

RIDDLE: Rule Induction with Deep Learning

Cosimo Persia* and Ricardo Guimarães

Department of Informatics, University of Bergen, Norway

Abstract

Numerous applications rely on the efficiency of Deep Learning models to address complex classification tasks for critical decisions-making. However, we may not know how each feature in these models contributes towards the prediction. In contrast, Rule Induction algorithms provide an interpretable way to extract patterns from data, but traditional approaches suffer in terms of scalability. In this work, we bridge Deep Learning and Rule Induction and define the RIDDLE (**R**ule **I**nduction with **D**eep **L**earning) architecture. We show that RIDDLE has state-of-the-art performance in Rule Induction via an empirical evaluation.

1 Introduction

The adoption of Rule Induction algorithms supports decisions in medicine [30, 27], fault and fraud detection [33, 7], and brings benefits to chemical, oil [4], and energy industries [32], among others. These algorithms express patterns found in data in the form of associative (‘if-then’) *rules* effectively aiding users in decision-making [5, 6, 31, 22].

Rule Induction approaches have to solve hard combinatorial problems, as the symbols available for constructing the rules form a discrete space [5, 6, 20]. Hence, their scalability suffers when compared to classification methods that can rely on techniques tailored for the optimisation of differentiable functions, such as Deep Learning with gradient-based optimisation [15, 22]. Indeed, Deep Learning approaches have excelled in many tasks, including image segmentation [24] and text generation [16]. The success of Deep Learning approaches have two well-known reasons: the flexibility to handle many different forms of data and scalability

powered by technological advances. Yet, the lack of interpretability of Deep Learning models limits their application in systems that directly or indirectly influence critical decisions [6, 29].

Our main contribution is an original Deep Learning architecture dubbed RIDDLE (**R**ule **I**nduction with **D**eep **L**earning) which learns rules using a differentiable error function. This means that our approach can employ the efficient optimisation methods, while providing interpretable rules. We show with formal arguments the reasons that led the development of the RIDDLE architecture. Since it is based on the formal framework of Possibility Theory, RIDDLE builds rules equipped with certainty degrees which express the reliability of each rule. This has two main advantages over most traditional Rule Induction methods: (1) they avoid having ‘sharp’ decision boundaries and (2) the orders of the rules is irrelevant, that is, they yield rule sets instead of lists [20]. Furthermore, we show that our method has state-of-the-art performance (accuracy) on well-known datasets, especially those with uncertain or missing information.

Related work. The most prominent Rule Induction algorithms belong to the decision tree family such as RIPPER [5] that outputs binary (or crisp) rules. FURIA is an extension that generates fuzzy rules [20]. To solve the issue of scalability of Rule Induction systems in Big Data settings, Elkan et al. [15] proposed CFM-BD, a distributed system for fuzzy Rule Induction using a MapReduce paradigm. While CFM-BD has shown promising results, it still solves a search problem on a large discrete space.

We can also find works on neuro-symbolic approaches for Rule Induction. DR-Net [28] employs a simple 2-layer neural network architecture to learn rule sets. DR-Net also controls the complexity of the rules learned via a sparsity term. In contrast, RIDDLE automatically prioritises simpler

*Corresponding Author: cosimo.persia@uib.no

rules over complex ones as a consequence of weight-decay. Kusters et al. [22] define a 3-layer neural network architecture for rule induction, named R2N, that can also identify potential new terms. R2N integrates neural networks and Rule Induction with a differentiable function, but it can only learn positive DNF, a restricted class of rules where variables cannot be negated [1]. RIDDLE instead can learn any propositional formula, while also providing a measure of the ‘reliability’ for each rule.

Glanois et al. [18] propose HRI, a hierarchical approach to Rule Induction designed for Inductive Logic Programming (ILP) [25]. The language of rules differs considerably from propositional rules. Their method also relies on pre-defined rule templates that determine the types of rules that can be learned, while we do not impose such restrictions with our method.

In Section 2, we present the theoretical foundations of our contribution. Then, we introduce RIDDLE in Section 3. In Section 4, we provide an empirical comparison between RIDDLE and FURIA, a prominent fuzzy Rule Induction algorithm. We conclude and mention future steps in Section 5.

2 Preliminaries

In this section, we present relevant notions regarding propositional logic [2], Possibility Theory [11, 12], and neural networks [19].

Propositional Logic. Let \mathbf{V} be a finite set of *boolean variables*. A *literal* over \mathbf{V} , denoted with the symbol l , is either a *variable* $v \in \mathbf{V}$ or its negation, in symbols, $\neg v$. The former is also called a *positive literal*, the latter a *negative literal*. A *clause* is a disjunction (\vee) of literals. A *formula*¹ is a conjunction (\wedge) or set of clauses. A clause c can also be expressed as a *rule* r of the form $\text{ant}(r) \rightarrow \text{con}(r)$ where $\text{ant}(r)$ (the *antecedent* of r) is a conjunction of all but one negated literals in c and $\text{con}(r)$ (the *consequent* of r) is a single literal. Example 1 illustrates these concepts.

Example 1. Let $\mathbf{V} = \{v_1, v_2, t\}$. For instance, $(\neg v_1 \vee v_2 \vee t)$ is a clause, $(v_1 \wedge v_2)$ or $\{(v_1), (v_2)\}$ is a formula and $((v_1 \wedge v_2) \rightarrow t)$ is a rule.

¹In general, formulas in propositional logic can combine \neg , \vee , and \wedge arbitrarily. However, we will assume that every formula is in conjunctive normal form (CNF).

A (partial) interpretation \mathcal{I} over \mathbf{V} is a function $\mathcal{I} : \mathbf{V} \rightarrow \{\top, \perp, ?\}$ that states which variables are regarded as ‘true’, ‘false’, and ‘unknown’. \mathcal{I} *falsifies* a variable $v \in \mathbf{V}$ if $\mathcal{I}(v) = \perp$, otherwise it *satisfies* it. \mathcal{I} satisfies a negative literal $\neg v$ iff $\mathcal{I}(v)$ is equal to ‘?’ or ‘ \perp ’. \mathcal{I} satisfies a clause c if it satisfies at least one literal in c . Intuitively, \mathcal{I} satisfies a clause if there is a way of replacing ‘unknown’ values (?) with ‘true’ (\top) or ‘false’ (\perp) such that the clause is satisfied. Also, \mathcal{I} satisfies a formula ϕ , in symbols $\mathcal{I} \models \phi$, if every clause in ϕ is satisfied by \mathcal{I} . We write $\mathcal{I} \not\models \phi$ instead, if \mathcal{I} does not satisfies ϕ . We clarify these notions with the following example.

Example 2. $\mathcal{I}_1 = \{(v_1, \top), (v_2, \top), (t, \top)\}$, and $\mathcal{I}_2 = \{(v_1, \top), (v_2, ?), (t, ?)\}$ satisfy the clause in Example 1. $\mathcal{I}_3 = \{(v_1, \top), (v_2, \perp), (t, \perp)\}$ does not.

For conciseness, we abuse the notation and write \perp and \top when referring to the empty disjunction and the empty conjunction, respectively.

Possibility Theory. A *possibilistic clause* over \mathbf{V} is a pair (ϕ, α) , where ϕ is a clause over \mathbf{V} and α is a real number with finite precision in the interval $(0, 1]$, called the *valuation* of ϕ . A *possibilistic formula* \mathcal{K} is a conjunction of possibilistic clauses. Given \mathcal{K} , and a set Ω of interpretations over \mathbf{V} , we define a *possibility distribution* $\pi_{\mathcal{K}} : \Omega \rightarrow [0, 1]$ as

$$\pi_{\mathcal{K}}(\mathcal{I}) := \min_{\mathcal{I} \in \Omega} (\{1\} \cup \{1 - \alpha \mid \mathcal{I} \not\models \phi, (\phi, \alpha) \in \mathcal{K}\}).$$

The *possibility* degree $\Pi_{\mathcal{K}}(\phi)$ of ϕ , indicates how much ϕ is coherent with $\pi_{\mathcal{K}}$ and $N_{\mathcal{K}}(\phi)$ expresses the *necessity* degree of ϕ being implied by $\pi_{\mathcal{K}}$. They are defined as follows

$$\begin{aligned} \Pi_{\mathcal{K}}(\phi) &:= \sup_{\mathcal{I} \in \Omega} \{\pi(\mathcal{I}) \mid \mathcal{I} \models \phi\} \\ N_{\mathcal{K}}(\phi) &:= 1 - \Pi(\neg\phi). \end{aligned}$$

A possibility distribution $\pi_{\mathcal{K}}$ *satisfies* a possibilistic clause (ϕ, α) , written $\pi_{\mathcal{K}} \models (\phi, \alpha)$, if $N_{\mathcal{K}}(\phi) \geq \alpha$, and it satisfies a possibilistic formula \mathcal{K} if it satisfies each $(\phi, \alpha) \in \mathcal{K}$. We have that (ϕ, α) is *entailed* by \mathcal{K} , written $\mathcal{K} \models (\phi, \alpha)$, if all possibility distributions that satisfy \mathcal{K} also satisfy (ϕ, α) . We recall key properties of this theory [14].

Lemma 1 ([9]). *For every possibilistic formulas \mathcal{K} and (ϕ, α) , we have that $N_{\mathcal{K}}(\phi) = \max\{\alpha \mid \mathcal{K} \cup \{(\neg\phi, 1)\} \models (\perp, \alpha)\}$.*

P1 $\Pi(\phi_1 \vee \phi_2) = \max(\Pi(\phi_1), \Pi(\phi_2))$; and

P2 $N(\phi_1 \wedge \phi_2) = \min(N(\phi_1), N(\phi_2))$.

Neural Networks. We assume an arbitrary but fixed ordering of the variables in \mathbf{V} : $(v_1, \dots, v_n, t_1, \dots, t_m)$. We will predict the certainty degree of the variables t_j with the information provided by v_i . In this work, a *neural network* model is a function $g : [0, 1]^{2n} \rightarrow [0, 1]^m$. It takes as input a vector denoting the possibilities of each literal and it outputs the possibilities of each t_j . The function g contains parameters to be optimised by iterative updates of the backpropagation algorithm. To take advantage of the optimisation steps, we employ the function log-sum-exp: $\text{LSE}_\alpha(x_1, \dots, x_n) := \frac{1}{\alpha} \ln(e^{\alpha x_1} + \dots + e^{\alpha x_n})$, as a smooth approximation of the max (min) function when $\alpha \rightarrow \infty$ ($\alpha \rightarrow -\infty$). We write LSE_{\max} for LSE_{30} and LSE_{\min} for LSE_{-30} .

3 Introducing RIDDLE

In this section, we describe the theoretical motivations that led the development of RIDDLE. Our goal is to predict $\Pi(\neg t)$ (or $\Pi(t)$), of a target variable t , from the possibility degrees of literals $\mathbf{V} := \{v_1, \neg v_1, \dots, v_n, \neg v_n\}$. If needed, we can compute how necessary the target variable is according to our input with $N(t) = 1 - \Pi(\neg t)$ (or $N(\neg t)$).

We assume the general setting in which we have a dataset of (partial) interpretations $\mathbf{I} := \{\mathcal{I}_1, \dots, \mathcal{I}_d\}$ where some rules of the form $\phi \rightarrow t$ generally hold. We can convert the statements in $\mathcal{I} \in \mathbf{I}$ to possibilistic degrees via the method proposed by Joslyn [21] to estimate possibilities from imprecise data. We first define the set $\mathcal{I}_{\mathcal{I}} := \{\mathcal{I}' \in \mathbf{I} \mid \forall v \in \mathbf{V}, \mathcal{I}(v) = ? \text{ or } \mathcal{I}(v) = \mathcal{I}'(v)\}$ that contains precisely the interpretations in \mathbf{I} that differ only on the unknown values of \mathcal{I} . Then, from $\mathcal{I}_{\mathcal{I}}$, we count the number of interpretations that satisfy a literal l which is given by $c_{\mathcal{I}_{\mathcal{I}}}(l) := |\{\mathcal{I} \in \mathcal{I}_{\mathcal{I}} \mid \mathcal{I} \models l\}|$. Finally, the possibility associated to a literal l , according to the facts in \mathcal{I} and \mathbf{I} , is defined as $\Pi^{\mathcal{I}_{\mathcal{I}}}(l) := c_{\mathcal{I}_{\mathcal{I}}}(l) / \max(c_{\mathcal{I}_{\mathcal{I}}}(l), c_{\mathcal{I}_{\mathcal{I}}}(\neg l))$. Therefore, from \mathbf{I} we can get the set of possibility degrees for each input literal $\mathbf{D} := \{\mathbf{x}^1, \dots, \mathbf{x}^d\}$, where for any $\mathcal{I}_i \in \mathbf{I}$, $\mathbf{x}^i := (\Pi^{\mathcal{I}_i}(v), \Pi^{\mathcal{I}_i}(\neg v_1), \dots, \Pi^{\mathcal{I}_i}(v_n), \Pi^{\mathcal{I}_i}(\neg v_n))$. The number j in \mathbf{x}_j^i corresponds to the value at

the j -th position. We denote the formula associated to \mathbf{x}^i w.r.t. $\mathcal{I}_i \in \mathbf{I}$ by $\mathcal{X}^i := \{(l, 1 - \Pi^{\mathcal{I}_i}(\neg l)) \mid \Pi^{\mathcal{I}_i}(\neg l) < 1\}$. Example 3 provides a clarification.

Example 3. Let $\beta = \{\mathcal{I}_1, \mathcal{I}_2\}$ as in Example 2. $\mathbf{x}^1 = (1, 0, 1, 0, 1, 0)$, and $\mathbf{x}^2 = (1, 0, 1, 1/2, 1, 1/2)$. \mathbf{x}_4^2 is $1/2$ and $\mathcal{X}^2 := \{(v_1, 1), (v_2, 1/2), (t, 1/2)\}$.

Theoretical Motivation. We assume that the unknown formula $\mathcal{K} = \{(\phi_i \rightarrow t, \alpha_i) \mid 1 \leq i \leq k\}$ holds in the dataset of possibility degrees \mathbf{D} , and that $N_{\mathcal{K}}(\mathbf{t}) = 0$.

Lemma 2. Given $\mathcal{K} = \{(\phi_i \rightarrow t, \alpha_i) \mid 1 \leq i \leq k\}$, and $\mathcal{F} = \mathcal{K} \cup \{(l_j, \alpha_j) \mid 1 \leq j \leq s\}$, we get $\Pi_{\mathcal{F}}(\neg t) = \min\{\max\{1 - \alpha_i, \Pi_{\mathcal{F}}(\neg \phi_i) \mid 1 \leq i \leq k\}\}$.

Proof. By Lemma 1 and definition of \mathcal{K} , we have

$$\begin{aligned} N_{\mathcal{F}}(\mathbf{t}) &= \max\{\alpha \mid \mathcal{F} \cup \{(\neg t, 1) \models (\perp, \alpha)\}\} \\ &= \max\{N_{\mathcal{F}}(\phi_i \wedge (\phi_i \rightarrow t)) \mid 1 \leq i \leq k\}. \end{aligned}$$

From **P2**, we get $N_{\mathcal{F}}(\mathbf{t}) = \max\{\min(N_{\mathcal{F}}(\phi_i), \alpha_i) \mid 1 \leq i \leq k\}$. Also, for any $\mathbf{x} \in [0, 1]^n$, it holds that $(1 - \max(\mathbf{x})) = \min(1 - \mathbf{x})$, hence by the relation between possibility and necessity:

$$\begin{aligned} \Pi_{\mathcal{F}}(\neg t) &= 1 - \max\{\min(N_{\mathcal{F}}(\phi_i), \alpha_i) \mid 1 \leq i \leq k\} \\ &= \min\{\max\{1 - \alpha_i, \Pi_{\mathcal{F}}(\neg \phi_i) \mid 1 \leq i \leq k\}\}. \quad \square \end{aligned}$$

For convenience, we denote by $\psi_i := l_{i,1} \vee \dots \vee l_{i,s_i}$ the clause $\neg \phi_i$ for any rule $(\phi_i \rightarrow t, \alpha_i) \in \mathcal{K}$. By Lemma 2 and **P1**, for any formula \mathcal{X} we can compute $\Pi_{\mathcal{K} \cup \mathcal{X}}(\neg t)$ with

$$\begin{aligned} &\min\{\max\{1 - \alpha_i, \Pi_{\mathcal{K} \cup \mathcal{X}}(\psi_i) \mid 1 \leq i \leq k\}\} \quad (1) \\ &\Pi_{\mathcal{K} \cup \mathcal{X}}(\psi_i) = \max(\Pi_{\mathcal{K} \cup \mathcal{X}}(l_{1,i}), \dots, \Pi_{\mathcal{K} \cup \mathcal{X}}(l_{1,s_i})) \end{aligned}$$

By definition of \mathcal{K} , $\Pi_{\mathcal{K} \cup \mathcal{X}}(\psi_i) = \Pi_{\mathcal{X}}(\psi_i)$, so we can propagate the known uncertainty of input $\mathbf{x} \in \mathbf{D}$ to obtain the certainty degrees of the unknown target variable t with min-max operations. In practice, we do not know what rules $\phi_i \rightarrow t$ hold in \mathcal{K} and their necessity degree α_i . But, such rules constrain every possibility degree in \mathbf{D} that we can use to induce ϕ_i and α_i , with $1 \leq i \leq k$. Now, we describe RIDDLE, a novel neural network architecture for Rule Induction leveraging the uncertainty propagation properties of Possibility Theory.

Architecture. We can alternatively compute $\Pi_{\mathcal{K} \cup \mathcal{X}}(\psi_i)$ as a parametrised combination of product and maximum operators:

$$\max(w_1^{\psi_i} \Pi_{\mathcal{X}}(v_1), w_2^{\psi_i} \Pi_{\mathcal{X}}(\neg v_1), \dots, w_{2n}^{\psi_i} \Pi_{\mathcal{X}}(\neg v_n))$$

where for odd (even) $1 \leq t \leq n$, $w_t^{\psi_i} \in [0, 1]$ selects to what degree v_t ($\neg v_t$) appears in ψ_i . In matrix notation with input $\mathbf{x} \in \mathbf{D}$, this operation becomes

$$f_{\psi}(\mathbf{x}) := \mathbf{x} \star \mathbf{w}^{\psi} = \mathbf{x} \star [w_1^{\psi}, w_2^{\psi}, \dots, w_{2n}^{\psi}]^T,$$

where \star denotes the matrix dot product with the sum replaced by the LSE_{\max} operator. Lemma 3 states that for any clause ψ and input \mathbf{x} , we can find $f_{\psi} : [0, 1]^n \rightarrow [0, 1]$ such that $\Pi_{\mathcal{X}}(\psi) = f_{\psi}(\mathbf{x})$.

Lemma 3. *For any clause ψ and formula \mathcal{X} over \mathbf{V} , there is $\mathbf{w}^{\psi} \in [0, 1]^{|V|}$, s.t. $f_{\psi}(\mathbf{x}) = \Pi_{\mathcal{X}}(\psi)$.*

Proof. Let $\psi := l_1 \vee \dots \vee l_s$. By **P1**, we have $\Pi_{\mathcal{X}}(\psi) = \max(\Pi_{\mathcal{X}}(l_1), \dots, \Pi_{\mathcal{X}}(l_s))$. We can assign for odd $1 \leq t \leq n$, the value $w_t^{\psi_i} = 1$ ($w_{2t}^{\psi_i} = 1$) if v_t ($\neg v_t$) appears as a top-level literal in the disjunct ψ_i , otherwise the value 0. By definition, we get $\max(w_1^{\psi} \Pi_{\mathcal{X}}(l_1), \dots, w_{2n}^{\psi} \Pi_{\mathcal{X}}(\neg v_n)) = f_{\psi}(\mathbf{x})$. \square

As a consequence, we can approximate the computation of $\Pi_{\mathcal{K} \cup \mathcal{X}}(\neg \mathbf{t})$ in Eq. (1) with

$$\begin{aligned} & \text{LSE}_{\min}(\text{LSE}_{\max}(\boldsymbol{\beta}, \mathbf{x} \star [\mathbf{w}^{\psi_1}, \dots, \mathbf{w}^{\psi_k}])) \\ &= \text{LSE}_{\min}(\text{LSE}_{\max}(\boldsymbol{\beta}_i, f_{\psi_i}(\mathbf{x})) \mid 1 \leq i \leq k). \end{aligned} \quad (2)$$

The vector $\boldsymbol{\beta} \in [0, 1]^k$ is the parameter that approximates $1 - \alpha_i$. Therefore, the rule induction problem of finding rules in \mathcal{K} is reduced to selecting the right value of each parameter w^{ψ_i} and β_i .

We can improve this method by exploiting the associative property of the LSE_{\max} operator and compute $f_{\psi}(\mathbf{x})$ as $\text{LSE}_{\max}(f_{\psi_1}(\mathbf{x}), \dots, f_{\psi_l}(\mathbf{x}))$, where each ψ_j is a subformula of ψ (Example 4). Moreover, some rule antecedents ϕ_i, ϕ_j (with $i \neq j$) in \mathcal{K} may share subformulas, so we can decrease the number of parameters by stratifying each $f_{\psi_i}(\mathbf{x})$.

Example 4. *Given $\psi_1 := v_1 \vee \neg v_2 \vee v_3$, $\psi_2 := v_1 \vee \neg v_2 \vee v_4$, and possibilities $\mathbf{x} \in [0, 1]^8$, we can compute $f_{v_1 \vee \neg v_2}(\mathbf{x}) = \mathbf{x} \star \mathbf{w}^{v_1 \vee \neg v_2}$, $f_{v_3}(\mathbf{x}) = \mathbf{x} \star \mathbf{w}^{v_3}$, and $f_{v_4}(\mathbf{x}) = \mathbf{x} \star \mathbf{w}^{v_4}$ at first and then compute $f_{\psi_1}(\mathbf{x}) = ([f_{v_1 \vee \neg v_2}(\mathbf{x}), f_{v_3}(\mathbf{x}), f_{v_4}(\mathbf{x})] \star [1, 1, 0]^T)$. Similarly, $f_{\psi_2}(\mathbf{x}) = ([f_{v_1 \vee \neg v_2}(\mathbf{x}), f_{v_3}(\mathbf{x}), f_{v_4}(\mathbf{x})] \star [1, 0, 1]^T)$.*

For $l \geq 1$, let $\text{HL}_i^l : [0, 1]^{l_h} \rightarrow [0, 1]$ be the layer that takes as input l_h arguments and for $\mathbf{x} \in \mathbf{D}$ computes $f_{\psi_i}(\mathbf{x})$ as follows

$$\begin{aligned} \text{HL}_i^l(\mathbf{x}) &:= \text{LSE}_{\max}\{\mathbf{w}_{i,j}^l \text{HL}_j^{l-1}(\mathbf{x}) \mid 1 \leq j \leq l_h\} \\ \text{HL}_i^0(\mathbf{x}) &:= \mathbf{x}_j, \end{aligned}$$

where each $1 \leq s \leq l$, $w_{i,j}^s \in [0, 1]$. In other words, for $1 \leq s \leq l$ and $1 \leq j \leq l_h$, the function $\text{HL}_j^s(\mathbf{x})$ computes $\Pi_{\mathcal{X}}(\Psi)$ for a subformula Ψ of ψ_i , in the same way $f_{\psi_i}(\mathbf{x})$ computes $\Pi_{\mathcal{X}}(\psi_i)$. Each layer $\text{HL}_j^s : [0, 1]^u \rightarrow [0, 1]^v$, $1 \leq s \leq l$, is a function with $u, v \geq 1$ freely chosen (hyperparameters) that obey the constraint posed by the standard matrix dot product. Finally, we can define $\text{RIDDLE}(\mathbf{x})$ as

$$\text{LSE}_{\min}(\text{LSE}_{\max}(\boldsymbol{\beta}_i, \text{HL}_i^l(\mathbf{x})) \mid 1 \leq i \leq k). \quad (3)$$

Theorem 1. *Given $\mathcal{K} = \{(\phi_i \rightarrow \mathbf{t}, \alpha_i) \mid 1 \leq i \leq k\}$ with $N_{\mathcal{K}}(\mathbf{t}) = 0$, we can find a configuration of the parameters in RIDDLE s.t. for any formula $\mathcal{X} = \{(l_j, \alpha_j) \mid 1 \leq j \leq s\}$, $\text{RIDDLE}(\mathbf{x}) = \Pi_{\mathcal{K} \cup \mathcal{X}}(\neg \mathbf{t})$.*

Proof. For all literals l , by definition $N_{\mathcal{K}}(l) = 0$. Hence, for any antecedent ψ , $\Pi_{\mathcal{K} \cup \mathcal{X}}(\psi_i) = \Pi_{\mathcal{X}}(\psi)$. By Lemma 3 and associativity of max, we can set the values of the parameters in $\text{HL}_i^l(\mathbf{x})$ so that it computes $\Pi_{\mathcal{X}}(\phi_i)$. By Eqs. (2) and (3), we get that $\text{RIDDLE}(\mathbf{x})$ computes $\Pi_{\mathcal{K} \cup \mathcal{X}}(\neg \mathbf{t})$ as in Eq. (1). \square

Multi-class tasks are modelled with many output nodes or with multiple RIDDLE instances in parallel. Theorem 1 shows the generality of our approach but it relies on Lemma 3 which requires parameters in $[0, 1]$. Thus, after updating the parameters with SGD [3], we replace negative values with 0 and values greater than 1 with 1.

Rule Extraction and Injection. When the optimisation procedure terminates, we can track the literals that are used to discriminate the target possibility value by inspecting the value of each parameter. Indeed, every parameter lies in the interval $[0, 1]$ and by the semantics given to their values, we can apply the argument in Lemma 3, to extract the literals included in the clause of each layer HL_i^l . We observed that the parameters always collapse to either 0 or 1 after a sufficient number of updates (Section 4). Also, the introduction of hidden layers can be considered a way of having predicate invention as in ILP settings [25].

We can also manually inject rules of the form $\phi \rightarrow t$ to a RIDDLE instance before or after training. Due to Lemma 3, we just need to append to the operation LSE_{\min} in Eq. (3) a layer $\text{HL}' : [0, 1]^{2n} \rightarrow [0, 1]$, corresponding to the function $f_{-\phi}$.

4 Experimental Results

We implemented the RIDDLE model in Python 3.9 that is fully integrated in the PyTorch [26] ecosystem. The implementation is freely available at the following link: <https://git.app.uib.no/Cosimo.Persia/riddle>. The gradient of the model parameters are computed with PyTorch’s automatic differentiation package and after the update, they are ‘clamped’ to the range $[0, 1]$ to preserve the correctness of the model. We conduct the experiments on an Ubuntu 18.04.5 LTS server with i9-7900X CPU at 3.30GHz, 32 physical cores, 8 GPUs NVIDIA A100 with 80GB, and 32GB RAM.

Test settings. Often, the features in the considered datasets include a mixture of nominal, continuous, and integer fields. Using feature discretization, we divide continuous or integer values in 8 bins such that all bins for each feature have the same number of points. Each bin will be associated with a new variable that it is going to be set to ‘true’ if the value of the original value belongs to the respective bin. Missing values assign the value ‘unknown’ to all related new binarised variables. In this way, we can generate a set of interpretations $D := \{\mathcal{I}_1, \dots, \mathcal{I}_d\}$. From D , we can derive the dataset $D := \{\mathbf{x}^i \mid \mathcal{I}_i \in D\}$ as explained in Section 3. D expresses the possibility values of variables and their negation for each interpretation in D . The first column in Table 1, shows the datasets that we considered for the benchmark. These are freely available at UCI machine learning repository [8]. Briefly, with ‘breast cancer’, ‘hepatitis’, ‘horse’, ‘hypothyroid’, ‘lymphography’, and ‘primary tumor’, the model should predict the type of disease or the patient’s survival. With ‘auto’, ‘credit’, ‘chess’, ‘glass’, ‘mushroom’, and ‘wine’, the model should classify the specific type of object under scrutiny; for example, whether a mushroom is edible. More details about these can be found at the UCI website². Most datasets have a substantial

amount of missing values.

Evaluation. RIDDLE minimises the MSE (mean squared error) during training (see column ‘MSE’ in Table 1). Compared with arbitrary linear/ReLU feedforward deep network architectures, RIDDLE performs slightly better (in addition to being explainable). Therefore, we will focus our comparison on the accuracy of the FURIA algorithm [20] that represents the state-of-the-art in propositional Rule Induction with confidence values. We use the FURIA implementation available on Weka [17], and compare it with RIDDLE on classification tasks using the standard definition of accuracy. To use the trained RIDDLE model for classification, we look at RIDDLE output $(\Pi(-\mathbf{t}_1), \dots, \Pi(-\mathbf{t}_m))$ and if $\Pi(-\mathbf{t}_i) \leq 0.4$, then the variable \mathbf{t}_i is preferred over its negation and we assume that the variable \mathbf{t}_i is predicted to be true ($N(\mathbf{t}) = 1 - \Pi(-\mathbf{t})$). We carried our tests on the same benchmark datasets used by the aforementioned rule induction system.

Model selection. We split each dataset in 70%, 10%, and 20% for training, validation and test, respectively. Finding the best combinations of hyperparameters (number of layers, nodes per layers, learning rate etc.) can be a time-consuming task. But, we noticed that, in general, RIDDLE performs well even with small networks and extremely well with deeper configurations. We adjust the hyperparameters on a grid search fashion using Tune [23]. The number of hidden layers varies from 1 to 10 and the number of nodes per layer from 2 to 500, the batch size from 8 to 64, and the learning rate from 0.1 to 0.001. The final models’ sizes correlate positively with the number of variables given as input and with the number of rules that hold in the dataset. On average, the resulting network has 6 layers with 50 nodes. Each model is trained with a batch of size 16 over 100 epochs with a learning rate of 0.01, and with early-stopping. That is, we stop the training routine if the validation loss has not decreased by more than 0.01 for 10 subsequent epochs. Moreover, we fixed a weight-decay factor of 0.001. This means that the gradients used for updating the parameters are summed to the constant value 0.001.

Results. Table 1 shows the results of our experiments. The column ‘Inst.’ shows the number of instances, #V shows the number of variables in the original dataset and #DV is the num-

²<http://archive.ics.uci.edu/ml>

Dataset	Inst.	#V	#DV	MSE	Accuracy		Model complexity			
					RIDDLE	FURIA	RIDDLE		FURIA	
							Size	Count	Size	Count
Anneal	798	39	286	1e-4	97	97	2.9	17.1	3.2	21.7
Audiology	226	71	708	4e-3	95	91	1.7	19.0	2.1	19.3
Auto	205	26	354	3e-4	98	85	2	3.4	2	3.3
Credit-A	690	15	158	2e-1	87	89	3.1	6.8	4.1	7.3
Credit-G	1000	21	166	2e-2	74	72	2.2	10.3	3.1	10.2
Breast cancer	699	21	160	3e-2	91	90	3.3	5.1	3.6	9.9
Chess	3196	37	146	1e-2	91	99	4.6	15.3	4.7	28.3
Glass	214	10	72	4e-3	94	68	1.8	5.9	2.0	16.7
Hepatitis	155	20	122	4e-4	86	75	3.1	3.4	2.2	4.1
Horse	368	28	252	5e-3	83	85	1.8	5.9	2.7	5.0
Hypothyroid	3163	30	140	8e-3	95	95	4.8	11.3	6.3	18.1
Lymphography	148	19	118	9e-4	86	86	2.1	5.4	2.2	5.3
Mushroom	8124	19	234	1e-4	97	98	2.8	5.3	2.4	6.0
Primary tumor	339	18	60	6e-5	94	75	1.9	3.6	1.9	4.0
Wine	178	14	112	1e-4	96	90	3.4	8.2	3.4	8.1

Table 1: Accuracy and model complexity of RIDDLE and FURIA compared with different datasets

ber variables after discretization and binarization. ‘MSE’ is the loss of RIDDLE on the test data. Often, RIDDLE converges to a minimum after only 30 epochs and soon after the early-stopping routine stops the training. The average training time with the largest datasets (mushrooms, chess, and hypothyroid) is 2 seconds for RIDDLE and 1 for FURIA. The ‘Accuracy’ columns in Table 1 show that in most cases RIDDLE generalises from training data better than FURIA. However, FURIA performs better in complete information scenarios as in the ‘chess’ dataset. But, even in such case RIDDLE outputs a simpler model. In contrast, RIDDLE yields better accuracy in datasets with more missing data, such as ‘glass’, or ‘hepatitis’.

The rightmost part of Table 1 reports the size and the number of induced clauses found by each model per dataset. For instance, rules found by RIDDLE in the ‘hepatitis’ are

$$((1.9 \leq \text{bilirubin}) \text{ and } \neg \text{has_ascites} \rightarrow \text{dies}, 0.8)$$

$$((3.7 \leq \text{albumin}) \text{ and } \text{firm_liver} \rightarrow \text{lives}, 1).$$

Moreover, RIDDLE has better model complexity with fewer and shorter rules, even when RIDDLE performed worse in terms of accuracy. This suggests a bias towards simpler models. When RIDDLE encodes known rules with additional layers (end of Section 3), the performance improves.

Dataset	Min	Max	Median	S.D.
Anneal	93	98	96	1.18
Audiology	94	96	95	0.63
Auto	95	99	98	2.48
Credit-A	66	92	86	4.12
Credit-G	68	82	76	5.65
Breast cancer	85	96	91	4.52
Chess	83	98	88	6.04
Glass	81	99	88	4.54
Hepatitis	71	89	86	6.83
Horse	78	86	84	5.71
Hypothyroid	91	98	94	4.37
Lymphography	83	88	87	2.91
Mushroom	94	98	97	1.17
Primary tumor	85	99	93	5.13
Wine	84	98	95	3.59

Table 2: Additional statistics from the empirical evaluation of RIDDLE obtained by running 40 training instances.

Table 2 shows the statistical results concerning the experiments. For each dataset, we stored the minimal and maximal accuracy. Then, we computed the median and standard deviation (S.D.). From Table 2 we can conclude that RIDDLE has consistent performance in most of the datasets considered.

We remark that another advantage of RIDDLE is that the values provided with the rules have a clear meaning, in terms of necessity. Meanwhile, fuzzy approaches such as FURIA provide a weaker foundation for the interpretation of the values associated with the rules. Additionally, the necessity values also distinguish RIDDLE from approaches such as decision trees which, usually, do not provide a measure of a rule’s reliability.

5 Conclusion

In this work, we introduced RIDDLE: a novel deep learning architecture specialised in performing Rule Induction in the presence of incomplete or uncertain data. RIDDLE is a white-box model as its trained weights have a clear meaning concerning the decisions that the model takes while performing inference on the input. These weights can be translated into propositionally complete rules that are simpler than the rules found by state-of-the-art algorithms. In addition, each rule is associated with certainty degree expressing the confidence of the model about the induced rule. Not only that, RIDDLE can also seamlessly incorporate background knowledge via rule injection. Thus, RIDDLE provides an efficient, flexible, and interpretable solution for Rule Induction.

Future work. The next step is to optimise the matrix computation in RIDDLE’s implementation and speed-up both training and inference time. Also, we will evaluate the effect of different methods of drawing possibilities distributions from imprecise data [10, 13] on accuracy.

Acknowledgements

Part of this work has been done in the context of CEDAS (Center for Data Science, University of Bergen, Norway). The first author is supported by the “L. Meltzers Høyskolefond” utdeling 2022.

The second author is supported by the ERC project “Lossy Preprocessing” (LOPRE), grant number 819416, led by Prof Saket Saurabh.

References

- [1] D. Angluin. Queries and concept learning. *Machine Learning*, 2(4):319–342, Apr 1988. doi: 10.1023/A:1022821128753. URL <https://doi.org/10.1023/A:1022821128753>.
- [2] J. L. J. L. Bell and M. Machover. *A course in mathematical logic / by J. L. Bell and M. Machover*. North-Holland, Amsterdam, 1977. ISBN 0720428440. doi: <https://doi.org/10.1007/978-1-4757-4385-2>.
- [3] L. Bottou. Stochastic learning. In O. Bousquet and U. von Luxburg, editors, *Advanced Lectures on Machine Learning*, Lecture Notes in Artificial Intelligence, LNAI 3176, pages 146–168. Springer Verlag, Berlin, 2004. doi: 10.1007/978-3-540-28650-9_7.
- [4] I. Bratko. Applications of machine learning: Towards knowledge synthesis. *New Generation Computing*, 11(3):343–360, Sep 1993. ISSN 1882-7055. doi: 10.1007/BF03037182.
- [5] W. W. Cohen. Fast effective rule induction. In A. Prieditis and S. Russell, editors, *Machine Learning, Proceedings of the Twelfth International Conference on Machine Learning, Tahoe City, California, USA, July 9-12, 1995*, pages 115–123. Morgan Kaufmann, 1995. doi: 10.1016/b978-1-55860-377-6.50023-2.
- [6] S. Dash, O. Günlük, and D. Wei. Boolean decision rules via column generation. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 4660–4670, 2018. doi: <https://doi.org/10.48550/arXiv.1805.09901>.
- [7] J. A. Dhanraj, M. Prabhakar, C. P. Ramaian, M. Subramaniam, J. M. Solomon, and N. Vinayagam. Increasing the Wind Energy Production by Identifying the State of Wind Turbine Blade. In *Lecture Notes*

- in *Mechanical Engineering*, pages 139–148. Springer Nature Singapore, 2022. doi: 10.1007/978-981-16-7909-4_13.
- [8] D. Dua and C. Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- [9] D. Dubois and H. Prade. Resolution principles in possibilistic logic. *Int. J. Approx. Reasoning*, 4(1):1–21, 1990. doi: [https://doi.org/10.1016/0888-613X\(90\)90006-N](https://doi.org/10.1016/0888-613X(90)90006-N).
- [10] D. Dubois and H. Prade. When upper probabilities are possibility measures. *Fuzzy Sets and Systems*, 49:65–74, 1992. doi: [https://doi.org/10.1016/0165-0114\(92\)90110-P](https://doi.org/10.1016/0165-0114(92)90110-P).
- [11] D. Dubois and H. Prade. *Possibility Theory*, pages 6927–6939. Springer New York, New York, NY, 2009. ISBN 978-0-387-30440-3. doi: 10.1007/978-0-387-30440-3_413.
- [12] D. Dubois and H. Prade. Possibility theory and its applications: Where do we stand? In *Springer Handbook of Computational Intelligence*, pages 31–60. 2015. doi: 10.1007/978-3-662-43505-2_3.
- [13] D. Dubois and H. Prade. Practical methods for constructing possibility distributions. *International Journal of Intelligent Systems*, 31(3): 215–239, 2016. doi: <https://doi.org/10.1002/int.21782>.
- [14] D. Dubois, J. Lang, and H. Prade. *Possibilistic Logic*, page 439–513. Oxford University Press, Inc., 1994.
- [15] M. Elkano, J. A. Sanz, E. Barrenechea, H. Bustince, and M. Galar. CFM-BD: A distributed rule induction algorithm for building compact fuzzy models in big data classification problems. *IEEE Transactions on Fuzzy Systems*, 28(1):163–177, jan 2020. doi: 10.1109/tfuzz.2019.2900856.
- [16] N. Fatima, A. S. Imran, Z. Kastrati, S. M. Daudpota, and A. Soomro. A systematic literature review on text generation using deep neural network models. *IEEE Access*, 10: 53490–53503, 2022. doi: 10.1109/access.2022.3174108.
- [17] E. Frank, M. A. Hall, G. Holmes, R. Kirkby, B. Pfahringer, and I. H. Witten. *Weka: A machine learning workbench for data mining.*, pages 1305–1314. Springer, Berlin, 2005. URL <http://researchcommons.waikato.ac.nz/handle/10289/1497>.
- [18] C. Glanois, Z. Jiang, X. Feng, P. Weng, M. Zimmer, D. Li, W. Liu, and J. Hao. Neuro-symbolic hierarchical rule induction. In *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 7583–7615. PMLR, 2022. doi: arxiv-2112.13418.
- [19] I. J. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, Cambridge, MA, USA, 2016. doi: <https://doi.org/10.1007/s10710-017-9314-z>.
- [20] J. Hühn and E. Hüllermeier. Furia: an algorithm for unordered fuzzy rule induction. *Data Mining and Knowledge Discovery*, 19(3): 293–319, Dec 2009. ISSN 1573-756X. doi: 10.1007/s10618-009-0131-8.
- [21] C. Joslyn. Towards an empirical semantics of possibility through maximum uncertainty. 1991.
- [22] R. Kusters, Y. Kim, M. Collery, C. de Sainte Marie, and S. Gupta. Differentiable rule induction with learned relational features. Jan. 2022. doi: arXiv:2201.06515.
- [23] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica. Tune: A research platform for distributed model selection and training. 2018. doi: <https://arxiv.org/abs/1807.05118>.
- [24] S. Minaee, Y. Y. Boykov, F. Porikli, A. J. Plaza, N. Kehtarnavaz, and D. Terzopoulos. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. doi: 10.1109/tpami.2021.3059968.
- [25] S. Muggleton, L. De Raedt, D. Poole, I. Bratko, P. Flach, K. Inoue, and A. Srinivasan. Ip turns 20. *Machine Learning*, 86

- (1):3–23, Jan 2012. ISSN 1573-0565. doi: 10.1007/s10994-011-5259-2.
- [26] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshain, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [27] V. Podgorelec, P. Kokol, B. Stiglic, and I. Rozman. Decision trees: An overview and their use in medicine. *Journal of Medical Systems*, 26(5):445–463, Oct 2002. ISSN 1573-689X. doi: 10.1023/A:1016409317640.
- [28] L. Qiao, W. Wang, and B. Lin. Learning accurate and interpretable decision rule sets from neural networks. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 4303–4311. AAAI Press, 2021. doi: <https://arxiv.org/abs/2103.02826>.
- [29] C. Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.*, 1(5):206–215, 2019. doi: 10.1038/s42256-019-0048-x.
- [30] G. Scala, A. Federico, V. Fortino, D. Greco, and B. Majello. Knowledge generation with rule induction in cancer omics. *International Journal of Molecular Sciences*, 21(1):18, dec 2019. doi: 10.3390/ijms21010018.
- [31] J. Vreeken, M. van Leeuwen, and A. Siebes. Krimp: mining itemsets that compress. *Data Mining and Knowledge Discovery*, 23(1):169–214, Jul 2011. ISSN 1573-756X. doi: 10.1007/s10618-010-0202-x.
- [32] L. WJ. A rule-based process control method with feedback. *Advances in Instrumentation* 41, 169–175, 1987.
- [33] P. Xu, Z. Ding, and M. Pan. A hybrid interpretable credit card users default prediction model based on RIPPER. *Concurrency and Computation: Practice and Experience*, 30(23):e4445, feb 2018. doi: 10.1002/cpe.4445.