCN's direct handling of continuous values led to superior imputation (e.g., 90% accuracy for thalach) and the discovery of clinically relevant rules with numerical thresholds (e.g., age > 60, chol > 250), as shown in Table 5. On the binary SPECT dataset, we randomly miss all 22 features, thus we report the diagnosis accuracy after imputation. NS-FCN achieved a 96% F1 score for diagnosis, outperforming all baselines (Table 6). Black-Box MLP models show the most promising results yet lack interpretability. Detailed rules and LLM assessments are in Appendix Tables 19, 21, and 22.

#### 7 Conclusion

Our NS-FCN framework effectively learns interpretable rules for missing value imputation, demonstrating strong performance across a diverse range of synthetic and real-world datasets. A key

strength is its ability to seamlessly reason over heterogeneous data, handling both binary predicates (e.g., Birds) and continuous features in complex domains like medical diagnosis (SPECT, Heart Disease). It successfully handles missing data and learns hierarchical rule structures, offering significant potential for trustworthy diagnostics and transparent decision-making.

Table 4: Comparison of imputation accuracy and learned rules on the Birds dataset.

Method	Imputation Acc.	Learned Rules
LEN	0.57	$abnormal\_bird \leftarrow (ostrich \land \neg wounded) \lor (bird \land wounded)$
LEN	0.55	$can\_fly \leftarrow (bird \land \neg ostrich) \lor (\neg ostrich \land \neg wounded)$
RRL	0.53	$abnormal\_bird \leftarrow (bird \land \neg wounded) \lor (bird \land ostrich)$
KKL	0.51	$can\_fly \leftarrow (\neg ostrich \land \neg wounded) \lor (bird \land \neg ostrich)$
BRCG	0.50	$abnormal\_bird \leftarrow bird \land ostrich$
BKCG	0.47	$can\_fly \leftarrow bird \land \neg abnormal\_bird$
DR-NET	0.56	$abnormal\_bird \leftarrow (bird \land \neg ostrich \land wounded) \lor (bird \land ostrich \land \neg wounded)$
DK-NET	0.53	$can\_fly \leftarrow (bird \land \neg ostrich \land \neg abnormal\_bird) \lor (bird \land \neg ostrich \land \neg wounded)$
MLP	0.99/0.99	_
NS-FCN	1.00	$abnormal\_bird \leftarrow ostrich \lor (bird \land wounded)$
NS-FCN	1.00	$can\_fly \leftarrow bird \land \neg abnormal\_bird$

Table 5: Comparison of imputation accuracy and learned rules on the Heart Disease dataset.

Method	Imputation Acc.	Learned Rules ()
	0.65	$trestbps\_high \leftarrow (\neg st\_mild \land cp\_atypical\_angina) \lor (chol\_low \land cp\_asymptomatic)$
LENI	0.53	$chol\_high \leftarrow (sex\_female \land ca\_2) \lor (bp\_normal \land cp\_asymptomatic)$
LEN	0.62	$hr\_high \leftarrow (cp\_asymptomatic \land target) \lor (chol\_low \land ca\_l)$
	0.70	$st\_severe \leftarrow (cp\_non\_anginal \land \neg fbs\_normal) \lor (age\_old \land chol\_low)$
	0.28	$trestbps\_high \leftarrow (sex\_female \land \neg cp\_typical\_angina) \lor (exang\_yes \land \neg thal\_normal)$
RRL	0.33	$hr\_high \leftarrow (age\_middle \land sex\_male) \lor (\neg restecg\_stt\_abnormality \land slope\_upsloping)$
KKL	0.33	$thalach \leftarrow (age < 60) \land (restecg = 0)$
	0.32	$st\_severe \leftarrow (\neg exang\_yes \land \neg slope\_flat) \lor (chol\_low \land cp\_asymptomatic)$
	0.53	$trestbps\_high \leftarrow \neg age\_young \land \neg ca\_4$
BRCG	0.35	$chol\_high \leftarrow \neg age\_young \land \neg restecg\_hypertrophy$
DKCG	0.33	$hr\_high \leftarrow \neg cp\_typical\_angina \land \neg ca\_4$
	0.32	$st\_severe \leftarrow \neg age\_young \land \neg slope\_upsloping$
	0.53	$trestbps\_high \leftarrow (chol\_low \land \neg hr\_low \land \neg fbs\_high) \lor (slope\_flat \land ca\_l \land thal\_normal)$
DR-NET	0.33	$chol\_high \leftarrow sex\_male \land slope\_upsloping \land ca\_3$
DK-NEI	0.33	$hr\_high \leftarrow \neg age\_old \land \neg cp\_typical\_angina \land fbs\_high$
	0.32	$st\_severe \leftarrow hr\_high \land \neg sex\_male \land \neg fbs\_normal$
MLP	0.72 / 0.52 / 0.60 / 0.68	_
	0.86	$trestbps\_high \leftarrow (age > 60) \land (chol > 250)$
NS-FCN	0.85	$chol\_high \leftarrow (sex = 1 \land age > 55) \lor (trestbps > 150)$
NS-FCN	0.90	$hr\_high \leftarrow (trestbps > 145) \lor (age > 57 \land cp = 3)$
	0.76	$st\_severe \leftarrow (slope = 2) \land (thalach < 150)$

Table 6: Medical diagnosis after missing value imputation.

Method	Heart I	Disease	SPECT		
11204104	Accu	F1	Accu	F1	
MLP	$0.80 \pm 0.053$	$0.80 \pm 0.055$	<b>0.92</b> ±0.007	$0.90\pm0.005$	
LEN (Barbiero et al., 2022)	$0.69 \pm 0.007$	$0.80 \pm 0.000$	$0.76 \pm 0.035$	$0.85 {\pm} 0.017$	
RRL (Wang et al., 2021)	$0.78 \pm 0.002$	$0.80 \pm 0.003$	$0.90 \pm 0.005$	$0.94 \pm 0.005$	
BRCG (Dash et al., 2018)	$0.77 \pm 0.006$	$0.74 \pm 0.034$	$0.85 \pm 0.046$	$0.90 \pm 0.035$	
DR-NET (Qiao et al., 2021)	$0.85{\pm}0.005$	$0.82 {\pm} 0.005$	$0.89 {\pm} 0.025$	$0.92 \pm 0.017$	
Ours	0.91±0.009	0.91±0.009	0.92±0.009	0.96±0.009	

#### REPRODUCIBILITY STATEMENT

We have made extensive efforts to ensure the reproducibility of our results. The complete description of both synthetic dataset generation and real-world dataset preprocessing methods are illustrated in Appendix D.1 and C.2. Details of the computational setup, including hardware configuration and software environment, as well as the choice of hyper-parameters are documented in Appendix D.6 and E.3. We will release our code in the camera-ready stage to facilitate replication and further research.

#### REFERENCES

- Pietro Barbiero, Gabriele Ciravegna, Francesco Giannini, Pietro Lió, Marco Gori, and Stefano Melacci. Entropy-based logic explanations of neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6046–6054, 2022.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems*, volume 26, 2013.
- Andres Campero, Aldo Pareja, Tim Klinger, Josh Tenenbaum, and Sebastian Riedel. Logical rule induction and theory learning using neural theorem proving. *arXiv preprint arXiv:1809.02193*, 2018.
- Zhi Chen, Sarah Tan, Urszula Chajewska, Cynthia Rudin, and Rich Caruna. Missing values and imputation in healthcare data: Can interpretable machine learning help? In *Conference on Health, Inference, and Learning*, pp. 86–99. PMLR, 2023.
- Xuening Feng Paul Weng Matthieu Zimmer Dong Li Wulong Liu Jianye Hao Claire Glanois, Zhaohui Jiang. Neuro-symbolic hierarchical rule induction. In *International Conference on Machine Learning (ICML)*, pp. 7583–7615. PMLR, 2022.
- William W Cohen. Fast effective rule induction. In *Machine learning proceedings* 1995, pp. 115–123. Elsevier, 1995.
- Andrew Cropper and Rolf Morel. Learning programs by learning from failures. *Machine Learning*, 110(4):801–856, 2021.
- Sanjeeb Dash, Oktay Gunluk, and Dennis Wei. Boolean decision rules via column generation. *Advances in neural information processing systems*, 31, 2018.
- Robert Detrano, Andras Janosi, Walter Steinbrunn, Matthias Pfisterer, Johann-Jakob Schmid, Sarbjit Sandhu, Kern H Guppy, Stella Lee, and Victor Froelicher. International application of a new probability algorithm for the diagnosis of coronary artery disease. *The American journal of cardiology*, 64(5):304–310, 1989.
- Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural logic machines. *arXiv preprint arXiv:1904.11694*, 2019.
- Bradley Efron. Missing data, imputation, and the bootstrap. *Journal of the American Statistical Association*, 89(426):463–475, 1994.
- Khaled M Fouad, Mahmoud M Ismail, Ahmad Taher Azar, and Mona M Arafa. Advanced methods for missing values imputation based on similarity learning. *PeerJ Computer Science*, 7:e619, 2021.
- Kelvin Guu, John Miller, and Percy Liang. Traversing knowledge graphs in vector space. *arXiv* preprint arXiv:1506.01094, 2015.
- Pascal Hitzler and Md Kamruzzaman Sarker. Neuro-symbolic artificial intelligence: The state of the art. 2022.
  - Lukasz A Kurgan, Krzysztof J Cios, Ryszard Tadeusiewicz, Marek Ogiela, and Lucy S Goodenday. Knowledge discovery approach to automated cardiac spect diagnosis. *Artificial intelligence in medicine*, 23(2):149–169, 2001.

- Linchao Li, Bowen Du, Yonggang Wang, Lingqiao Qin, and Huachun Tan. Estimation of missing values in heterogeneous traffic data: Application of multimodal deep learning model. *Knowledge-Based Systems*, 194:105592, 2020.
  - Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. Learning entity and relation embeddings for knowledge graph completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
  - Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. Deepproblog: Neural probabilistic logic programming. *Advances in neural information processing systems*, 31, 2018.
  - Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for learning large incomplete matrices. *The Journal of Machine Learning Research*, 11:2287–2322, 2010.
  - Alessandro Oltramari, Jonathan Francis, Cory Henson, Kaixin Ma, and Ruwan Wickramarachchi. Neuro-symbolic architectures for context understanding. In *Knowledge Graphs for Explainable Artificial Intelligence: Foundations, Applications and Challenges*, pp. 143–160. IOS Press, 2020.
  - Leonardo Pellegrina and Fabio Vandin. Scalable rule lists learning with sampling. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 2352–2363, 2024.
  - Litao Qiao, Weijia Wang, and Bill Lin. Learning accurate and interpretable decision rule sets from neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 4303–4311, 2021.
  - J. Ross Quinlan. Learning logical definitions from relations. *Machine learning*, 5:239–266, 1990.
  - Md Geaur Rahman and Md Zahidul Islam. A decision tree-based missing value imputation technique for data pre-processing. In *The 9th Australasian Data Mining Conference: AusDM 2011*, pp. 41–50. Australian Computer Society Inc, 2011.
  - Md Geaur Rahman and Md Zahidul Islam. Fimus: A framework for imputing missing values using co-appearance, correlation and similarity analysis. *Knowledge-Based Systems*, 56:311–327, 2014.
  - Matthew Richardson and Pedro Domingos. Markov logic networks. *Machine learning*, 62:107–136, 2006.
  - Tim Rocktäschel and Sebastian Riedel. End-to-end differentiable proving. *Advances in neural information processing systems*, 30, 2017.
  - Hikaru Shindo, Masaaki Nishino, and Akihiro Yamamoto. Differentiable inductive logic programming for structured examples. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 5034–5041, 2021.
  - Lena Stempfle and Fredrik Johansson. Minty: Rule-based models that minimize the need for imputing features with missing values. In *International Conference on Artificial Intelligence and Statistics*, pp. 964–972. PMLR, 2024.
  - Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. Proofwriter: Generating implications, proofs, and abductive statements over natural language. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pp. 3621–3634, 2021.
  - Hind Taud and Jean-Franccois Mas. Multilayer perceptron (mlp). *Geomatic approaches for modeling land change scenarios*, pp. 451–455, 2018.
  - Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings for simple link prediction. In *International conference on machine learning*, pp. 2071–2080. PMLR, 2016.
    - Xinyu Wang, Isil Dillig, and Rishabh Singh. Synthesis of data completion scripts using finite tree automata. *Proceedings of the ACM on Programming Languages*, 1(OOPSLA):1–26, 2017.

Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. Knowledge graph embedding by trans-lating on hyperplanes. In Proceedings of the AAAI Conference on Artificial Intelligence, vol-ume 28, 2014. doi: 10.1609/aaai.v28i1.8870. URL https://doi.org/10.1609/aaai. v28i1.8870. Zhuo Wang, Wei Zhang, Ning Liu, and Jianyong Wang. Scalable rule-based representation learning for interpretable classification. Advances in Neural Information Processing Systems, 34:30479-30491, 2021. Dennis Wei, Sanjeeb Dash, Tian Gao, and Oktay Gunluk. Generalized linear rule models. In International conference on machine learning, pp. 6687–6696. PMLR, 2019. Yang Yang, Chao Yang, Boyang Li, Yinghao Fu, and Shuang Li. Neuro-symbolic temporal point processes. arXiv preprint arXiv:2406.03914, 2024. Chengqi Zhang, Yongsong Qin, Xiaofeng Zhu, Jilian Zhang, and Shichao Zhang. Clustering-based missing value imputation for data preprocessing. In 2006 4th IEEE International Conference on Industrial Informatics, pp. 1081–1086. IEEE, 2006. 

## A RELATED WORK SUPPLEMENT

**Traditional Inductive Logic Programming (ILP) Methods.** Inductive Logic Programming learns logical rules from relational data. Cohen (1995) proposed RIPPER, a fast rule induction algorithm using separate-and-conquer strategy. Quinlan (1990) developed FOIL, which generates clauses iteratively. Dash et al. (2018) introduced Boolean decision rules using column generation. Wei et al. (2019) proposed GLRM integrating decision rules into linear models. Cropper & Morel (2021) presented LFF implemented in Popper. These approaches rely on heuristics but may not guarantee optimal solutions. Pellegrina & Vandin (2024) proposed SamRuLe for near-optimal rule lists via sampling.

**Differentiable ILP Methods.** Traditional ILP models struggle with noisy data and scalability. Differentiable approaches address these issues by integrating continuous relaxation, which allows gradient descent for optimization. Shindo et al. (2021) proposed  $\partial ILP$ , which represents logic rules in a differentiable form and combines neural networks with symbolic logic. Manhaeve et al. (2018) introduced DeepProbLog, extending ProbLog with neural predicates. Neural Logic Machines (NLMs) (Dong et al., 2019) combine MLPs with logic programming to improve computational efficiency but reduce interpretability.

Broader Missing Data Imputation Methods. Missing data imputation methods range from global model-based techniques to localized and hybrid strategies, extending to deep and ensemble frameworks. At the global end, nonparametric bootstrap methods (Efron, 1994) provide bias-corrected estimates via repeated sampling, while spectral regularization approaches like SOFT-IMPUTE (Mazumder et al., 2010) solve a nuclear-norm minimization through iterative soft-thresholded SVD. Moving toward local adaptation, decision tree-based EM (DMI) (Rahman & Islam, 2011) partitions complete cases via C4.5 and imputes within each leaf, and clustering-based random imputation (CRI). (Zhang et al., 2006) applies kernel-weighted estimation in the nearest k-means cluster. Hybrid similarity learners, such as KI and its fuzzy extension FCKI (Fouad et al., 2021), refine this idea by dynamically selecting neighborhood sizes before multivariate imputation. For high-dimensional or heterogeneous data, deep architectures like MMDL (Li et al., 2020) align stacked autoencoder embeddings across modalities, and ensemble schemes like FIMUS (Rahman & Islam, 2014) combine co-appearance, correlation, and similarity in a weighted-voting framework. Despite their varied focuses—ranging from global inference to localized and multimodal learning—these methods uniformly rely on statistical patterns and do not leverage explicit logical rules to govern inter-variable relationships.

#### B MODEL SUPPLEMENT DESCRIPTION

#### **B.1** OPTIMIZATION DETAILS

Throughout all training stages, each rule embedding (or set of embeddings during joint fine-tuning) is optimized using the Adam optimizer. A crucial step following each gradient update is the normalization of the rule embeddings. This involves applying a Rectified Linear Unit (ReLU) activation to the embedding data (ensuring non-negative values, which can aid interpretability for positive predicate contributions) followed by  $L_2$  normalization of each row vector within the rule embedding matrix. This normalization helps stabilize the training process and maintains consistent magnitudes for the embedding components.

#### C DATASETS AND BASELINES

#### C.1 DATASETS

**Heart Disease.** We use the widely-cited Cleveland Clinic dataset from the UCI Heart Disease database (Detrano et al., 1989). This dataset contains 303 patient records, each with 13 features—a mix of continuous and categorical variables—such as age, cholesterol level, and resting blood pressure. The task is to predict the presence of heart disease, which is indicated by the target variable on a scale from 0 (absence) to 4 (severe). Following standard practice, we simplify this into a binary classification problem: predicting presence (values 1-4) versus absence (value 0).

**SPECT.** The SPECT (Single Proton Emission Computed Tomography) dataset presents a binary classification task to diagnose cardiac conditions (normal/abnormal) based on 22 binary patient features. The dataset describes the diagnosis of cardiac SPECT images. Each of the patients is classified into two categories: normal and abnormal. The 267 SPECT image sets (patients) database were processed to extract features that summarize the original SPECT images. As a result, 44 continuous

 feature patterns were created for each patient. The pattern was further processed to obtain 22 binary feature patterns. The CLIP3 algorithm was used to generate classification rules from these patterns (Kurgan et al., 2001). The CLIP3 algorithm generated rules that were 84.0% accurate (as compared with cardiologists' diagnoses). A key challenge in this domain is the prevalence of missing data, making it an ideal testbed for our model's imputation and rule-learning capabilities.

**Birds.** Bird's Rulebase is a well-known logic problem designed to assess an AI's ability to learn and reason with hierarchical logical rules that mimic common-sense knowledge (Tafjord et al., 2021). It has the ground truth single theory of six rules <sup>1</sup> as follows.

```
\begin{aligned} \operatorname{can\_fly}(X) \leftarrow \operatorname{bird}(X), \operatorname{not\ abnormal\_bird}(X) \\ \operatorname{bird}(X) \leftarrow \operatorname{ostrich}(X) \\ \operatorname{abnormal\_bird}(X) \leftarrow \operatorname{ostrich}(X) \\ \operatorname{not\ can\_fly}(X) \leftarrow \operatorname{ostrich}(X) \\ \operatorname{abnormal\_bird}(X) \leftarrow \operatorname{bird}(X), \operatorname{wounded}(X) \\ \operatorname{not\ can\_fly}(X) \leftarrow \operatorname{wounded}(X) \end{aligned}
```

Figure 5 further illustrates the structure of these rules.

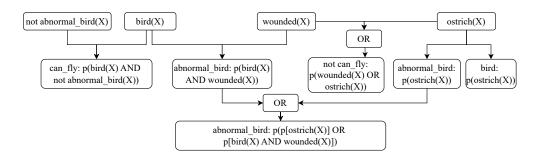


Figure 5: Ground truth rules for Bird dataset.

#### C.2 Preprocessing of Datasets

**Heart Disease.** The UCI Heart Disease dataset contains a mix of 13 continuous and categorical features with 303 samples. To create a challenging imputation task, we introduced a 30% missing ratio independently into four key continuous variables: resting blood pressure (trestbps), cholesterol (chol), maximum heart rate (thalach), and ST depression (oldpeak). For our NS-FCN framework, the task is to directly impute these missing continuous values. However, to accommodate the baseline models which only support binary inputs, we first discretized these four variables into three categorical bins based on clinical thresholds: blood pressure (< 120, 120 - 140, > 140), cholesterol (<  $200, 200 - 240, \ge 240$ ), max heart rate (<  $100, 100 - 160, \ge 160$ ), and ST depression (100, 100 - 100, 100, 100). The baselines were then tasked with imputing the correct category. Consequently, we evaluate the imputation accuracy on the discretized bins.

**SPECT Heart.** The dataset's 22 binary features were randomly masked with a 30% probability to simulate missing data. Our framework was then applied to a two-stage task: first, to impute the missing features, and second, to perform the final patient diagnosis based on the completed feature set. The diagnostic performance is compared against five baseline methods, including four rule-based approaches and an MLP.

**Birds.** Following the ground truth logical rules, we generated a dataset of 1,500 samples. To create a difficult logical reasoning challenge, we introduced a 90% missing ratio for two crucial latent predicates: can\_fly and abnormal\_bird. The task for all models was to impute these missing binary values based on the observed predicates. The imputation accuracy is compared against the same set of baselines.

<sup>&</sup>lt;sup>1</sup>https://www.doc.ic.ac.uk/ mjs/teaching/KnowledgeRep491/ExtendedLP 491-2x1.pdf, p5

#### C.3 BASELINE MODELS

- 758 760
- 761
- 762 763 764
- 765 766
- 768 769 770 771
- 772 773 774 775
- 776 777 779 780 781
- 782 783 784 785
- 786 787 788 789
- 792 793 794

796

797

798

791

799 800 801

804 805

808

809

- BRCG (Dash et al., 2018) is an integer program designed to trade classification accuracy for rule simplicity. It uses column generation to search over an exponential number of candidate clauses efficiently.
- LEN (Barbiero et al., 2022) is an end-to-end differentiable method for extracting logical explanations from neural networks using First-Order Logic.
- DR-NET (Qiao et al., 2021) is a method for learning independent logical rules in disjunctive standard form as an interpretable model for classification.
- RRL (Wang et al., 2021) learns interpretable non-fuzzy rules for data representation and classification using a novel training method called Gradient Grafting.
- Black-Box Model A two-layer Multilayer Perceptron (MLP) is used as the black-box model (Taud & Mas, 2018).

#### ADDITIONAL SYNTHETIC EXPERIMENTS RESULTS

## PERFORMANCE UNDER DIFFERENT MISSINGNESS MECHANISMS

We compare three general missingness mechanisms for dataset generation:

- MCR (Missing Completely at Random): The probability of being missing is the same for all cases, which is the missingness mechanism in other experiments on our paper.
- MAR (Missing at Random): Missingness depends on observed variables. We can indicate which observed variable to use for missingness; the default is  $X_0$ . Then, we set a higher probability of missing when the dependency variable is 1.
- MNAR (Missing Not at Random): Missingness depends on unobserved variables or the missing values themselves. Take  $X_3$  for example, we set it is more likely to be missing when  $X_3 = 1$  (positive values are harder to observe).

We show an observation ratio = 0.2 and a sample size = 50,000 as a representative case in Table 7. We run 20 random seeds. Since the seeds are different from those used in Tables 11 and 12, the results are slightly different.

Table 7: Comparison of inference accuracy and rule accuracy under different missing mechanisms.

	MCAR		MA	R	MNAI	₹
	Imputation Accu.	Rule Accu.	Imputation Accu.	Learned Rules	Imputation Accu.	Rule Accu.
$X_3$	$1.00 \pm 0.00$	1.0	$1.00 \pm 0.00$	1.0000	$1.00 \pm 0.00$	1.0000
$X_4$	$1.00 \pm 0.00$	1.0	$1.00 \pm 0.00$	1.0000	$1.00 \pm 0.00$	1.0000
$X_5$	$0.95 \pm 0.07$	0.6	$0.95 \pm 0.07$	0.6000	$0.93 \pm 0.05$	0.4000

The results show that MAR and MNAR show comparable results to MCAR, which demonstrates our method's effectiveness across the full spectrum of missing data scenarios.

#### D.2 RUNNING TIME AND MEMORY COST ANALYSIS

While coordinate descent requires different cycle numbers (Table 3), our method demonstrates efficient performance on standard CPU configurations. We conducted experiments using an Apple M4 chip with 10 cores and 16GB memory, taking observation ratio = 0.2 as an example. Results over 20 runs on setting (b) of Figure 3.

Table 8: Running time and memory cost of our model with varying sample sizes.

Sample size	2500	5000	10,000	25,000	50,000	100,000
Running time (s)	15.66±3.48	30.17±2.12	54.49±15.98	130.36±45.80	194.59±98.53	493.99±152.81
Memory cost (MB)	$64.84 \pm 10.72$	$71.92 \pm 0.82$	$78.55 \pm 1.96$	$95.81 \pm 12.13$	$126.99 \!\pm\! 26.97$	$175.64 \pm 33.93$

Overall, we observe minimal time and memory costs. Time complexity scales near-linearly with increasing sample size, while memory requirements remain modest even for large datasets. Processing 100,000 samples in under 9 minutes demonstrates strong efficiency for CPU-based execution.

#### D.3 MAIN RESULTS SUPPLEMENT OF EXAMPLE (B) OF FIGURE 3.

**Dataset Generation.** The base variables  $\{X_0, X_1, X_2, X_6, X_7\}$  are independently generated from a Bernoulli distribution, each with p=0.5. Subsequently, the values for  $\{X_3, X_4, X_5\}$  are deterministically derived using the ground truth logical rules depicted in Figure 3. Specifically, these rules are:

$$X_3 \leftarrow X_0 \wedge X_1$$

$$X_4 \leftarrow X_2 \wedge X_7$$

$$X_5 \leftarrow (X_3 \wedge X_6) \vee (X_4 \wedge X_0)$$

Finally, to introduce missing information, a portion of the values for  $X_3, X_4$ , and  $X_5$  are randomly masked. These masked variables become the targets for imputation. In our experiments, we vary the level of missingness, applying masking probabilities of 70%, 80%, and 90% to these target variables (corresponding to observation ratios of 30%, 20%, and 10%, respectively).

**Main Results.** As demonstrated in a previous case study (Table 3, which shows three runs using the same seed but different internal rule optimization orders), variations in the rule optimization sequence within a single seed can affect training efficiency. We thus show the coordinate descent training progress under a different random optimization order from Figure 4 here in Figure 7. In this run, the optimization order is  $[X_5, X_4, X_3]$  for cycle 1 and  $[X_4, X_5, X_3]$  for cycle 2. Given such different learning trajectories, our model still discovers the correct rules successfully.

Furthermore, random initialization across different seeds can lead to the discovery of varied rule sets, and occasionally, the model might converge to a local optimum. However, as the analysis of convergence before, performing multiple runs with different initializations enhances the probability of identifying the global optimal solution. Our findings indicate that the model can finds global optima several times within 20 random seeds (Tables 11 and 12).

**Learning Efficiency.** As the observation ratio decreases, the guidance signal becomes less informative, reducing both rule structure recovery and missing value imputation. We also evaluated the model's performance with a smaller training set of 10,000 samples. The results, detailed in Tables 9 and 10, demonstrate that our model maintains high accuracy for simple AND rule learning and predicate inference. Even for challenging OR rule learning, the model successfully identifies most body predicates. We further investigated the impact of dataset sample size, varying it from 1,000 to 20,000 samples. As shown in Figure 6, the most efficient setting we can recover the OR rule for  $X_5$  is to use an observation ratio of 0.1 and a dataset of 4,000 samples. For the simpler AND rules governing  $X_3$  and  $X_4$ , correct rule structures could be learned with 1,000 or even smaller samples and a 0.1 observation ratio.

Table 9: Summary of synthetic data experiment results for example (b) of the Figure 3. Each observation ratio is evaluated using 10,000 samples and results are averaged over 20 random seeds.

Obs. Ratio	Avg. Imputation Accu. (Before Fine-tune)	Avg. Imputation Accu. (After Fine-tune)	Train Loss (Before Fine-tune)	Train Loss (After Fine-tune)
0.3	$X_3: 0.91 \pm 0.014$	/	$X_3: 0.069 \pm 0.007$	/
	$X_4:0.93\pm0.013$	/	$X_4:0.054\pm0.006$	/
	$X_5: 0.87 \pm 0.003$	$X_5:0.88\pm0.002$	$X_5: 0.110 \pm 0.002$	$X_5: 0.103 \pm 0.000$
0.2	$X_3: 0.91 \pm 0.014$	/	$X_3:0.067\pm0.007$	
	$X_4:0.90\pm0.015$	/	$X_4:0.072\pm0.007$	/
	$X_5:0.86\pm0.003$	$X_5: 0.87 \pm 0.003$	$X_5: 0.116 \pm 0.002$	$X_5: 0.105 \pm 0.001$
0.1	$X_3: 0.90 \pm 0.015$	/	$X_3:0.075\pm0.007$	
	$X_4:0.91\pm0.014$	/	$X_4:0.063\pm0.006$	/
	$X_5: 0.85 \pm 0.003$	$X_5: 0.88 \pm 0.002$	$X_5: 0.124 \pm 0.002$	$X_5:0.107\pm0.001$

#### D.4 RESULTS OF EXAMPLE (A) OF FIGURE 3

**Dataset Generation.** The base variables  $\{X_0, X_1, X_2, X_6\}$  are independently generated from a Bernoulli distribution, each with p=0.5. Subsequently, the values for  $\{X_3, X_4, X_5\}$  are deterministically derived using the ground truth logical rules depicted in Figure 3. Specifically, these rules

Table 10: Summary of learned rule structures and accuracy for example (b) of Figure 3. Each observation ratio is evaluated using 10,000 samples, with results averaged over 20 random seeds. We present the top 3 learned rule structures in order of discovery accuracy. Rule accuracy indicates the percentage of 20 runs in which a rule was learned completely correctly.

Obs. Ratio	Learned Rule Structure	Rule Accu.
0.3	$X_3: X_0 \wedge X_1, X_0 \wedge X_6, X_1 \wedge X_7$	$X_3:0.65$
	$X_{\overline{4}}: \overline{X_2} \wedge \overline{X_7}, X_6 \wedge X_7, X_7$	$X_4:0.70$
	$X_5: (X_0 \wedge X_4) \vee (X_3 \wedge X_6), (\overline{X_3 \wedge X_4}) \vee (X_3 \wedge X_6), (X_0 \wedge X_1) \vee (X_0 \wedge X_4)$	$X_5:0.10$
0.2	$X_3: X_0 \wedge X_1, X_0, X_0 \wedge X_2$	$X_3:0.65$
	$X_4:\overline{X_2\wedge X_7},X_2,X_6\wedge X_7$	$X_4:0.60$
	$X_5: \underline{(X_0 \wedge X_4) \vee (X_3 \wedge X_6), (X_3 \wedge X_4)} \vee (X_3 \wedge X_6), (X_0) \vee (X_1 \wedge X_6)$	$X_5:0.10$
0.1	$X_3: X_0 \wedge X_1, X_0 \wedge X_2, X_0 \wedge X_6$	$X_3:0.60$
	$X_{\overline{4}}: X_2 \wedge X_7, X_2, X_6 \wedge X_7$	$X_4:0.65$
	$X_5: (X_3 \wedge X_4) \vee (X_3 \wedge X_6), \overline{(X_0 \wedge X_1)}) \vee (X_0 \wedge X_4), \underline{(X_0 \wedge X_4) \vee (X_3 \wedge X_6)}$	$X_5:0.10$

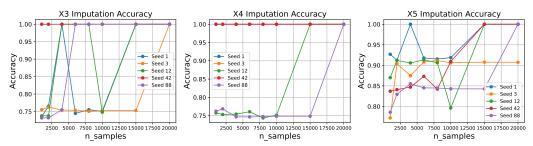


Figure 6: Imputation accuracy versus dataset sample size for Figure 3 (b). For these experiments, 10% of the data was observed (i.e., a 90% missing ratio) for predicates in  $X_3$ ,  $X_4$ , and  $X_5$ .

are:

$$X_3 \leftarrow X_0 \land X_1$$

$$X_4 \leftarrow X_2 \land X_3$$

$$X_5 \leftarrow X_4 \land X_6$$

Finally, to introduce missing information, a portion of the values for  $X_3, X_4$ , and  $X_5$  are randomly masked. These masked variables become the targets for imputation. In our experiments, we vary the level of missingness, applying masking probabilities of 70%, 80%, and 90% to these target variables (corresponding to observation ratios of 30%, 20%, and 10%, respectively).

**Main Results.** We show the coordinate descent training progress under different random optimization order. Figure 8 demonstrates the convergence in two cycles, while Figure 9 requires three cycles to complete training.

Table 11: Summary of synthetic data experiment results for example (b) of the Figure 3. Evaluated on 50,000 samples and results are averaged over 20 random seeds.

Obs. Ratio	Avg. Imputation Accu. (Before Fine-tune)	Avg. Imputation Accu. (After Fine-tune)	Train Loss (Before Fine-tune)	Train Loss (After Fine-tune)
0.3	$X_3: 0.98 \pm 0.006$ $X_4: 1.00 \pm 0.000$ $X_5: 0.94 \pm 0.003$	$X_5: 0.96 \pm 0.003$	$X_3: 0.024 \pm 0.003$ $X_4: 0.005 \pm 0.000$ $X_5: 0.054 \pm 0.002$	$X_5:0.065\pm0.001$
0.2	$X_3: 1.00 \pm 0.000$ $X_4: 0.95 \pm 0.010$ $X_5: 0.93 \pm 0.003$	$X_5: 0.96 \pm 0.002$	$X_3: 0.005 \pm 0.000$ $X_4: 0.041 \pm 0.005$ $X_5: 0.063 \pm 0.003$	$X_5:0.067\pm0.001$
0.1	$X_3: 1.00 \pm 0.000$ $X_4: 1.00 \pm 0.000$ $X_5: 0.93 \pm 0.003$	$X_5: 0.94 \pm 0.002$	$X_3: 0.005 \pm 0.000 \ X_4: 0.005 \pm 0.000 \ X_5: 0.056 \pm 0.002$	$X_5: 0.073 \pm 0.001$

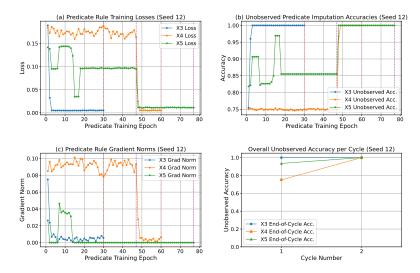


Figure 7: Training dynamics for a representative run (Obs. Ratio = 0.3) of Figure 3 (b). The optimization order:  $[X_5, X_4, X_3]$  for Cycle 1;  $[X_4, X_5, X_3]$  for Cycle 2. Subplots display: (a) training losses, (b) unobserved imputation accuracies, and (c) gradient norms for rule embeddings; (d) overall imputation accuracies each cycle. Red dashed lines indicate the conclusion of training blocks for  $X_3$  or  $X_4$  (each allocated 30 epochs when active within a cycle). Purple dashed lines delineate training phases for  $X_5$  (Rule 1, Rule 2, and Fine-tune); the epoch count for these  $X_5$  phases can vary per cycle due to the dynamic nature of the hard covering mechanism. Correct rule structures were learned for  $X_3$  by the end of Cycle 1, and for  $X_4$  and  $X_5$  by the end of Cycle 2.

Table 12: Summary of learned rule structures and accuracy for example (b) of Figure 3. Each observation ratio was evaluated using 50,000 samples, with results averaged over 20 random seeds. We present the top 3 learned rule structures in order of discovery accuracy. Rule accuracy indicates the percentage of 20 runs in which a rule was learned completely correctly.

Obs.Ratio	Learned Rule Structure	Rule Accu.
0.3	$X_3:X_0\wedge X_1,X_0\wedge X_2$	$X_3:0.90$
	$\overline{X_4:X_2}\wedge X_7$	$X_4:1.00$
	$X_5: \underline{(X_0 \wedge X_4) \vee (X_3 \wedge X_6)}, (X_3 \wedge \overline{X_4}) \vee (\overline{X_3} \wedge X_6), (X_0 \wedge X_4) \vee (X_1 \wedge X_3)$	$X_5:0.50$
0.2	$X_3: \underline{X_0 \wedge X_1}$	$X_3:1.00$
	$X_4: X_2 \wedge X_7, X_0 \wedge X_7, X_2$	$X_4:0.80$
	$X_5: \underline{(X_0 \wedge X_4) \vee (X_3 \wedge X_6)}, \overline{(X_3 \wedge X_4)} \vee (X_3 \wedge X_6), (X_0 \wedge X_1) \vee (X_0 \wedge X_4)$	$X_5:0.40$
0.1	$X_3:X_0\wedge X_1$	$X_3:1.00$
	$X_4:X_2\wedge X_7$	$X_4:1.00$
	$X_5: (X_0 \wedge X_4) \vee (X_3 \wedge X_6), (X_3 \wedge X_4) \vee (X_3 \wedge X_6), (X_0 \wedge X_1) \vee (X_0 \wedge X_4)$	$X_5:0.30$

We summarize the results for example (a) of the Figure 3 in Tables 13 and 14, which demonstrate both the effectiveness of our rule discovery approach and the precision of missing variables imputation. Our analysis reveals that learning the multi-step chain structure presents significant challenges, primarily because the algorithm uses inferred predicate values  $v^t$  from previous steps to update the current values by Eq. 3. This creates a dependency chain where suboptimal rule embeddings learned at earlier optimization steps can propagate errors to subsequent steps, potentially degrading overall performance. Despite these challenges, our model successfully identifies the correct rules in the majority of experimental runs. This robustness indicates that with multiple random initializations, the algorithm reliably converges to the optimal rule structures like the results from Figures 8 and 9, which effectively overcome the inherent difficulties of sequential dependency learning in chain-like logical structures.

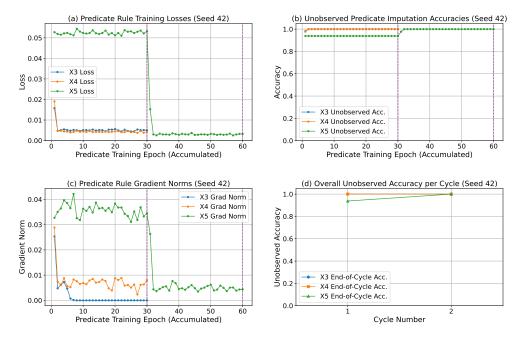


Figure 8: Training dynamics for a representative run (Observation Ratio = 0.2) of Figure 3 (a). Subplots display: (a) training losses, (b) unobserved imputation accuracies, and (c) gradient norms for rule embeddings; (d) overall imputation accuracies each cycle. Purple dashed lines indicate the conclusion of training blocks for one cycle (each allocated 30 epochs). The rule embedding optimization order: [ $X_5, X_3, X_4$ ] for Cycle 1; [ $X_5, X_4, X_3$ ] for Cycle 2. Correct rule structures were learned for  $X_3$  and  $X_4$  by the end of Cycle 1, for  $X_5$  by the end of Cycle 2. The learned rules:  $X_3 \leftarrow X_0 \wedge X_1, X_4 \leftarrow X_2 \wedge X_3, X_5 \leftarrow X_4 \wedge X_6$ .

Table 13: Summary of synthetic data experiment results for example (a) of the Figure 3. Each observation ratio is evaluated using 50,000 samples and results are averaged over 20 random seeds. No fine-tune phase since we assume no disjunctive rules.

Obs. Ratio	Avg. Imputation Accu.	Train Loss
0.3	$X_3:0.86\pm0.13$	$X_3:0.09\pm0.08$
	$X_4:0.91\pm0.06$	$X_4:0.07\pm0.04$
	$X_5:0.95\pm0.03$	$X_5:0.04\pm0.02$
0.2	$X_3:0.85\pm0.13$	$X_3:0.10\pm0.08$
	$X_4:0.90\pm0.05$	$X_4:0.08\pm0.04$
	$X_5: 0.94 \pm 0.02$	$X_5:0.04\pm0.02$
0.1	$X_3:0.82\pm0.12$	$X_3:0.10\pm0.07$
	$X_4:0.90\pm0.05$	$X_4:0.08\pm0.04$
	$X_5: 0.94 \pm 0.02$	$X_5:0.05\pm0.02$

#### D.5 RESULTS OF EXAMPLE (C) OF FIGURE 3

**Dataset Generation.** The base variables  $\{X_0, X_1, X_2, X_6, X_7\}$  are independently generated from a Bernoulli distribution, each with p=0.5. Subsequently, the values for  $\{X_3, X_4, X_5, X_8\}$  are deterministically derived using the ground truth logical rules depicted in Figure 3. Specifically, these rules are:

$$X_3 \leftarrow X_0 \wedge X_1$$

$$X_4 \leftarrow X_2 \wedge X_7$$

$$X_8 \leftarrow X_4 \wedge X_0$$

$$X_5 \leftarrow (X_3 \wedge X_6) \vee (X_8) \vee (X_6 \wedge X_7)$$

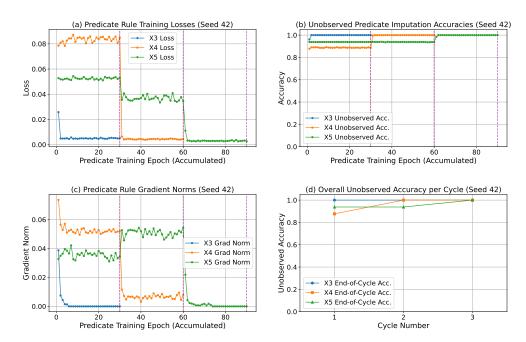


Figure 9: Training dynamics for a representative run (Observation Ratio = 0.2) of Figure 3 (a). Subplots display: (a) training losses, (b) unobserved imputation accuracies, and (c) gradient norms for rule embeddings; (d) overall imputation accuracies each cycle. Purple dashed lines indicate the conclusion of training blocks for one cycle (each allocated 30 epochs). The rule embedding optimization order:  $[X_5, X_4, X_3]$  for Cycle 1,2;  $[X_3, X_5, X_4]$  for Cycle 3. Correct rule structures were learned for  $X_3$  by the end of Cycle 1, for  $X_4$  by the end of Cycle 2, and for  $X_5$  by the end of Cycle 3. The learned rules:  $X_3 \leftarrow X_0 \wedge X_1, X_4 \leftarrow X_2 \wedge X_3, X_5 \leftarrow X_4 \wedge X_6$ .

Table 14: Summary of learned rule structures and accuracy for example (a) of Figure 3. Each observation ratio is evaluated using 50,000 samples, with results averaged over 20 random seeds. We present the top 3 learned rule structures in order of discovery accuracy. The rules that are truth rules are indicated by underline. Rule accuracy indicates the percentage of 20 runs in which a rule was learned completely correctly.

Obs. Ratio	<b>Learned Rule Structure</b>	Rule Accuracy
0.3	$X_3: X_0 \wedge X_1, X_0 \wedge X_4, X_5$	$X_3:0.60$
	$X_4: X_2 \overline{\wedge X_3, X_3} \wedge X_5, X_0 \wedge X_5$	$X_4:0.40$
	$X_5: \overline{X_4 \wedge X_6}, X_3 \wedge X_4, X_1 \wedge X_4$	$X_5:0.40$
0.2	$X_3: X_0 \wedge X_1, X_5, X_0 \wedge X_5$	$X_3:0.40$
	$X_4:\overline{X_2\wedge X_3},X_5\wedge X_3,X_5$	$X_4:0.30$
	$X_5: \overline{X_4 \wedge X_6}, X_4, X_3 \wedge X_4$	$X_5:0.20$
0.1	$X_3: X_0 \wedge X_1, X_0 \wedge X_4, X_0 \wedge X_5$	$X_3:0.40$
	$\overline{X_4:X_2\wedge X_3}, X_5, X_2\wedge X_5$	$X_4:0.40$
	$X_5: \underline{X_4 \wedge X_6}, \overline{X_4} \wedge X_2, X_3 \wedge X_6$	$X_5:0.10$

Finally, to introduce missing information, a portion of the values for  $X_3, X_4, X_8$  and  $X_5$  are randomly masked. These masked variables become the targets for imputation. In our experiments, we vary the level of missingness, applying masking probabilities of 70%, 80%, and 90% to these target variables (corresponding to observation ratios of 30%, 20%, and 10%, respectively).

**Main Results.** We summarize the results for example (c) of the Figure 3 in the Tables 15 and 16, showcasing the effectiveness of rule discovery and the precision of missing variables imputation. We have random coordinate descent training order for rule optimization.

This task is more challenging due to the chain-like structure of the disjunctive rules, particularly with three clauses for  $X_5$ , resulting in lower learning accuracy than in example (b). Nonetheless, our method achieves the highest rule discovery accuracy for the ground-truth rules while maintaining acceptable imputation accuracy. For the most difficult prediction task  $(X_5)$ , we obtain over 80% accuracy across all three observation ratios. Other predicate predictions reach  $\sim 90\%$  accuracy, including the chain-derived predicate  $X_8$ . For the learned rules in Table 16, we can find most body predicates are correct even for the complex three-clause rules governing  $X_5$ , which include the chain-derived predicate  $X_8$ . We also show the loss plot for one run in Figure 10.

Table 15: Summary of synthetic data experiment results for example (c) of the Figure 3. Each observation ratio is evaluated using 50,000 samples and results are averaged over 20 random seeds.

Obs. Ratio	Avg. Imputation Accu. (Before Fine-tune)	Avg. Imputation Accu. (After Fine-tune)	Train Loss (Before Fine-tune)	Train Loss (After Fine-tune)
0.3	$X_3: 0.86 \pm 0.12$ $X_4: 0.87 \pm 0.13$ $X_5: 0.79 \pm 0.09$ $X_8: 0.91 \pm 0.06$	$X_5: 0.84 \pm 0.08$	$X_3: 0.102 \pm 0.09$ $X_4: 0.093 \pm 0.09$ $X_5: 0.111 \pm 0.06$ $X_8: 0.060 \pm 0.04$	$X_5: 0.147 \pm 0.08$
0.2	$X_3: 0.89 \pm 0.12$ $X_4: 0.88 \pm 0.12$ $X_5: 0.78 \pm 0.08$ $X_8: 0.93 \pm 0.06$	$X_5:0.82\pm0.09$	$X_3: 0.083 \pm 0.09$ $X_4: 0.087 \pm 0.08$ $X_5: 0.119 \pm 0.06$ $X_8: 0.046 \pm 0.04$	$X_5:0.159 \pm 0.08$
0.1	$X_3: 0.89 \pm 0.12$ $X_4: 0.89 \pm 0.12$ $X_5: 0.78 \pm 0.10$ $X_8: 0.93 \pm 0.06$	$X_5: 0.80 \pm 0.11$	$X_3: 0.084 \pm 0.09$ $X_4: 0.084 \pm 0.09$ $X_5: 0.133 \pm 0.07$ $X_8: 0.048 \pm 0.04$	$X_5: 0.169 \pm 0.10$

Table 16: Summary of learned rule structures and accuracy for example (c) of Figure 3. Each observation ratio is evaluated using 50,000 samples, with results averaged over 20 random seeds. We present the top 3 learned rule structures in order of discovery accuracy. The rules that are truth rules are indicated by underline. Rule accuracy indicates the percentage of 20 runs in which a rule was learned completely correctly.

Obs. Ratio	Learned Rule Structure	Rule Accu.
0.3	$X_3: \underline{X_0 \wedge X_1}, X_0 \wedge X_2, X_1 \wedge X_2$	$X_3:0.45$
	$X_4: \underline{X_2 \wedge X_7}, X_1 \wedge X_2, X_2$	$X_4:0.5$
	$X_5: (X_3 \wedge X_6) \vee X_8 \vee \overline{(X_6 \wedge X_7)}, X_4 \vee (X_3 \wedge X_6) \vee (X_6 \wedge X_7),$	$X_5:0.15$
	$\overline{(X_3 \wedge X_4) \vee X_3} \vee (X_3 \wedge X_7)$	
	$X_8: \underline{X_4 \wedge X_0}, X_3 \wedge X_7, X_2 \wedge X_0$	$X_8:0.3$
0.2	$X_3: X_0 \wedge X_1, X_1 \wedge X_2, X_0 \wedge X_2$	$X_3:0.55$
	$X_4:\overline{X_2\wedge X_7},X_0\wedge X_2,X_1\wedge X_2$	$X_4:0.55$
	$X_5: (X_3 \wedge X_6) \vee X_8 \vee (X_6 \wedge X_7), (X_2 \wedge X_3) \vee X_4 \vee (X_6 \wedge X_7),$	$X_5:0.10$
	$X_4 \vee X_8 \vee (X_3 \wedge X_6)$	
	$X_8: \underline{X_4 \wedge X_0}, X_3 \wedge X_4, X_3 \wedge X_7$	$X_8:0.25$
0.1	$X_3: X_0 \wedge X_1, X_0, X_0 \wedge X_7$	$X_3:0.55$
	$X_4: X_2 \overline{\wedge X_7, X_0} \wedge X_2, X_1 \wedge X_7$	$X_4:0.55$
	$X_5: X_4 \vee X_8 \vee (\overline{X_6 \wedge X_7}), (X_3 \wedge X_6) \vee X_8 \vee (X_6 \wedge X_7),$	$X_5:0.10$
	$X_3 \vee X_4 \vee (X_6 \wedge X_7)$	
	$X_8: X_4 \wedge X_0, X_2 \wedge X_4, X_3 \wedge X_7$	$X_8:0.2$

#### D.6 HYPER-PARAMETERS SETTING AND COMPUTING RESOURCE

Our model operates efficiently in a CPU environment utilizing the PyTorch library. The hyperparameters are configured as follows:

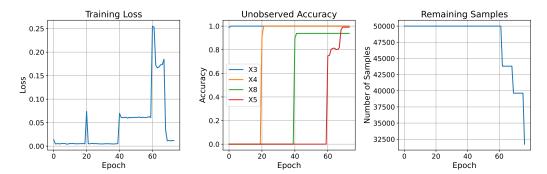


Figure 10: An example of loss and imputation accuracy during coordinate optimization (Obs. Ratio = 0.1, seed = 88, from example (c) of Figure 3). The training order is  $[X_3, X_4, X_8, X_5]$ . Epochs 0–19 correspond to rule learning for  $X_3$ ; epochs 20–39 for  $X_4$ ; epochs 40–59 for  $X_8$ , and epochs 60-end for  $X_5$ . Remaining samples identified how many samples are "well-explained" during hard covering phase. As the imputation accuracy for missing  $X_5$  is 1.00, we do not go to the fine-tune phase. The learned rules:  $X_3 \leftarrow X_0 \land X_1, X_4 \leftarrow X_2 \land X_7, X_5 \leftarrow (X_3 \land X_6) \lor X_8 \lor (X_6 \land X_7), X_8 \leftarrow X_3 \land X_4$ .

- Rule Embedding and Fine-tuning Optimizer: Adam, learning rate: 0.01.
- Temperature (Softmin/Softmax): 10.0 (for Eq. 4 and Eq. 5).
- "Well-explained" Threshold: 0.99 (for sequential hard covering in disjunctive rule learning).
- Batch Size: 64.

### E ADDITIONAL REAL WORLD DATA EXPERIMENTS RESULTS

#### E.1 SPECT

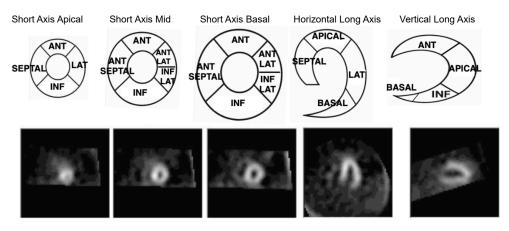


Figure 11: The five slices consists of 22 regions of interest (ROI) for SPECT Diagnosis. The slices are chosen according to the following: Three slices for short axis view-one slice near heart's apex, one in middle of the LV and one near the heart base; One slice corresponds to the center of the LV cavity for horizontal long axis view; One slice corresponds to the center of the LV cavity for vertical long axis view (Kurgan et al., 2001).

We ask for an expertise from cardiovascular surgery of a hospital to give us domain knowledge, and then we try to explain the learned rules. We select several meaningful rules to demonstrate. The domain knowledge are as follows.

R1: The anterior wall and the septum of the left ventricle are adjacent and often simultaneously affected by the Left Anterior Descending artery (LAD). If both anterior wall and septum show infarction, it strongly suggests an issue with the LAD. If both apical anterior and mid-anterior show infarction, it indicates a more extensive problem within the LAD territory, affecting both the apical and mid-portions of the anterior wall.

- R2: The apical lateral wall (typically LCX territory) and the apical inferior wall (typically RCA or LCX territory) are adjacent. Infarction in both suggests a problem in this combined region.
- R3: If both apical septal and apical septal show infarction, it indicates a more extensive problem in the LAD territory, involving ischemia in multiple myocardial segments.
- R4: If apical lateral and apical lateral show infarction, it indicates a more extensive ischemic problem in the Left Circumflex artery (LCX) territory.
- R5: The apical anterior (ANT) and apical septal (SEPTAL) regions are primarily supplied by the Left Anterior Descending artery (LAD); the apical lateral (LAT) region is primarily supplied by the Left Circumflex artery (LCX); the apical inferior (INF) region is primarily supplied by the Right Coronary Artery (RCA), but can sometimes be supplied by the LCX, depending on the coronary artery dominance pattern.

Refer to Figure 11, we can give some explanations of rules learned in Table 17 based on the domain knowledge R1 to R5. For example,

- F<sub>5</sub> ← F<sub>1</sub> ∧ F<sub>2</sub>: F<sub>1</sub> and F<sub>2</sub> are features from the first slice near the heart's apex, while F<sub>5</sub> is from the second slice at the middle of the left ventricle (LV). According to clinical knowledge R1 and R5, the anterior and septal regions are primarily supplied by the Left Anterior Descending (LAD) artery. Therefore, this rule is clinically plausible: if partial diagnosis (labeled as 1) is present in both F<sub>1</sub> and F<sub>2</sub>, it strongly suggests an LAD artery problem. Since F<sub>5</sub> is in the mid-anterior region, also supplied by the LAD, it has a high probability of being affected as well.
- $Diagnosis \leftarrow F_5 \land F_6$ : From R1, we know that the anterior wall and the septum of the left ventricle are adjacent.  $F_5$  and  $F_6$  both from middle of the LV (left ventricular), and they are adjacent. Thus, if both these adjacent mid-ventricular regions ( $F_5$  and  $F_6$ ) show signs of infarction, it significantly increases the likelihood of an overall positive diagnosis.

Table 17: Example rules learned by NS-FCN for SPECT feature imputation and diagnosis.

#### Selected Feature Imputation Rules Learned by NS-FCN

```
F_5 \leftarrow F_1 \wedge F_2: partial diagnosis of segment 1 and 2 causes the partial diagnosis of segment 5. F_6 \leftarrow F_{11} \wedge F_{19} F_{13} \leftarrow F_{22} \wedge F_{12}
```

#### **Learned Diagnosis Rule Structure**

```
Diagnosis \leftarrow (F_5 \wedge F_6) \vee (F_2 \wedge F_{11}) \vee (F_4 \wedge F_{13})
```

As detailed in Table 19, the learned rules for diagnosing cardiac abnormalities correspond closely with established domain knowledge from cardiovascular surgery experts. For instance, the model identified that infarcts in adjacent regions like  $F_1$  and  $F_2$  are indicative of an issue in the Left Anterior Descending (LAD) artery territory. Furthermore, the model learned a composite rule for the final diagnosis, logically aggregating signals from multiple infarcted regions across different coronary artery territories (LAD, LCX, RCA). This ability to synthesize information from disparate features into a coherent diagnostic rule highlights the model's capacity for complex reasoning. The clinical relevance of these rules was further validated by a Large Language Model (LLM), which confirmed their consistency with expert knowledge on ischemia propagation patterns.

Table 18: Performance on the SPECT dataset with varying observation ratios.

Obs. Ratio	Imputation Acc.	Diagnosis Acc.	Diagnosis F1
0.3	0.501	0.679	0.751
0.5	0.630	0.765	0.808
0.7	0.763	0.920	0.958
0.9	0.791	0.929	0.960

**Performance with varying missing ratios.** To assess the model's robustness under different levels of data scarcity, we evaluated its performance on the SPECT dataset while varying the observation

Table 19: Analysis of learned rules for the SPECT dataset, evaluated by human experts and LLM.

Rules	Evaluation with Human Expert Knowledge	LLM Evaluation
$F_6 \leftarrow F_1 \land F_2$	Matches R1 & R5: $F_1$ and $F_2$ are in LAD territory. Infarction in both suggests LAD issue affecting apical and mid-anterior LV.	<b>Plausible:</b> Both regions are LAD-supplied and adjacent; mid-anterior $(F_3)$ likely also affected if $F_1$ & $F_2$ show infarction. Clinically consistent.
$F_0 \leftarrow F_{11} \land F_{19}$	Related to R2 & R4: $F_{11}$ and $F_{19}$ are adjacent. Infarction implies LCX or RCA/LCX combined territory issue.	<b>Valid:</b> Matches adjacency and vascular territory logic (LCX-lateral, RCA-inferior). Supports ischemia propagation in midventricular slices.
$F_{13} \leftarrow F_{22} \land F_{12}$	Partial link to R3 & R5: Likely involves basal/apical septal $(F_{22})$ and adjacent basal regions. Indicates LAD or multi-segment ischemia.	Reasonable: Suggests ischemia spread in basal-septal regions (LAD) adjacent to basal/anterior. Fits multi-segment LAD pathology.
$\begin{array}{c} \text{Diagnosis} \leftarrow \\ (F_1 \land F_0) \lor \\ (F_2 \land F_{11}) \lor \\ (F_6 \land F_{13}) \end{array}$	Consistent with R1 & R4. Combines LAD $(F_0)$ , LCX/RCA $(F_6)$ , and adjacent mixed regions. Multiple adjacent infarct pairs increase diagnosis likelihood.	<b>Strong:</b> Logical aggregation of adjacent infarcted regions across LAD, LCX, RCA territories. Matches expert ischemia propagation patterns.

ratio from 0.3 to 0.9. As shown in Table 18, the model's accuracy remains strong and improves consistently as more data becomes available. Notably, even with only 30% of the data observed (a 70% missing ratio), the model maintains a high F1 score of 0.751, demonstrating its capability to learn meaningful diagnostic rules from highly incomplete datasets.

Table 20: Comparison of running time and memory cost across different methods.

Method	Running time (s)	Memory cost (MB)
MLP	$0.16 \pm 0.02$	$158.60 \pm 0.10$
LEN	$0.20 \pm 0.00$	$102.78 \pm 0.01$
RRL	$16.23 \pm 0.01$	$132.59 \pm 0.01$
BRCG	$2.65 \pm 0.27$	$135.11 \pm 0.08$
DR-NET	$89.01 \pm 0.06$	$45.93 \pm 0.30$
NS-FCN (Ours)	$10.34 \pm 0.30$	$61.42 \pm 1.02$

Running time and memory cost analysis. We conducted a comparative analysis of our proposed NS-FCN model against baseline methods, focusing on computational efficiency. The results in Table 20 demonstrate that NS-FCN achieves a competitive balance between performance and resource consumption. While methods like MLP and LEN offer the fastest execution times, they use higher memory costs. Our NS-FCN, though not the fastest, maintains a considerably minimal memory cost and running time.

#### E.2 HEART DISEASE

For feature imputation, as shown in Table 21, our model discovers rules with clinically relevant numerical thresholds by directly modeling continuous data. For instance, it learns to impute resting blood pressure (trestbps) based on conditions like age >60 and chol >250. Similarly, it links high cholesterol to factors like age >55 in males or very high blood pressure (trestbps >150). The learned rule for ST depression (oldpeak) combines the slope of the ST segment with a maximum heart rate threshold (thalach <150), demonstrating the model's ability to capture complex, nonlinear relationships within the data.

Beyond imputation, NS-FCN learns interpretable rules for the final diagnosis, classifying patients into low-risk or high-risk categories.

Table 22 presents several of these diagnostic rules. For example, the model learns that a combination of factors such as an upsloping ST segment (slope\_upsloping), a fixed thallium defect (thal\_fixed\_defect), and exercise-induced angina (exang\_yes) is strongly indicative of high risk. Conversely, it identifies that factors like the absence of exercise-induced angina (exang\_no) and a flat ST slope (slope\_flat) in female patients suggest a low risk of coronary artery disease. These diagnostic rules were also evaluated by an LLM and deemed "Excellent" or "Strong," underscoring their consistency with clinical practice.

Table 21: Learned rules for feature imputation on the Heart dataset, with LLM assessments.

Feature	Imputation Acc.	Learned Rule	LLM Assessment
trestbps	0.86	$\begin{array}{c} \textit{trestbps\_high} \; \leftarrow \\ (\text{age} \; > \; 60) \; \land \\ (\text{chol} > 250) \end{array}$	<b>Excellent:</b> This rule captures the well-established link between age, high cholesterol, and hypertension. Both are primary risk factors for cardiovascular disease and often co-occur.
chol	0.85	$\begin{array}{l} \textit{chol\_high} & \leftarrow \\ (\text{sex} = 1 \land \text{age} > \\ 55) \lor (\text{trestbps} > \\ 150) \end{array}$	<b>Excellent:</b> The rule correctly identifies two key risk profiles for high cholesterol: middle-aged to elderly males, and individuals with significant hypertension. This aligns perfectly with clinical understanding of metabolic syndrome.
thalach	0.90	$\begin{array}{lll} & hr\_high & \leftarrow \\ & (trestbps & > \\ & 145) \lor (age & > \\ & 57 \land cp = 3) \end{array}$	Strong: This rule insightfully links factors that limit exercise capacity to the maximum heart rate achieved. Both hypertension and severe asymptomatic coronary disease can prevent a patient from reaching a higher peak heart rate.
oldpeak	0.76	$\begin{array}{ll} \textit{st\_severe} & \leftarrow \\ (\text{slope} = 2) \land \\ (\text{thalach} < 150) \end{array}$	<b>Excellent:</b> This rule identifies a classic high-risk pattern. A downsloping ST segment is a strong positive finding, and its occurrence at a sub-maximal heart rate indicates ischemia at a low workload, a sign of severe coronary artery disease.

Table 22: Learned rules for disease prediction on the Heart dataset, with LLM assessments.

Learned Rule	LLM Assessment
$\begin{array}{ll} \textit{high\_risk} & \leftarrow \\ \textit{restecg\_stt\_abnormality} \ \land \ \textit{ca} & = \\ 3 \land \textit{oldpeak} > 1.49 \end{array}$	<b>Excellent:</b> This rule identifies a high-risk profile by combining three critical indicators of severe coronary artery disease: significant ST depression, an abnormal resting ECG, and extensive vessel blockage.
$high\_risk \leftarrow slope\_downsloping \land restecg\_normal \land trestbps > 145.68$	<b>Strong:</b> A downsloping ST segment is a powerful predictor of ischemia. Combining this with hypertension identifies patients at high risk, even if their resting ECG appears normal, highlighting the importance of stress-test indicators.
$\begin{array}{l} \textit{high\_risk} \leftarrow \textit{slope\_flat} \land \textit{oldpeak} > \\ 1.49 \land \textit{restecg\_hypertrophy} \end{array}$	<b>Excellent:</b> This rule effectively combines signs of acute ischemia (a flat ST slope with significant depression) with evidence of chronic cardiac stress (left ventricular hypertrophy). This profile is strongly indicative of advanced coronary artery disease.

# E.3 HYPER-PARAMETERS SETTING AND COMPUTING RESOURCE

#### For NS-FCN (Ours):

- Rule embedding optimizer: Adam with learning rate of 0.01.
- Fine-tine optimizer: Adam with learning rate of 0.01.
- Temperature of softmin and softmax for Eq. 4 and Eq. 5: 10.0.
- Our model can run efficiently on CPU environment with PyTorch package.

#### Baselines:

- BRCG (Dash et al., 2018), LEN (Barbiero et al., 2022), DR-NET (Qiao et al., 2021), RRL (Wang et al., 2021) are trained with the default hyperparameter settings specified in the original paper.
- All baseline models can run efficiently on CPU environment with PyTorch package.

#### F LIMITATION

While our model shows promising performance, the ethical implications, such as potential overreliance or misuse for inferring sensitive information, require careful consideration.

Despite its strengths, NS-FCN has limitations. While effective, the asynchronous coordinate gradient descent optimization can be computationally intensive. Besides, the negative predicates are not well explored (we consider negative predicates as an independent predicate from positive predicates). Furthermore, while our model can derive predicates from continuous features, the current implementation learns a single threshold per feature, which may not capture more complex relationships (e.g., intervals). Extending the framework to learn more expressive predicates from continuous data is a promising direction for future work.

#### G USE OF LLMS

In this paper, LLMs were used solely for writing polishing. The key idea, the model design, research study, and all substantive writing are completed by human authors.

In the assessment of discovered rules, we use LLM to write the evaluation of rule quality, which we have mentioned in the paper.