

---

# LOGIC EXPLAINED NETWORKS

---

Gabriele Ciravegna<sup>\*1,2</sup> , Pietro Barbiero<sup>\*3</sup> , Francesco Giannini<sup>\*2</sup>,  
 Marco Gori<sup>2,4</sup> , Pietro Lió<sup>3</sup>, Marco Maggini<sup>2</sup>, Stefano Melacci<sup>2</sup> 

<sup>1</sup>Università di Firenze (Italy), <sup>2</sup>Università di Siena (Italy),

<sup>3</sup>University of Cambridge (UK), <sup>4</sup>Université Côte d'Azur (France)

gabriele.ciravegna@unifi.it, pb737@cam.ac.uk, francesco.giannini@unisi.it

marco.gori@unisi.it, pl219@cam.ac.uk, marco.maggini@unisi.it, mela@diism.unisi.it

## ABSTRACT

The large and still increasing popularity of deep learning clashes with a major limit of neural network architectures, that consists in their lack of capability in providing human-understandable motivations of their decisions. In situations in which the machine is expected to support the decision of human experts, providing a comprehensible explanation is a feature of crucial importance. The language used to communicate the explanations must be formal enough to be implementable in a machine and friendly enough to be understandable by a wide audience. In this paper, we propose a general approach to Explainable Artificial Intelligence in the case of neural architectures, showing how a mindful design of the networks leads to a family of interpretable deep learning models called Logic Explained Networks (LENs). LENs only require their inputs to be human-understandable predicates, and they provide explanations in terms of simple First-Order Logic (FOL) formulas involving such predicates. LENs are general enough to cover a large number of scenarios. Amongst them, we consider the case in which LENs are directly used as special classifiers with the capability of being explainable, or when they act as additional networks with the role of creating the conditions for making a black-box classifier explainable by FOL formulas. Despite supervised learning problems are mostly emphasized, we also show that LENs can learn and provide explanations in unsupervised learning settings. Experimental results on several datasets and tasks show that LENs may yield better classifications than established white-box models, such as decision trees and Bayesian rule lists, while providing more compact and meaningful explanations.

## 1 Introduction

The application of deep neural networks in safety-critical domains has been strongly limited [Chander et al., 2018], since neural networks are generally considered black-boxes whose decision processes are opaque or too complex to be understood by users. Employing black-box<sup>2</sup> models may be unacceptable in contexts such as industry, medicine or courts, where the potential economical or ethical repercussions are calling for lawmakers to discourage from a reckless application of non-interpretable models [EUGDPR, 2017, Law, 10, Goddard, 2017, Gunning, 2017]. As a consequence, research in Explainable Artificial Intelligence (XAI) has become strategic and has been massively encouraged, leading to the development of a variety of techniques that aim at explaining black-box models [Das and Rad, 2020, Brundage et al., 2020] or at designing interpretable models [Carvalho et al., 2019, Rudin, 2019].

The notions of *interpretability* and *explainability* were historically used as synonyms. However the distinction between interpretable models and models providing explanations has now become more evident, as recently discussed by different authors [Gilpin et al., 2018, Lipton, 2018, Marcinkevičs and Vogt, 2020]. Even if there are no common accepted formal definitions, a model is considered interpretable when its decision process is generally

---

\*Equal contribution

<sup>2</sup>In the context of this paper, a black-box classifier is any classifier that cannot provide human understandable explanations about its decision.

transparent and can be understood directly by its structure and parameters, such as *linear models* or *decision trees*. On the other hand, the way an existing (black-box) model makes predictions can be explained by a surrogate interpretable model or by means of techniques providing intelligible descriptions of the model behaviour by e.g. formal rules, saliency maps, question-answering. In some contexts, the use of a black-box model may be unnecessary or even not preferable [Doshi-Velez and Kim, 2017, Doshi-Velez and Kim, 2018, Ahmad et al., 2018, Rudin, 2019, Samek et al., 2020, Rudin et al., 2021]. For instance, a proof of concept, a prototype, or the solution to a simple classification problem can be easily based on standard interpretable by design AI solutions [Breiman et al., 1984, Schmidt and Lipson, 2009, Letham et al., 2015, Cranmer et al., 2019, Molnar, 2020]. However, interpretable models may generally miss to capture complex relationships among data. Hence, in order to achieve state-of-the-art performance in more challenging problems, it may be necessary to leverage black-box models [Battaglia et al., 2018, Devlin et al., 2018, Dosovitskiy et al., 2020, Xie et al., 2020] that, in turn, may require an additional explanatory model to gain the trust of the user.

In the literature, there is a wide consensus on the necessity of having an explanation for machine learning models that are employed in safety-critical domains, while there is not even agreement on what an *explanation* actually is, nor if it could be formally defined [Doshi-Velez and Kim, 2017, Lipton, 2018]. As observed by Srinivasan and Chander, explanations should serve cognitive-behavioral purposes such as engendering trust, aiding bias identification, or taking actions/decisions [Srinivasan and Chander, 2020]. An explanation may help humans understand the black-box, may allow for a deeper human-machine interaction [Koh et al., 2020], and may lead to more trustworthy fully automated tasks. All of this is possible as long as the explanations are useful from a human perspective. In a nutshell, an explanation is an answer to a “why” question and what makes an explanation good or bad depends on “the degree to which a human can understand the cause of a decision” [Miller, 2019]. The goodness of an explanation is intimately connected to how humans collect evidences and eventually make decisions. In Herbert Simon’s words we may say that a good explanation is *satisficing* when “it either gives an optimal description of a simplified version of the black-box (e.g. a surrogate model) or a satisfactory description for the black-box itself” [Simon, 1979]. However, the notion itself of explanation generally depends on both the application domain and whom it is aimed at [Carvalho et al., 2019].

The need for human-understandable explanations is one of the main reasons why concept-based models are receiving ever-growing consideration, as they provide explanations in terms of human-understandable symbols (the *concepts*) rather than raw features such as pixels or characters [Kim et al., 2018, Ghorbani et al., 2019, Koh et al., 2020]. As a consequence, they seem more suitable to serve many strategic human purposes such as decision making tasks. For instance, a concept-based explanation may describe a high-level category through its attributes as in “a *human* has *hands* and a *head*”. While concept ranking is a common feature of concept-based techniques, there are very few approaches formulating hypotheses on how black-boxes combine concepts to arrive to a decision and even less provide synthetic explanations whose validity can be quantitatively assessed [Das and Rad, 2020].

A possible solution to provide human-understandable explanations is to rely on a formal language that is very expressive, closely related to reasoning, and somewhat related to natural language expressions, such as First-Order Logic (FOL). A FOL explanation can be considered a special kind of a concept-based explanation, where the description is given in terms of logic predicates, connectives and quantifiers, such as “ $\forall x : is\_human(x) \rightarrow has\_hands(x) \wedge has\_head(x)$ ”, that reads “being human implies having hands and head”. However, FOL formulas can generally express much more complex relationships among the concepts involved in a certain explanation. Compared to other concept-based techniques, logic-based explanations provide many key advantages, that we briefly describe in what follows. An explanation reported in FOL is a rigorous and unambiguous statement (clarity). This formal clarity may serve cognitive-behavioral purposes such as engendering trust, aiding bias identification, or taking actions/decisions. For instance, dropping quantifiers and variables for simplicity, the formula “ $snow \wedge tree \leftrightarrow wolf$ ” may easily outline the presence of a bias in the collection of training data. Different logic-based explanations can be combined to describe groups of observations or global phenomena (modularity). For instance, for an image showing only the face of a person, an explanation could be “ $(nose \wedge lips) \rightarrow human$ ”, while for another image showing a person from behind a valid explanation could be “ $(feet \wedge hair \wedge ears) \rightarrow human$ ”. The two local explanations can be combined into “ $(nose \wedge lips) \vee (feet \wedge hair \wedge ears) \rightarrow human$ ”. The quality of logic-based explanations can be quantitatively measured to check their correctness and completeness (measurability). For instance, once the explanation “ $(nose \wedge lips) \vee (feet \wedge hair \wedge ears)$ ” is extracted for the class *human*, this logic formula can be applied on a test set to check its generality in terms of quantitative metrics like accuracy, fidelity and consistency. Further, FOL-based explanations can be rewritten in different equivalent forms such as in *Disjunctive Normal Form* (DNF) and *Conjunctive Normal Form* (CNF) (versatility). Finally, techniques such as the Quine–McCluskey algorithm can be used to compact and simplify logic explanations [McCull, 1878, Quine, 1952, McCluskey, 1956] (simplifiability). As a toy example, consider the explanation “ $(person \wedge nose) \vee (\neg person \wedge nose)$ ”, that can be easily simplified in “*nose*”.

This paper presents a unified framework for XAI allowing the design of a family of neural models, the *Logic Explained Networks* (LENs), which are trained to *solve-and-explain* a categorical learning problem integrating elements

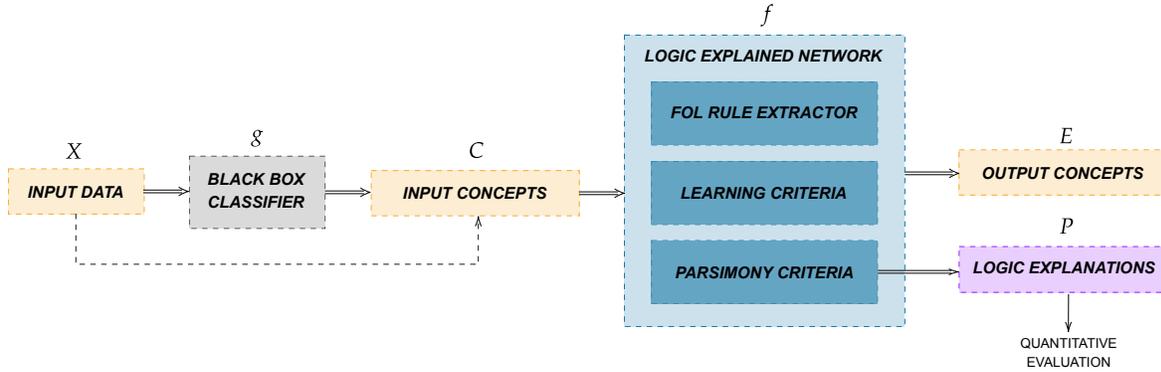


Figure 1: Logic Explained Networks (LENs,  $f$ ) are neural networks capable of making predictions of a set of output concepts (activation scores belonging to  $E$ ) and providing First-Order Logic explanations (belonging to  $P$ ) in function of the LEN inputs. Inputs might be other concepts (activation scores belonging to  $C$ ) either computed by a neural network ( $g$ ) or directly provided within the available data (each data sample belongs to  $X$ ). There are several different ways of instantiating this generic model into real-world problems (Section 3). Within the blue box, the key components of LENs are listed (Section 4).

from deep learning and logic. Differently from vanilla neural architectures, LENs can be directly interpreted by means of a set of FOL formulas. In order to implement such a property, LENs require their inputs to represent the activation scores of human-understandable concepts. Then, specifically designed learning objectives allow LENs to make predictions in a way that is well suited for providing FOL-based explanations that involve the input concepts. In order to reach this goal, LENs leverage parsimony criteria aimed at keeping their structure simple. There are several different computational pipelines in which a LEN can be configured, depending on the properties of the considered problem and on other potential experimental constraints. For example, LENs can be used to directly classify data in an explainable manner, or to explain another black-box neural classifier. Moreover, according to the user expectations, different kinds of logic rules may be provided. Due to this intrinsic versatility of LENs, what we propose can also be thought as a generic framework that encompasses a large variety of use cases.

Fig. 1 depicts a generic view on LENs, in which the main components are reported. A LEN (blue box – function  $f$ ) provides FOL explanations (purple box) of a set of output concepts (rightmost yellow-box) in function of the LEN inputs. Inputs might be other concepts (mid yellow box) computed by a neural network classifier (gray box – function  $g$ ) and/or concepts provided within the available data (leftmost yellow box). This generic structure can be instantiated in multiple ways, depending on the final goal of the user and on the properties of the considered problem. In order to provide the reader with an initial example/use-case (different configurations are explored in the paper), we consider an image classification problem with concepts organized into a two level hierarchy. In Fig. 2 we report an instance of Fig. 1 in which a CNN-based neural classifier gets an input image, predicting the activations of a number of low-level concepts. The LEN  $f$  processes such concepts, and it predicts the activation of higher-level output concepts. The LEN can provide a FOL description of each (high-level) output concept with respect to the (low-level) input ones. Another possible instance of the proposed framework consists in using the LEN to directly classify and explain input data (that is basically the case in which the input concepts are immediately available in the data themselves, and not the output of another neural model), thus the LEN itself becomes an interpretable machine. Moreover, LENs can be paired with a black-box classifier operating on the same input data, and forced to mimic as much as possible the behaviour of the black-box, implementing an additional explanation-oriented module.

We investigate three different use-cases that are inspired by the aforementioned instances, comparing different ways of implementing the LEN models. While most of the emphasis of this paper is on supervised classification, we also show how LEN can be leveraged in fully unsupervised settings. Additional human priors could be eventually incorporated into the learning process [Ciravegna et al., 2020b], in the architecture [Koh et al., 2020], and, following Ciravegna et al. [Ciravegna et al., 2020b, Ciravegna et al., 2020a], what we propose can be trivially extended to semi-supervised learning (out of the scope of this paper). Our work contributes to the XAI research field in the following ways.

- It generalizes existing neural methods for solving and explaining categorical learning problems [Ciravegna et al., 2020a, Ciravegna et al., 2020b] into a broad family of neural networks i.e., the *Logic Explained Networks* (LENs).
- It describes how users may interconnect LENs in the classification task under investigation, and how to express a set of preferences to get one or more customized explanations.

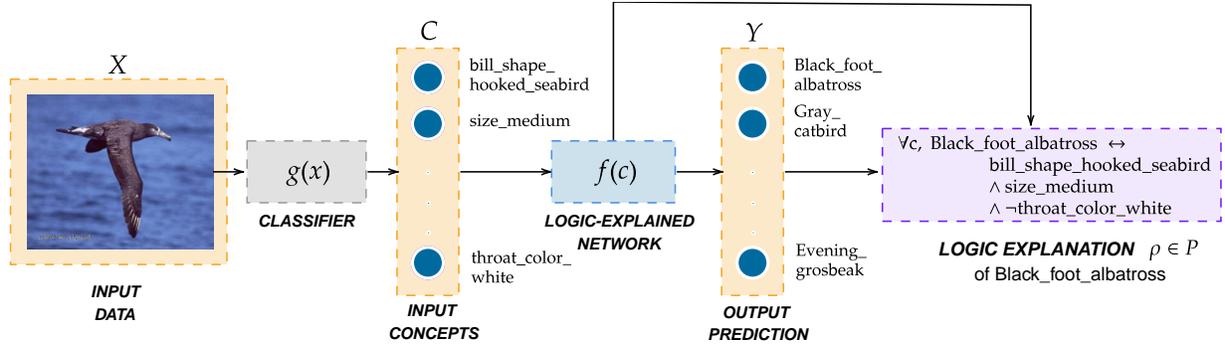


Figure 2: An example of a possible instance of the generic model of Fig. 1, inspired by the CUB 200-2011 fine-grained classification dataset. Classes are divided into a two-level hierarchy. A LEN is placed on top of a convolutional neural network  $g(\cdot)$  in order to (i) classify the species of the bird in input and (ii) provide an explanation on why it belongs to this class. The logic explanation in the example showcases the predicted output class (all the output concepts can be explained), dropping the argument of the predicates for compactness.

- It shows how to get a wide range of logic-based explanations, and how logic formulas can be restricted in their scope, working at different levels of granularity (explaining a single sample, a subset of the available data, etc.).
- It reports experimental results using three out-of-the-box preset LENs showing how they may generalize better in terms of model accuracy than established white-box models such as decision trees on complex Boolean tasks (in line with Tavares’ work [Tavares et al., 2020]).
- It advertises our public implementation of LENs through a Python package<sup>3</sup> with an extensive documentation about LENs models, implementing different trade-offs between interpretability/explainability and accuracy.

The paper is organized as follows (see also Fig. 3). Related works are described in Section 2. Section 3 gives a formal definition of a LEN and it describes its underlying assumptions, key paradigms and design principles. The methods used to extract logic formulas and to effectively train LENs are described in Section 4. Three out-of-the-box LENs, presented in Section 5, are compared on a wide range of benchmarks in terms of classification performance and quality of the explanations, in Section 6, including an evaluation of the LEN rules in an adversarial setting. Finally, Section 7 outlines the social and economical impact of this work as well as future research directions.

<sup>3</sup><https://pypi.org/project/torch-explain/>

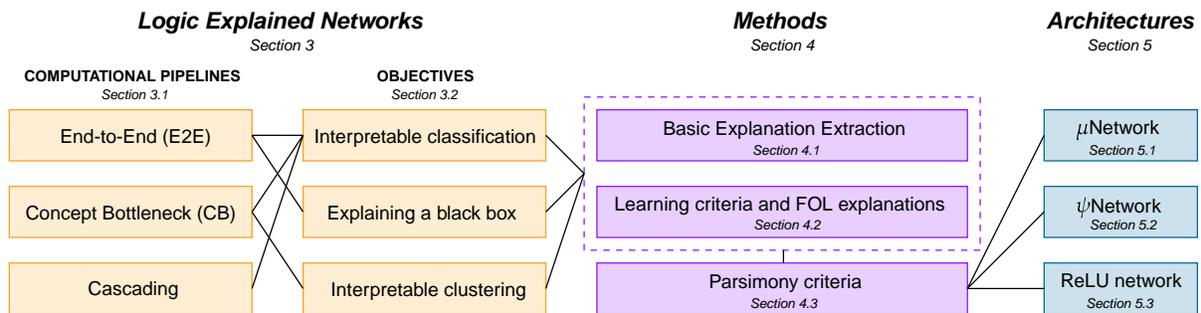


Figure 3: Visual overview of the organization of the paper for those sections that are about describing the whole LEN framework and the specific use-cases we selected. Starting from the nodes on the left, each path that ends to one of the nodes on the right creates a specific instance of the LEN framework. The proposed framework is generic enough to create several other instances than the ones we study in this paper.

|                           | Result Type     |            | Scope |        | Role             |              |
|---------------------------|-----------------|------------|-------|--------|------------------|--------------|
|                           | feature-scoring | rule-based | local | global | interpret. model | expl. method |
| LIME                      | ✓               |            | ✓     |        |                  | ✓            |
| SHAP                      | ✓               |            | ✓     |        |                  | ✓            |
| Activation Maximization   | ✓               |            | ✓     |        |                  | ✓            |
| Saliency Maps             | ✓               |            | ✓     |        |                  | ✓            |
| SP-LIME                   | ✓               |            |       | ✓      |                  | ✓            |
| Class Model Visualization | ✓               |            |       | ✓      |                  | ✓            |
| LORE                      |                 | ✓          | ✓     |        |                  | ✓            |
| Anchors                   |                 | ✓          | ✓     |        |                  | ✓            |
| DeepRED                   |                 | ✓          | ✓     | ✓      |                  | ✓            |
| GAM                       | ✓               |            |       | ✓      | ✓                |              |
| Decision Trees            |                 | ✓          | ✓     | ✓      | ✓                |              |
| BRL                       |                 | ✓          | ✓     | ✓      | ✓                |              |
| LENs                      |                 | ✓          | ✓     | ✓      | ✓                | ✓            |

Table 1: Summary of related work. The first column lists a number of approaches in the context of XAI. The other columns are about different properties. See the paper text for more details.

## 2 Related work

In the last few years, the demand for human-comprehensible models has significantly increased in safety-critical and data-sensitive contexts. This popularity is justified by the emerging need for unveiling the decision process of pitch-black models like deep neural networks. To this aim, the scientific community has developed a variety of XAI techniques, with different properties and goals. Several taxonomies have been proposed to categorize the XAI models, with partial overlapping and some ambiguities on the referred terminology. Without pretending to be exhaustive, in the following we focus on some key properties that are relevant to compare LENs with other existing approaches. In particular, we will describe related XAI methods considering the type of RESULT they provide (feature-scoring vs. rule-based), their specific ROLE (interpretable models vs. explanation methods), and the SCOPE of the provided explanations (local vs. global). A brief summary of a few XAI works in terms of these features is reported in Tab. 1 – for more details on existing taxonomies we refer to the recent surveys [Adadi and Berrada, 2018, Carvalho et al., 2019, Marcinkevičs and Vogt, 2020, Molnar, 2020].

XAI models can be differentiated according to the RESULT of the produced explanation. Most of the methods in literature usually focus on scoring or providing summary statistics of features [Erhan et al., 2010, Simonyan et al., 2013, Zeiler and Fergus, 2014, Ribeiro et al., 2016b, Ribeiro et al., 2016a, Lundberg and Lee, 2017, Selvaraju et al., 2017]. However, *feature-scoring* techniques can be of scarce utility in decision support cases, whereas a comprehensible language may bring light on the black-box behaviour by identifying concept-based correlations. On the other hand, *rule-based* methods are generally more comprehensible [Breiman et al., 1984, Angelino et al., 2018, Letham et al., 2015], since they usually rely on a formal language, such as FOL. Further, the learned rules can be directly applied to perform the learning task in place of the explained model. Existing approaches have different ROLES in the XAI landscape, acting as intrinsically *interpretable models* or as *explanation methods*. As a matter of fact, the interpretability of a model can be achieved either by constraining the model to be interpretable per se (*interpretable models*) or by applying some methods to get an explanation of an existing model (*explanation methods*). Interpretable models are by definition model specific, while explanation methods can rely on a single model or be, as more commonly happens, model agnostic. In principle, interpretable models are the best suited to be employed as decision support systems. However, their decision function often is not smooth which makes them quite sensible to data distribution shifts and impairs their generalization ability on unseen observations [Tavares et al., 2020, Molnar, 2020]. On the other hand, explanation methods can be applied to get approximated interpretations of state-of-the-art models. These methods are known as “post hoc” techniques, as explanations are produced once the training procedure is concluded. Finally, a fundamental feature to distinguish among explainable methods is the SCOPE of the provided explanations. *Local* explanations are valid for a single sample, while *global* explanations hold on the whole input space. However, several models consider local explanations together with aggregating heuristics to get a global explanation. Some of the most famous XAI techniques share the characteristics of being feature-scoring, local, post hoc methods.

Sometimes the easiest way to solve the explanation problem is simply to treat the model as a black-box and, sample-by-sample, determine which are the most important features for a prediction. Prominent examples of algorithms falling in this area accomplish this task by perturbing the input data. Local Interpretable Model-agnostic Explanations (LIME) [Ribeiro et al., 2016a] trains a white-box model (e.g. a logistic regression) to mimic the predictions of a given model in the neighbourhood of the sample to explain. By analyzing the weights of the white-box model it identifies the most important group of pixels (superpixel). Differently, SHapley Additive exPlanations (SHAP) [Lundberg and Lee, 2017] computes the Shapley value of each feature by removing each of them in an iterative manner. Other important techniques provide the same type of explanations through gradient analysis. For instance, Erhan, Courville, and Bengio introduced the Activation Maximization [Erhan et al., 2010] framing this as an optimization problem. They explain the behaviour of a hidden/output unit by slightly modifying a given input data such that it maximizes its activation. An easier way for identifying the most relevant input features consists in computing Saliency Maps [Simonyan et al., 2013], which backtrack the classification loss back to the input image. Other post hoc methods have been devised to extract global feature-scoring explanations, often starting from local explanations. A submodular-pick algorithm extends the LIME approach (SP-LIME) to provide global explanations [Ribeiro et al., 2016a]. SP-LIME first finds superpixels of all input samples with the standard LIME procedure. Successively, it identifies the minimum number of common superpixels covering most of the images. Starting from the idea of Activation Maximization, Class Model Visualization [Simonyan et al., 2013] searches the input sample maximizing class probability scores in the whole input space. While inheriting several of the above properties, LENs take a different perspective, as they aim at providing human-comprehensible explanations in terms of FOL formulas (see Tab. 1).

Existing XAI methods can also supply FOL explanations. Local Rule-based Explanations of black-box decision systems (LORE) [Guidotti et al., 2018] extracts FOL explanations via tree induction. Here the authors focus on local rules extracted through input perturbation. For each sample, a simple decision tree is trained to mimic the behaviour of the black-box model in the neighbourhood of the sample. Both standard and counterfactual explanations can then be extracted following the branches of the tree. Also the Anchors method [Ribeiro et al., 2018], based on LIME, provides local rules of black-box model predictions via input perturbation. In this case, rules are not extracted from a decision tree but are generated by solving a Multi-Armed-Bandit beam search problem. Extending the CRED algorithm [Sato and Tsukimoto, 2001], DeepRED [Zilke et al., 2016] employs a decomposition strategy to globally explain neural networks. Starting from the output nodes, the predictions of each neuron are explained in terms of the activations of the neurons in the previous layer by training a decision tree. In the end, all rules for each class are merged in a single global formula in terms of input features. For a given sample, a unique local rule is extracted following the firing path. The approach proposed in this paper resembles the DeepRED algorithm. As it will become clear in the following sections, LENs can explain the behaviour of a neural network both end-to-end and layer by layer. However, LENs can be used to explain any black-box model as far as its input and output correspond to human-interpretable categories (Section 3). Furthermore, LENs support different forms of FOL rules, with different scopes and goals.

Interpretable models are capable of providing both local and global explanations. Such explanations may be based either on feature rankings, as in Generalized Additive Models (GAM) [Hastie and Tibshirani, 1987], or on logic formulas, as in Decision Trees [Breiman et al., 1984] or Bayesian Rule Lists (BRL) [Letham et al., 2015]. GAMs overcome the linearity assumption of linear regression by learning target categories disjointly from each feature. Caruna et al. proposed GA<sup>2</sup>M where pairs of features are allowed to interact [Caruana et al., 2015]. A further extension employing neural networks as non-linear functions has recently been proposed by Agarwal et al. [Agarwal et al., 2020], which, however, loses interpretability at feature-level. Decision trees [Breiman et al., 1984] are a greedy learning algorithm partitioning input data into smaller subsets until each partition only contains elements belonging to the same class. Different pruning algorithms have been proposed to simplify the final structure to get simple explanations [Quinlan, 1987]. Each path of a decision tree is equivalent to a decision rule, i.e. an *IF-THEN* statement. Other approaches focus on generating sets of decision rules, either with sequential covering [Cohen, 1995] or by selecting the best rules from a pre-mined set via Bayesian statistics [Letham et al., 2015]. Interestingly, LENs can be used to solve a learning task directly, as they are interpretable per se. Employing an interpretable-by-design neural network has the advantage that the extracted explanations will perfectly match the classifier predictions. In addition, relying on a neural network, LEN generalization performance can be much better than standard interpretable-by-design methods, whose representation capacity is usually limited.

To sum up, the proposed family of neural networks can be used both as interpretable classifiers and as surrogate models. Indeed, the flexibility of the proposed framework allows the user to find an appropriate trade-off between interpretability/explainability and accuracy that is well suited for the task at hand. Concerning the explanation type, LENs provide explanations as FOL formulas. The inner mechanism to extract such explanations is general enough to cover rules with different scopes, from local rules valid on a single sample to global rules that hold for an entire class.

### 3 Logic Explained Networks

This work presents a special family of neural networks, referred to as Logic Explained Networks (LENs), which are able to both make predictions and provide explanations. LENs may explain either their own predictions or the behaviour of another classifier, being it a neural network or a generic black-box model. Moreover, a LEN might also be used to explain relationships among some given data. In this section, we describe how LENs can be instantiated in the framework of Fig. 1. We start by introducing the notation (following paragraphs). Then we discuss a variety of computational pipelines (Section 3.1) and illustrate different types of explanations LENs can generate in function of different learning objectives (Section 3.2). In Fig. 3 (leftmost part) we provide an overview of the contents of this section. In the same figure, we also show how the following sections will cover the specific methods (Section 4) and the considered neural architectures (Section 5).

A LEN is a function  $f$ , implemented with a neural network, that maps data falling within a  $k$ -dimensional Boolean hypercube onto data falling within another  $r$ -dimensional Boolean hypercube. Each dimension is about the activation strength of what we refer to as a *concept*, with the strict requirement of having a human-understandable description of each *input* dimension/concept [Kim et al., 2018]. Formally, a LEN is a mapping  $f : C \rightarrow E$ , where  $C = [0, 1]^k$  is the *input concept space*, and  $E = [0, 1]^r$  is the so called *output concept space*. FOL explanations produced by a LEN are about the relationships between the output and the input concepts, and they belong to the generic *rule space*  $P$ . In particular, whenever a LEN has been trained, the  $i$ -th output  $f_i$  can be directly translated into a logic rule  $\varphi_i$  that involves the input concepts and that is leveraged to devise a FOL formula  $\rho_i \in P$ , such as the one shown in Fig. 2. The FOL extraction process is general enough to offer logic rules with different levels of granularity, from a local to a more global coverage of the available data.

In order to provide a concrete meaning to LEN models, the source of the input concept activations needs to be defined as well as the learning criteria. Amongst a variety of possible configurations that are instance of Fig. 1, we will focus on a limited set of computational pipelines in which the input and output concept spaces are defined with respect to real-world use cases. In some of them, LENs communicate with a black-box classifier, with the goal of providing explanations of its predictions, or with the goal of leveraging its activations as input concepts. In all cases, for each input sample, LENs provide  $r$  FOL explanations that might either be directly associated to  $r$  categories of a classification problem or they could be interpreted as  $r$  explanations of unknown relationships among input data. In the former case, LENs are trained in a supervised manner, while in the latter case they are trained in an unsupervised setting.

Before going into further details, we introduce the main entities and the notation that will be used throughout the paper, paired with a short description to help the reader in following the paper and to have a quick reference to the main symbols.

- $f, C, E$ : The function computed by a LEN is  $f : C \rightarrow E$ , where  $C = [0, 1]^k$  is the space of the activations of the  $k$  input concepts, and  $E = [0, 1]^r$  is about the activations of the  $r$  output concepts.
- $X, \mathcal{X}, C$ : We consider a scenario in which a finite collection  $\mathcal{X}$  of data samples that belong to  $X \subseteq \mathbb{R}^d$  is available to train LENs. There exists a mapping from  $X$  to the space  $C$  of input concept activations. The notation  $C$  indicates the finite set of concept activations obtained by mapping each sample of  $\mathcal{X}$  onto  $C$ .
- $q$ : We use the notation  $q$  to indicate the total number of classes in a classification problem. Each data sample can be associated with one or more classes. No strict conditions are posed on the potential mutual exclusivity of the different classes, thus we consider the most generic setting that spans from multi-label to single-label classification. Moreover, classes could also be organized in a structured manner, such as in hierarchy.
- $\bar{y}, \bar{y}_i$ : Whenever we consider classification problems, classes are assumed to be encoded with binary targets in  $\{0, 1\}^q$ . The function  $\bar{y}(\cdot)$  returns the target vector associated to the data sample passed as its argument, while  $\bar{y}_i(\cdot)$  returns the  $i$ -th component of such a vector, that is about the  $i$ -th class. What we propose holds also in the case in which targets are encoded with scores in the unit interval, so that we will frequently refer to generic data encoded in  $Y = [0, 1]^q$ .
- $Y^{(a:b)}$ : We will use the notation  $Y^{(a:b)}$  to indicate a view of  $Y$  limited to the dimensions going from index  $a \geq 1$  to index  $b \leq q$ , included.
- $g$ : The function  $g(x)$  is about a generic black-box neural classifier that computes the membership scores of  $x$  to a set of categories. Such categories could be exactly the ones of the considered classification problem (encoded in  $Y$ ) or a subset of them (encoded in  $Y^{(a:b)}$ ). In detail,  $g : X \rightarrow Y$  (resp.  $g : X \rightarrow Y^{(a:b)}$ ), and  $g(x)$  defines the membership scores of  $x$  to the considered categories. Of course,  $g_i(x) = 1$  when  $x$  strongly belongs to class  $i$  (resp.  $a + i - 1$ ). No special conditions are enforced in the definition of  $g$ .

- $\bar{f}, \bar{f}_i$ : The notation  $\bar{f}$  is used to indicate a Boolean instance of the main LEN function  $f$ , in which each output of  $f$  is projected to either 0 or 1.<sup>4</sup> In order to refer to a single output of the vector function  $f$  (or  $\bar{f}$ ), the subscript will be used, e.g.  $f_i$  (or  $\bar{f}_i$ ).
- $c_j, \bar{c}_j$ : Similarly, for a vector  $c \in C$  with activation scores of  $k$  concepts,  $c_j$  is the  $j$ -th score, while  $\bar{c}_j$  is the Boolean instance of it.
- $\bar{c}_j$  (name): Any dimension of the space of input concepts  $C$  includes a human-readable label. Logic rules generated by LENs will leverage such labels to give understandable names to predicates. In order to simplify the notation, whenever we report a logic rule, we will use the already introduced symbol  $\bar{c}_j$  also to refer to the human-understandable name of the  $j$ -th input concept. This notation clash makes the presentation easier.
- $\bar{y}_i$  (name): Any dimension of the space of output concepts  $E$  might or might not include a human-readable label, whether we consider supervised or unsupervised learning when training LENs, respectively. In the former case, the already introduced notation  $\bar{y}_i$  will also refer to the human-understandable name of the  $i$ -th output class, following the same simplification we described in the case of  $\bar{c}_j$ .
- $\varphi_i$ : Any output  $f_i$  is associated with a logic rule  $\varphi_i$  given in terms of (the names of) the input concepts.
- $\rho_i$ : We indicate with  $\rho_i \in P$  the FOL formula that explains the  $i$ -th output concept leveraging  $\varphi_i$ . The precise way in which  $\varphi_i$  is used to build  $\rho_i$  depends on whether we are considering supervised or unsupervised output concepts.

The listed elements will play a precise role in different portions of the LEN framework.

### 3.1 Computational pipelines

Amongst a large set of feasible input/output configurations of the LEN block, here we consider three different computational pipelines fitting three concrete scenarios/use-cases where LENs may be applied, that consist in what we refer to as END-TO-END, CONCEPT BOTTLENECK, and CASCADING pipelines.

END-TO-END (E2E). The most immediate instance of the LEN framework is the one in which the LEN block aims at directly explaining, and eventually solving, a supervised categorical learning problem. In this case, we have  $C = X$ , and  $r = q$ , as shown in Fig. 4. In order to make the first equality valid with respect to the LEN requirements, the  $d$  input features of the data in  $X$  must score in  $[0, 1]$  (or in  $\{0, 1\}$  in the most extreme case), so that they can be considered as activations of  $d$  human-interpretable concepts. In order to create these conditions in different types of datasets, generic continuous features can be either discretized into different bins or Booleanized by comparing them to a reference ground truth (e.g. for gene expression data, the gene signature obtained from control groups can be used as a reference). (i) An E2E LEN can perform a fully *interpretable classification* task, having the double role of main classifier and explanation provider (Fig. 4, top), or (ii) it can work in parallel with another neural black-box classifier  $g$ , thus providing explanations of the predictions of  $g$  (Fig. 4, bottom) – what we refer to as *explaining a black-box*. In both the cases (i) and (ii) we have  $E = Y$ . The first solution is preferred when a white-box classifier is of utmost importance, while a loss in terms of classification accuracy is acceptable, due to the constrained nature of the LENs. The second solution is useful in contexts where the classification performance plays a key role, and approximate explanations of black-box decisions are sufficient.

CONCEPT-BOTTLENECK (CB). A Concept-Bottleneck LEN is a computational pipeline in which the LEN  $f$  aims at explaining, and eventually solving, a categorical learning problem whose features do not correspond to human-interpretable concepts and, as a consequence, they are not suitable as LEN inputs. In this case, a black-box model  $g$  computes the activations of an initial set of  $z$  concepts out of the available data  $\mathcal{X}$ . The LEN solves the problem of predicting  $r$  new concepts from the activations of the  $z$  ones, as shown in Fig. 5 (top/bottom). Differently from the E2E case, there is always a neural network  $g$  processing the data, and the outcome of such processing is the input of the LEN. Formally, we have  $g : X \rightarrow Y^{(1:z)}$ , while  $f : C \rightarrow E$  with  $C = Y^{(1:z)}$ , with  $z \leq q$ ,<sup>5</sup> and where the meaning of  $E$  varies in function of what we describe in the following. (i) This two-level scenario can be implemented in a way that is coherent with the already discussed example of Fig. 2, that is what we also show in Fig. 5 (top). In this case, each example is labeled with  $q$  binary targets divided into two disjoint sets, where the black-box network  $g$  predicts the first  $z < q$  ones, and the LEN  $f$  predicts the remaining  $r = q - z$ , i.e.  $E = Y^{(z+1:q)}$ . The LEN will be both responsible of predicting the higher level concepts and of explaining them in function of the lower level concepts yielded by the black-box model, thus still falling within the context of *interpretable classification*. (ii) We also consider the case in which the outcome of the LEN is not associated to any known categories. In this case, the

<sup>4</sup>When not directly specified, we will assume 0.5 to be used as threshold value to compute the projection.

<sup>5</sup>We implicitly assumed the axes of  $Y$  to be sorted so that the first  $z$  dimensions are the ones we want to predict with  $g$ , and we will make this assumption in the whole paper.

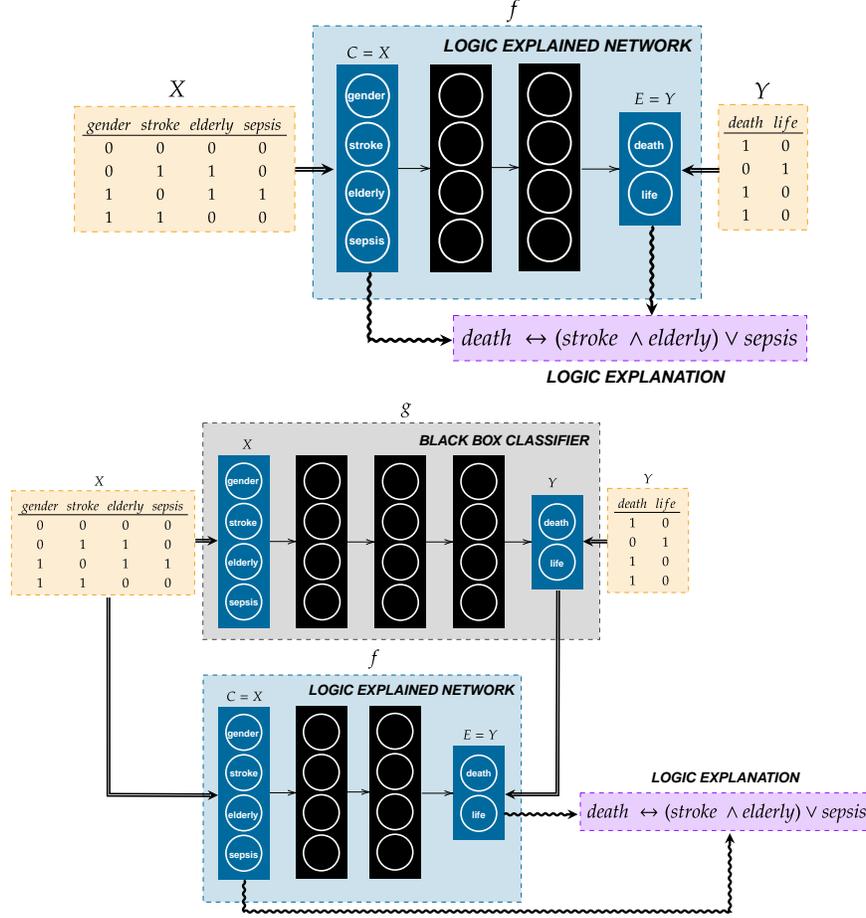


Figure 4: End-to-end (E2E) LENs directly work on input data that are interpretable per se and treated as concepts ( $C = X$ ), while the output concepts are the activation scores of the classes of the dataset ( $r = q = 2$ ). This is a real-world example from the MIMIC II dataset (see the experimental section), where patient features are used to classify if the patient will survive 28 days after hospital recovery. Top: the LEN solves the classification problem and provides explanation, also referred to as *interpretable classification*. Bottom: the LEN *provides explanations* of a black-box classifier. The universal quantifier and the argument of the predicates have been dropped for simplicity.

LEN is trained in an unsupervised manner, as shown in Fig. 5 (bottom), thus implementing a form of *interpretable clustering*. Formally,  $g : X \rightarrow Y$ ,  $C = Y$ , and  $E = [0, 1]^r$ , with customizable  $r \geq 1$ . In both the cases (i) and (ii), the black-box classifier  $g$  is trained in a supervised manner.

**CASCADING.** Cascading LENs can be used to provide explanations by means of a hierarchy of concept layers. In particular, multiple LENs are used to map concepts into a higher level of abstraction (without the need of any black-box  $g$ ), as shown in Fig. 6. Each LEN provides explanations of its outputs with respect to its inputs, allowing LENs to handle multiple levels of granularity. This structure provides a hierarchy of explanations and upholds human interventions at different levels of abstraction, providing a far deeper support for human-machine interactions. We indicate with  $(C_1, E_1), (C_2, E_2), \dots, (C_s, E_s)$  the input/output concept spaces of the  $s$  cascading LENs, with  $C_j = E_{j-1}$ ,  $j > 1$ . In the experiments, we will consider the case in which LENs are trained in a supervised manner, thus implementing multiple *interpretable classifications*.

### 3.2 Objectives

In Section 3.1, when describing the selected computational pipelines, we made an explicit distinction among INTERPRETABLE CLASSIFICATION, EXPLAINING A BLACK-BOX, and INTERPRETABLE CLUSTERING, as they represent three different “objectives” the user might consider in solving the problem at hand. These objectives impose precise constraints in the selection of the learning criteria and on the form of the FOL formulas that can be obtained by the

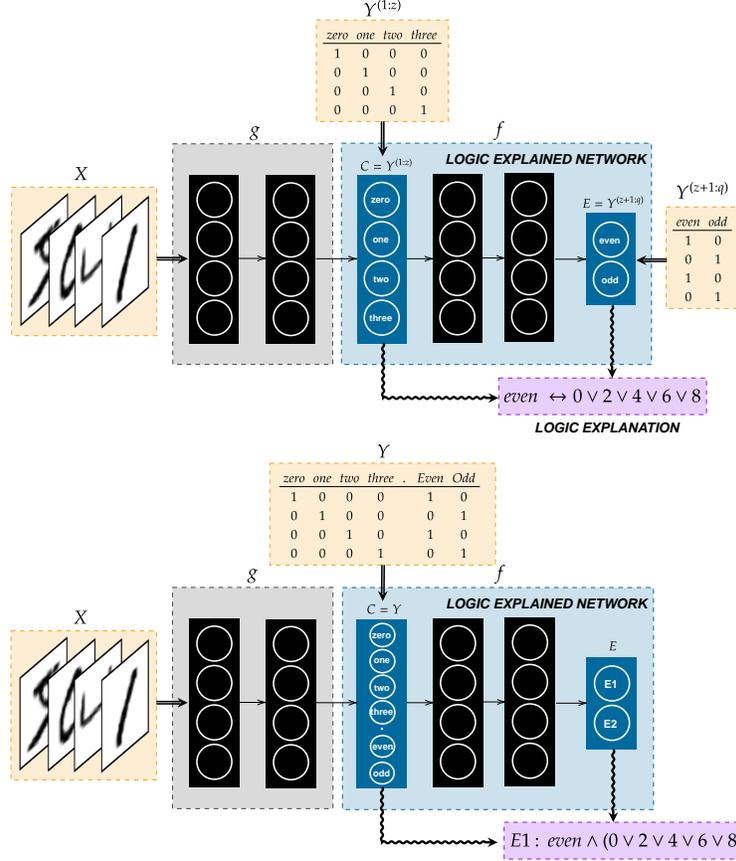


Figure 5: Concept-bottleneck (CB) LEN. Left: the LEN is placed on top of a black-box model  $g$  which maps the input data into a first set of interpretable concepts. An MNIST-based experiment is shown (see the experimental section). Handwritten digits are first classified by  $g$ . A LEN is then employed to classify and explain whether the digit is even or odd – *interpretable classification*. Top: MNIST digits are classified as belonging to one of the 10 classes and whether they are even or odd by a black-box  $g$  (see the experimental section). Bottom: A LEN groups these predictions in an unsupervised manner, and analyzes the relations within each cluster – *interpretable clustering*. Supervision labels are not provided for the LEN output. The universal quantifier and the argument of the predicates have been dropped for simplicity.

LENs. The form of FOL formulas will be described in detail in Section 4.1 (generic process of logic rule extraction), and in Section 4.2 (learning criteria and FOL rules). The key difference among the aforementioned objectives is about the criterion that drives the learning dynamics of  $f$ , as each LEN module can be trained in a supervised or an unsupervised fashion.

**INTERPRETABLE CLASSIFICATION.** Whenever LENs are leveraged to both solve the classification problem and provide explanations, i.e. in INTERPRETABLE CLASSIFICATION, learning is driven by supervised criteria. Following Ciravegna et al. [Ciravegna et al., 2020a], supervised criteria leverage the available data so that: (i)  $f$  learns class labels as in classic supervised learning, and (ii)  $f$  is constrained to be coherent with the target type of FOL rules. Each output neuron of the LEN is associated to a target category named  $\bar{y}_i$  of the considered learning problem. LENs can provide a FOL explanation  $\rho_i \in P$  of such category, so that the class-predicate  $\bar{y}_i(\cdot)$  will be involved in the extracted FOL formula. For example, if *person* is one of the output categories, LENs can explain the reasons behind the prediction of class  $\bar{y}_i = person$  by means of an automatically discovered relationship  $\varphi_i$  involving the activations/not-activations of some ( $m_i$ ) of the  $k$  input concepts ( $m_i < k$ ). If such input concepts are  $\bar{c}_j = head$ ,  $\bar{c}_z = hands$ ,  $\bar{c}_h = body$ , the system could learn that  $\rho_i = \forall c \in C : \bar{y}_i(c) \leftrightarrow \varphi_i(c)$ , where  $\varphi_i(c) = \bar{c}_j(c) \wedge \bar{c}_z(c) \wedge \bar{c}_h(c)$ , i.e. (discarding the quantifier)  $person \leftrightarrow head \wedge hands \wedge body$ . Interestingly, in the case of interpretable classification, LENs basically become white-box classifiers, as is showcased by the concrete examples of Fig. 4 (top), Fig. 5 (top), and Fig. 6.

**EXPLAINING A BLACK-BOX.** In case the user aims at EXPLAINING A BLACK-BOX, LENs act in parallel to black-box classifiers with the goal of explaining the decisions of such black-box models. In this scenario, LENs are constrained

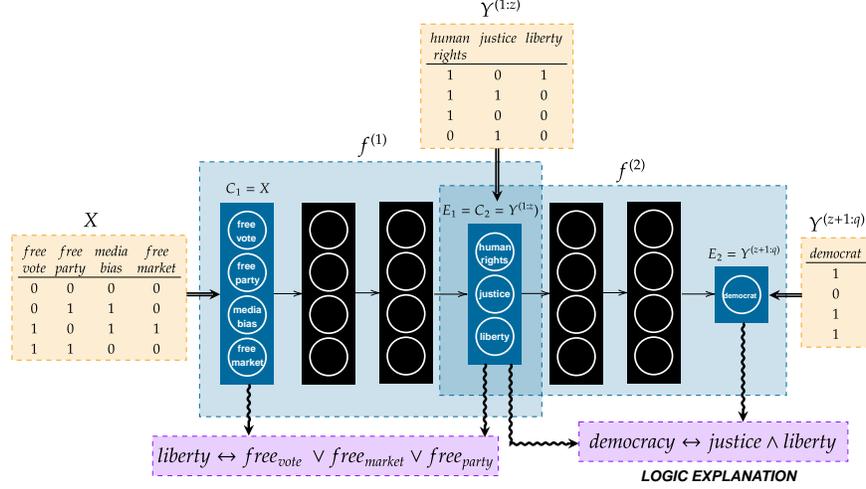


Figure 6: Cascading LENS, an example taken from the V-Dem classification dataset (see the experimental section). The final classification on the status of the democracy is divided into two steps. First, the input data/concepts  $C_1$  are mapped into high-level concepts  $C_2$ . High-level concepts are then used to compute the final classification. Cascading LENS can provide explanations at different levels of granularity, implementing multiple *interpretable classifications*. The universal quantifier and the argument of the predicates have been dropped for simplicity.

to mimic the predictions of the black-boxes, i.e.  $f$  is forced to be close to  $g$  when evaluated on the available data samples. This is what is shown in Fig. 4 (bottom) and it may not require any supervision, since learning can be driven by a coherence criterion between the outputs of  $g$  and those of  $f$  when processing the same data. The type of FOL rules LENS can extract are the same of the case of interpretable classification, thus the same principles are followed.

**INTERPRETABLE CLUSTERING.** Differently, LENS can be trained using unsupervised criteria, whenever the user is not aiming at explaining the target classes of a classification problem, but is interested in discovering generic relations among input concepts. We refer to this objective as **INTERPRETABLE CLUSTERING**. The user might be interested in discovering co-occurrences  $\varphi_i$  of input concept activations in not-defined-before subsets  $O_i$  of the concept space. This leads to FOL rules such as  $\rho_i = \forall c \in O_i: \varphi_i(c)$ , where, for example,  $\varphi(c) = \bar{c}_e(c) \vee \bar{c}_t(c)$  and  $\bar{c}_e = soccer\_ball$ ,  $\bar{c}_t = foot$ , and the set  $O_i$  can be thought as a cluster. In this scenario, LENS’s output neurons are not associated to any known categories (in contrast to previous objectives), but to unspecified generic symbols, as shown in Fig. 5 (bottom). The activation score of the output concept represents a cluster membership score and it is used to define whether an input pattern belongs to the subset  $O_i$  or not.

In all the discussed objectives, logic formulas are extracted from LENS using the same principles. Rules are then instantiated into FOL formulas that well cope with objective-specific learning criteria. As it will become clear in the following section, due to such generality of the extraction mechanisms, the user might decide to automatically discover definite regularities focused on single examples, groups of data points, or, more generally, the whole dataset, moving from local to more global explanations.

## 4 Methods

This section presents the fundamental methods used to implement Logic Explained Networks introduced in Section 3. We start by describing the procedure that is used to extract logic rules out of LENS for an individual observation or a group of samples (Section 4.1). This procedure is common to all LENS’ objectives/use-cases. Then, we provide the formal definition of the learning objectives constraining LENS to provide required types of FOL explanations (Section 4.2). Finally, we discuss how to constrain LENS to yield concise logic formulas. To this aim, ad-hoc parsimony criteria (Section 4.3) are employed in order to bound the complexity of the explanations. In Fig. 3 (middle) we provide an overview of the contents of this section.

### 4.1 Extraction of logic explanations

Once LENS are trained, a logic formula can be associated to each output concept  $f_i$ . As it will become clear shortly, extracting logic formulas out of trained LENS can be done by inspecting its inputs/outputs accordingly to their Boolean

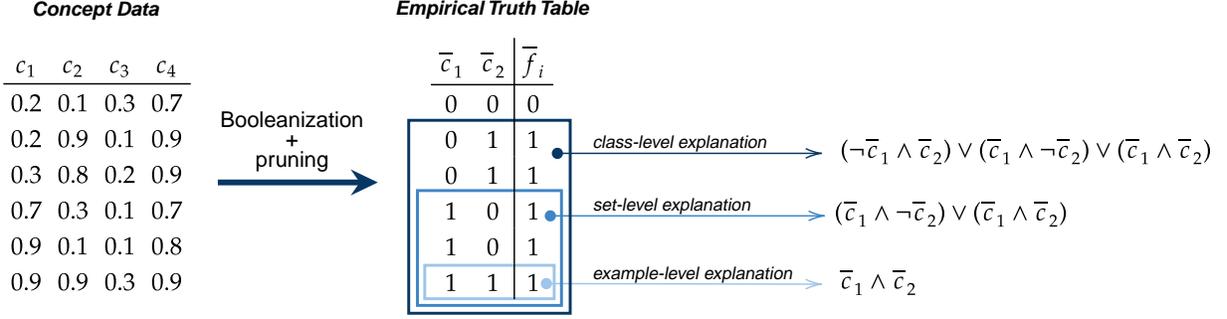


Figure 7: Empirical truth table  $\mathcal{T}^i$  of the  $i$ -th LEN output  $f_i$ , with  $k = 4$  and  $m_i = 2$  (assuming that only the first two input concepts are kept). The aggregation of concept tuples with the same Booleanization yields a common example-level explanation, and do not complicate the class-level explanation. For any example-level explanation, we may count how many samples it explains, to discard the most infrequent cases.

interpretation. We have already introduced the notation  $\varphi_i$  to indicate the logic explanation of the output concept  $i$ . This generic notation will be properly formalized in the following. We overload the symbol  $\varphi_i$  to explicitly indicate, when needed, the data subset where the logic explanation holds true, using the notation  $\varphi_{i,\dots}$ . Here the second subscript can refer either to a single data sample  $c$ ,  $\varphi_{i,c}$ , or to a set  $S$  of data samples,  $\varphi_{i,S}$ . In practice,  $S$  denotes the region of the concept space that is covered by the  $i$ -th explanation, i.e. the set of concept tuples for which the formula  $\varphi_{i,S}$  is true. By aggregating over multiple samples, the scope of the logic formula may be tuned from strictly local EXAMPLE-LEVEL EXPLANATIONS ( $S = \{c\}$ ) to SET-LEVEL EXPLANATIONS ( $S \subseteq \mathcal{C}$ ), where the latter can be focused on a precise class, i.e., CLASS-LEVEL EXPLANATIONS. Eventually, for  $S = \mathcal{C}$ , global logic formulas holding on the whole concept space  $\mathcal{C}$  can be extracted.

To allow the extraction of FOL formulas, any LEN  $f = (f_1, \dots, f_r)$  requires both its inputs and outputs to belong to the real-unit interval. This apparent limitation allows any  $f_i$ , for  $i = 1, \dots, r$ , to correspond to a logic formula. First,  $f$  maps the data in  $\mathcal{C}$  into the rule space  $\mathcal{E}$ . After this forward pass, both the input data  $\mathcal{C}$  and the predictions of  $f$  are thresholded, e.g. with respect to 0.5, to obtain their Boolean values. Then, for each output neuron  $i$ , an *empirical truth-table*  $\mathcal{T}^i$  is built by concatenating the  $k$ -columns of Booleanized input concept tuples  $\{\bar{c}, \forall c \in \mathcal{C}\}$ , with the column of the corresponding LEN's predictions  $\bar{f}_i(c)$  (left-side of Fig. 7). The truth-table  $\mathcal{T}^i$  can be converted into a logic formula  $\varphi_i$  in Disjunctive Normal Form (DNF) as commonly done in related literature [Mendelson, 2009]. However, the rationale behind LENs is to extract formulas that are compact, emphasizing the most relevant relationships among the input concepts, according to specific parsimony indexes (that will be the subject of Section 4.3). Thus, any  $f_i$  will depend only on a proper subset of  $m_i \leq k$  concepts, and the formula  $\varphi_i$  will be built according to the restriction of  $\mathcal{T}^i$  to  $m_i \leq k$  columns (see e.g. Fig. 7). Notice that, for convenience in the notation, we assumed the first  $m_i$  columns to be the ones playing a role in the explanation, even if they could be any set of  $m_i$  columns of  $\mathcal{T}^i$ . In order to give more details about the rule extraction, we formally introduce the set  $O_i = \{c \in \mathcal{C} : \bar{f}_i(c) = 1\}$  as the set of all the sampled concept tuples that make true the  $i$ -th output explanation, i.e. the *support* of  $f_i$ .

EXAMPLE-LEVEL EXPLANATIONS. Given a sample  $c \in O_i \subseteq \mathcal{C}$ , the Booleanization  $\bar{c}$  of its continuous features may provide a natural way to get an example-level logic explanation  $\varphi_{i,c}$ . To make logic formulas more interpretable, the notation  $\tilde{c}$  denotes human-interpretable strings representing the concept names or their negation,

$$\varphi_{i,c} = \tilde{c}_1 \wedge \dots \wedge \tilde{c}_{m_i} \quad \text{where } \tilde{c}_j := \begin{cases} \bar{c}_j, & \text{if } c_j \geq 0.5 \\ \neg\bar{c}_j, & \text{if } c_j < 0.5 \end{cases}, \text{ for } j = 1, \dots, m_i. \quad (1)$$

SET-LEVEL AND CLASS-LEVEL EXPLANATIONS. By considering Eq. 1 for any  $c \in S$ , with  $S \subseteq O_i$ , and aggregating all the example-level explanations, an explanation for a set of samples can be generated as follows:

$$\varphi_{i,S} = \bigvee_{c \in S} \varphi_{i,c} = \bigvee_{c \in S} \tilde{c}_1 \wedge \dots \wedge \tilde{c}_{m_i} \quad (2)$$

As some  $\varphi_{i,c}$ 's might be equivalent for different  $c$ 's, repeated instances can be discarded keeping only one of them, without loss of generality. In case  $S = O_i$ , we simply write  $\varphi_i$  in place of  $\varphi_{i,S}$  and we refer to such set-level explanation as the CLASS-LEVEL EXPLANATION of the  $i$ -th output concept  $f_i$ .

In the rest of this section we will explore further details of what has been described so far. We will start by the following example.

**Example 1.** Let’s consider the Boolean XOR function, defined by  $xor(0, 0) = xor(1, 1) = 0$ ,  $xor(1, 0) = xor(0, 1) = 1$ . Let  $f = [f_1]$  be a LEN that has been trained to approximate the XOR function in the input space  $C = [0, 1]^2$ . Then, if we consider, e.g. the inputs  $c^1 = (0.2, 0.7)$ ,  $c^2 = (0.6, 0.3)$ , we get  $\bar{c}^1 = (0, 1)$ ,  $\bar{c}^2 = (1, 0)$ , and therefore  $\bar{f}_1(c^1) = \bar{f}_1(c^2) = 1$ . These examples yield the example-level explanations  $\varphi_{1,c^1} = \neg\bar{c}_1 \wedge \bar{c}_2$  and  $\varphi_{1,c^2} = \bar{c}_1 \wedge \neg\bar{c}_2$  respectively. As a result, the class-level explanation for  $f_1$  is given by  $\varphi_1 = (\neg\bar{c}_1 \wedge \bar{c}_2) \vee (\bar{c}_1 \wedge \neg\bar{c}_2)$ , which correctly matches the truth-table of the Boolean XOR function.

It is worth noting that both example and set-level explanations can be recursively applied in case of LENs with multiple  $[0, 1]$ -valued hidden layers or in case of cascading LENs. For instance, given a cascading LEN, as the one in Fig. 6, we may get both example and class-level explanations of the concepts in  $E_1 = C_2$  with respect to the ones in  $C_1$ , so as the concepts in  $E_2$  with respect to the ones in  $C_2$  and, in turn, in  $C_1$ . The modularity of logic formulas allows the composition of categories at different levels that may express arbitrary complex relationships among concepts.

Logic explanations  $\varphi_{i,S}$  generally hold only on a sub-portion  $S$  of the sampled concept space  $C$ . However, we may get a logic formula providing an explanation holding everywhere by means of the disjunction  $\varphi = \varphi_1 \vee \dots \vee \varphi_r$ , if we assume  $O_1 \cup \dots \cup O_r = C$ . Since  $\varphi$  can turn out to be of little significance, if we are interested in simpler (and clearer) global explanations we may convert  $\varphi$  into an equivalent  $\varphi'$  in *Conjunctive Normal Form* (CNF). In particular, there always exist  $r'$  and some clauses  $\varphi'_1, \dots, \varphi'_{r'}$  such that:

$$\varphi = \bigvee_{i=1}^r \varphi_i \equiv \bigwedge_{i=1}^{r'} \varphi'_i = \varphi'. \quad (3)$$

As a result, we get a set of  $r'$  explanations holding on the whole  $C$ , indeed  $\varphi'_i(c) = 1$  for every  $c \in C$ ,  $i = 1, \dots, r'$ . Unfortunately, converting a Boolean formula from DNF into CNF can lead to an exponential explosion of the formula. However, after having converted  $\varphi_i$  in CNF, the conversion can be computed in polynomial time with respect to the number of minterms in  $\varphi_i$  [Russell and Norvig, 2016].

The methodologies described so far illustrate how logic-based explanations can be aggregated to produce a wide range of explanations, from the characterization of individual observations to formulas explaining model predictions for all the samples leading to the same output concept activation. The formula for a whole class can be obtained by aggregating all the minterms corresponding to example-level explanations of all the observations having the same concept output. In theory, this procedure may lead to overly long formulas as each minterm may increase the complexity of the explanation. In practice, we observe that many observations share the same logic explanation, hence their aggregation may not change the complexity of the class-level formula (right-side Fig. 7). In general, “satisficing” class-level explanations can be generated by aggregating the most frequent explanations for each output concept, avoiding a sort of “explanation overfitting” with the inclusion of noisy minterms which may correspond to outliers [Simon, 1956].

As a final remark, a possible limitation of LENs can be the readability of logic rules. This may occur when (i) the number of input concepts (the length of any minterm)  $k \gg 1$ , or (ii) the size of the support  $|O_i|$  is very large (possibly getting too many different minterms for any  $f_i$ ). In these scenarios, viable approaches to generate concise logic rules are needed to provide interpretable explanations. More details on how to generate concise explanations are in Section 4.3.

## 4.2 Learning criteria

In this section, we describe some of the loss functions allowing LENs to provide FOL explanations  $\rho_i \in P$ , according to the objectives introduced in Section 3.2 (INTERPRETABLE CLASSIFICATION, EXPLAINING A BLACK-BOX and INTERPRETABLE CLUSTERING). For simplicity, here we will not distinguish among the different computational pipelines, and we will refer to a generic LEN with  $r$  output units. We just saw how a logic rule  $\varphi_i$  can be associated to an output concept  $f_i$ . However to improve the expressiveness of logic explanations, in this section we will promote  $\varphi_i$  to a FOL formula  $\rho_i$ , depending on the selected learning criterion. This follows the usual approach adopted when a model combining logic and machine learning, trained on a finite collection of data, is then applied to out-of-sample inputs, following related studies [Ciravegna et al., 2020b, Ciravegna et al., 2020a, Gnecco et al., 2015]. As a result, any  $\bar{c}_j$  that composes  $\varphi_i$  is thought of as a logic predicate defined on the concept space  $C$ , and such that  $\bar{c}_j(c) = 1$  if and only if  $c_j > 0.5$ , for any  $c \in C$ . In addition, if for instance  $\varphi_i = \bar{c}_2 \wedge \neg\bar{c}_5$ , we will write  $\varphi_i(c)$  for  $\bar{c}_2(c) \wedge \neg\bar{c}_5(c)$ . Then, the specific choice on the loss function that drives the learning criteria of LENs introduces a link between  $\varphi_i(c)$  and the final FOL formulas  $\rho_i$  that is produced by the network.

INTERPRETABLE CLASSIFICATION. Supervised learning is needed to extract explanations for specific categories from LENs. This learning approach binds a logic explanation  $\varphi_i$  to a specific output class of a classification problem. By denoting with  $\bar{y}_i$  the binary predicate associated to the output class  $i$ , we will consider three kinds of FOL explanations

$\rho_i$ , between  $\varphi_i$  and  $\bar{y}_i$ , expressed as the universal closure of an IF, Only IF or IFF rule. These rules can be imposed in case of interpretable classification according to the following learning criteria.

IF rules mean that, in the extreme case, for each sample of class  $i$  we want the  $i$ -th output of the LEN to score 1 (but not necessarily the opposite). In other words, the set of concept tuples  $c$  belonging to the  $i$ -th class, i.e. such that  $\bar{y}_i(c) = 1$  has to be included in the support of  $f_i$ , while no conditions are imposed when  $\bar{y}_i(c) = 0$  (recall that  $f_i(c) \in [0, 1]$ ). This behavior can be achieved by minimizing a hinge loss,

$$L_{\rightarrow}(\bar{y}_i, f_i, \mathcal{C}) = \sum_{c \in \mathcal{C}} \max\{0, \bar{y}_i(c) - f_i(c)\} \quad i \in [1, r]. \quad (4)$$

Following a symmetric approach, in Only IF rules a class is explained in terms of lower-level concepts. This principle is enforced by swapping the two terms in the loss function in Eq. 4,

$$L_{\leftarrow}(\bar{y}_i, f_i, \mathcal{C}) = \sum_{c \in \mathcal{C}} \max\{0, f_i(c) - \bar{y}_i(c)\} \quad i \in [1, r]. \quad (5)$$

Both in IF and Only IF rules, further conditions on  $f$  must be included in order to make the learning problem well posed. To this aim, we need to define the behavior of the model in regions of the concept space not covered by training samples in order to avoid trivial solutions with constant  $f$ . For example,  $f$  could have a fixed bias equal to 1 or no biases at all.

At last, LENs can learn double implication rules (IFF) which completely characterize a certain class. Our experiments are focused on this type of explanations. Any function penalizing points for which  $f_i(c) \neq \bar{y}_i(c)$  may be employed in this scenario, such as the the classic Cross-Entropy loss,

$$L_{\leftrightarrow}(\bar{y}_i, f_i, \mathcal{C}) = \sum_{c \in \mathcal{C}} \bar{y}_i(c) \log(f_i(c)) + (1 - \bar{y}_i(c)) \log(1 - (f_i(c))). \quad (6)$$

For all the above loss functions (Eq. 4-6), LENs provide logic explanations in First-Order Logic by means of the following equations:

$$\text{IF-rule :} \quad \rho_i = \forall c \in \mathcal{C} : \bar{y}_i(c) \rightarrow \varphi_i(c). \quad (7)$$

$$\text{Only IF-rule :} \quad \rho_i = \forall c \in \mathcal{C} : \varphi_i(c) \rightarrow \bar{y}_i(c) \quad (8)$$

$$\text{IFF-rule :} \quad \rho_i = \forall c \in \mathcal{C} : \bar{y}_i(c) \leftrightarrow \varphi_i(c) \quad (9)$$

where each  $\rho_i$ , for  $i = 1, \dots, r$  corresponds to a FOL formula, ranging on the whole concept space  $\mathcal{C}$ , thus generalizing the relationships discovered on the data samples. In the same way, in case of a concept-bottleneck pipeline, a FOL explanation can be derived from the function  $g : X \rightarrow \mathcal{C}$  extracting LEN's input concepts from raw features. For instance for the IFF-rule, we will get

$$\rho_i = \forall x \in X : \bar{y}_i(g(x)) \leftrightarrow \varphi_i(g(x)) \quad (10)$$

We remark that the logic predicate  $\bar{y}_i$  appearing in the Eq. 7-10 is simply a *virtual* predicate denoting the membership of a certain concept tuple  $c \in \mathcal{C}$  to the  $i$ -th output class. Despite the correspondence between  $f_i$  and  $\bar{y}_i$  can be enforced only on the sampled concept space  $\mathcal{C}$ , we assume that any  $\rho_i$  can generalize to unseen concept tuples in the whole concept space  $\mathcal{C}$ , in line with previous works [Ciravegna et al., 2020a, Ciravegna et al., 2020a, Gnecco et al., 2015].

EXPLAINING A BLACK-BOX. In this case, LEN's outputs are forced to mimic the predictions of a black-box  $g : \mathcal{C} \rightarrow Y$ , instead of ground-truth labels. This behaviour can be imposed by considering loss functions analogous to the ones in Eq. 4-6, with  $g_i$  in place of  $\bar{y}_i$ . For instance, for the IFF rule the coherence loss of Eq. 6 can be leveraged, allowing LENs to mimic the behaviour of the black-box,

$$L_{\leftrightarrow}(g_i, f_i, \mathcal{C}) = \sum_{c \in \mathcal{C}} g_i(c) \log(f_i(c)) + (1 - g_i(c)) \log(1 - (f_i(c))) \quad (11)$$

As in the case of interpretable classification, IF, Only IF and IFF rules can be expressed according to Eq. 7-10. However, here FOL explanations will hold assuming that  $\bar{y}_i(c)$  is not a virtual predicate, but it is explicitly associated to the (Booleanized) black-box predictions  $g_i(c)$ .

INTERPRETABLE CLUSTERING. Generic explanations can be obtained by means of fully unsupervised principles we borrow from Information Theory. As a matter of fact, the maximization of the Mutual Information (MI) index between the input concept space  $\mathcal{C}$  and the output concept space  $E$  allows LENs to be trained in a fully unsupervised way [Ciravegna et al., 2020b]. More specifically, a max-MI criterion (see [Melacci and Gori, 2012] for further details) leads to LENs leaning towards 1-hot activation scores, such that  $\forall c \in \mathcal{C}$  only one  $f_i(c) \simeq 1$ , while the others are close

to zero. This encourages LENSs to cluster input data such that each input sample belongs to a single cluster. In order to define the MI index, we have to model the probability distribution of each  $f_i$  to be active (close to 1) on a given sample  $c$ , that we implemented using the softmax operator on LENSs’ outputs. The learning criterion to minimize in order to train LENSs is minus the MI index,

$$L_{MI}(f, \mathcal{C}) = -H_E(f, \mathcal{C}) + H_{E|C}(f, \mathcal{C}), \quad (12)$$

where  $H_E$  and  $H_{E|C}$  denote the entropy and the conditional entropy functions associated to the aforementioned probability distribution, respectively, and measured over the whole  $\mathcal{C}$ , as described by Ciravegna et al. [Ciravegna et al., 2020b]. In this case, the support  $O_i$  of  $f_i$  is exactly the cluster of data points where the  $i$ -th output of the LENS is active. Leaving  $\varphi_i$  free to relate concepts in a purely unsupervised manner, we naturally get the FOL explanation

$$\rho_i = \forall c \in O_i : \varphi_i(c). \quad (13)$$

As a final remark, in all the computational pipelines in which a classifier  $g$  is employed, the available supervision is enforced on  $g$  as well, e.g. by means of the Cross-Entropy loss. As a side note, we mention that what we are describing in a fully supervised setting can be trivially extended to the semi-supervised one, as investigated in previous works [Ciravegna et al., 2020b, Ciravegna et al., 2020a].

### 4.3 Parsimony

When humans compare a set of explanations outlining the same outcomes, they tend to have an implicit bias towards the simplest one [Aristotle, nd, MacKay and Mac Kay, 2003]. In the case of LENSs, the notion of simplicity is implemented by reducing the dependency of each output unit by all of the  $k$  input concepts, encouraging only a subset of them to have a major role in computing LENSs’ outputs. Such subset is of size  $m_i \leq k$  for the  $i$ -th output unit.

Over the years, researchers have proposed many approaches to integrate “*the law of parsimony*” into learning machines. These approaches can be eventually considered as potential solutions to implement parsimony criteria in LENSs, in order to find a valid way to fulfill the end-user requirements on the quality of the explanations and on the classification performance. For instance, Bayesian priors [Wilson, 2020] and weight regularization [Kukačka et al., 2017] are two of the most famous techniques to put in practice the Occam’s razor principle in the fields of statistics and machine learning. Among such techniques,  $L1$ -regularization [Santosa and Symes, 1986] has been recently shown to be quite effective for neural networks providing logic-based explanations [Ciravegna et al., 2020a] as it encourages weight sparsity by shrinking the less important weights to zero [Tibshirani, 1996]. This allows the model to ignore some of the input neurons, that, in the case of the first layer of a LENS, corresponds to discarding or giving negligible weight to some of the input concepts, opening to simplified FOL explanations. If  $W$  collects (a subset of) the weights of the LENS that are subject to this regularization, the learning criterion is then augmented by adding  $\lambda \|W\|_1$ . Notice that the precise weights involved in  $W$  might vary in function of the selected neural architecture to implement the LENS, that is the subject of Section 5. All the out-of-the-box LENSs of the next section are trained using this parsimony criterion, that acts in different portions of the network in the considered instances of LENSs. The parsimony criterion is usually combined with a pruning strategy amongst the ones that are defined in the following.

#### 4.3.1 Pruning strategies

The action of parsimony criteria, such as regularizers or human priors, influences the learning process towards specific local minima. Once the model has finally converged to the desired region of the optimization space, the effort can be speed up and finalized by pruning the neural network [LeCun et al., 1989, Hassibi and Stork, 1993] i.e., removing connections whose likelihood of carrying important information is low. The choice of the connections to be pruned depends on the selected pruning strategy. Such a strategy has a profound impact both on the quality of the explanations but also on the classification performance [Frankle and Carbin, 2018]. Here we present three effective pruning strategies specifically devised for LENSs, whose main goal is to keep FOL formulas compact for each LENS’s output, namely NODE-LEVEL PRUNING, NETWORK-LEVEL PRUNING, EXAMPLE-LEVEL PRUNING.

**NODE-LEVEL PRUNING.** The most “fine-grained” pruning strategy considers each neuron of the network independently. This strategy requires the user to define in advance the maximum *fan-in*  $\zeta \in \mathbb{Z}^+$  for each neuron of a feed-forward neural network i.e., the number of non-pruned incoming connections each neuron can support. In this case, the pruning strategy removes all the connections associated to the smallest weights entering the neuron, until the target fan-in is matched. We refer to this strategy as *node-level pruning*. In detail, the node-level approach prunes one by one the weights with the smallest absolute value, such that each neuron in the network has a fixed number  $\zeta$  of incoming non-pruned weights. The parameter  $\zeta$  determines the computational capabilities of the pruned model as well as the complexity of the logic formulas which can be extracted, as it reduces the number of paths connecting the network inputs to each output neuron. To get simple explanations out of each neuron, the parameter  $\zeta$  may range between 2 and 9

[Miller, 1956, Cowan, 2001, Ma et al., 2014]). Recent work on explainer networks has shown how node-level pruning strategies may lead to fully explainable models [Ciravegna et al., 2020b]. However, we will show in the experimental section that this pruning strategy strongly reduces the classification performances.

**NETWORK-LEVEL PRUNING.** In order to overcome the heavy reduction in learning capacity of node-level pruned models, we introduce the so-called *network-level pruning*. This pruning operation aims at reducing the number of concepts involved in FOL explanations by limiting the availability of input concepts. In detail, the  $L_2$  norm of the connections departing from each input concept is computed. If  $w = [w_1, \dots, w_k]$  is the vector that collects the resulting  $k$  norms, then we re-scale it in the interval  $[0, 1]$ ,

$$w' = \frac{w}{\max_j \{w_j : j = 1, \dots, k\}} \quad (14)$$

where the division is intended to be applied in an element-wise fashion. In this way,  $w'$  gives a normalized score to rank input concepts by importance. The *network-level strategy* consists in pruning all the input features for which  $w'_j < \tau$ , where  $\tau$  is a custom threshold (in our experiments,  $\tau = 0.5$ ). Alternatively, we can retain the  $\zeta$  most important input concepts discard all the others. This can be achieved by pruning all the connections departing from the least relevant concepts similarly to *node-level pruning*. Anyway, this pruning strategy is far less stringent compared to node-level pruning as it affects only the first layer of the network and it does not prescribe a fixed fan-in for each neuron.

**EXAMPLE-LEVEL PRUNING.** Example-level pruning is a particular strategy leveraging the Voronoi tessellation generated by neural networks whose activation functions in all hidden layers are Rectified Linear Units (ReLU Networks) [Hahnloser et al., 2000]. If the LEN is implemented as a ReLU Network, for any input  $c \in C$ , the Directed Acyclic Graph (DAG)  $\mathcal{G}$  describing the structure of the connections in the LEN, can be reduced to  $\mathcal{G}_c$ , which only keeps the units corresponding to active neurons (the ones for which the ReLU activation is non-zero) and the corresponding arcs (referred to as the “firing path”). Since all neurons operate in “linear regime” (affine functions) in the reduced  $\mathcal{G}_c$ , as stated in the following, the input-to-output transformation computed by the multi-layer feed-forward ReLU network with structure  $\mathcal{G}_c$  is a composition of affine functions over the network hidden layers that, in turn, can be simplified with a single affine function, leading to

$$f(c) = \sigma(\hat{W}^{(c)}c + \hat{b}^{(c)}), \quad (15)$$

being  $\sigma$  the activation of the output layer and  $\hat{W}^{(c)}$ ,  $\hat{b}^{(c)}$  a weight matrix and biases (respectively) computed as described in the following.

**Theorem 4.1.** *Let  $\{\xi_1, \dots, \xi_L\}$  be a collection of affine functions, where  $\xi_\kappa : U_\kappa \mapsto V_\kappa : u \mapsto W_\kappa u + b_\kappa$  and  $\forall \kappa = 1, \dots, L-1 : U_{\kappa+1} \subset V_\kappa$ . If a multi-layer network computes the last layer activations as  $\xi_L \circ \xi_{L-1} \circ \dots \circ \xi_2 \circ \xi_1$ , then such transformation is affine and we can re-write it as  $\hat{W}c + \hat{b}$ , where*

$$\hat{W} = \prod_{\kappa=1}^L W_\kappa, \quad \hat{b} = \sum_{\kappa=0}^{L-1} b_{L-\kappa} \prod_{h=0}^{\kappa} W_{L-h+1} \quad (16)$$

and  $W_{L+1} := I$ .

The proof of Theorem 4.1 is straightforward. Such a theorem perfectly applies to the case of ReLU networks once  $\mathcal{G}_c$  has been fixed, since they boil down to simple networks with linear activations. The theorem does not introduce any conditions on how input is represented or on the activation functions in the output layer. Given  $c \in C$ , once  $\mathcal{G}_c$  has been determined, the function computed by the (deep) LEN has the following form,  $f(c) = \sigma(\xi_L^{(c)} \circ \xi_{L-1}^{(c)} \circ \dots \circ \xi_2^{(c)} \circ \xi_1^{(c)}(c))$ , that reduces to  $f(c) = \sigma(\hat{W}^{(c)}c + \hat{b}^{(c)})$ , where the superscript  $(c)$  is added to the symbols of Theorem 4.1 to highlight the fact they are specifically instantiated in the case of the network with structure  $\mathcal{G}_c$ . The reduced form of  $f$  is much easier than considering the one of the original deep network. However,  $\mathcal{G}_c$  might vary for each input  $c$ , thus we might get different transformations for different input samples. Indeed, for each sample we can compute the corresponding  $\mathcal{G}_c$  and, in turn,  $\hat{W}^{(c)}$  and  $\hat{b}^{(c)}$ , that we can prune as described in the case of network-level pruning, keeping only the connections associated with the most important concepts (for that specific sample).

## 5 Out-of-the-box LENS

Crafting state-of-the-art fully interpretable models is not an easy task; rather, there is a trade-off between interpretability and performances. The framework introduced in Section 3, with the methods described in Section 4, is designed to provide the building-blocks to create a wide range of models having different interpretability vs. accuracy trade-offs. Here, we showcase three out-of-the-box neural networks implementing different LENS, whose key properties are about different ways of leveraging the parsimony strategies, as visually anticipated in Fig. 3 (right). In Fig. 8 we sketch

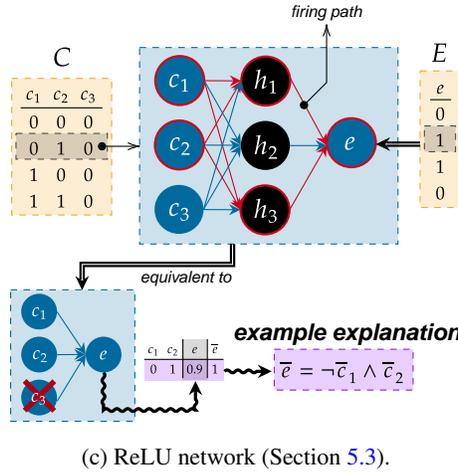
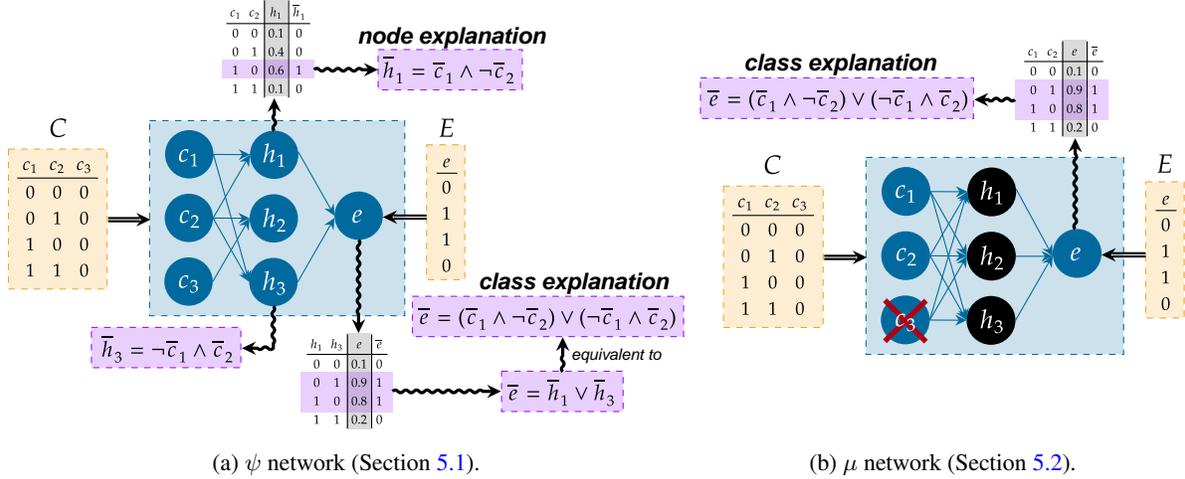


Figure 8: Out-of-the-box LENs showcased in the case of interpretable classification, with examples of logic rules extracted using the procedure of Section 4.1. (a)  $\psi$  networks only admit  $[0, 1]$ -valued neurons with low fan-in, hence we can associate a Boolean formula to each neuron by analysing its I/O truth-Table (b)  $\mu$  networks do not have interpretable hidden neurons but the input concepts are drastically pruned which allows to provide simple network-level explanations. (c) ReLU networks supply explanations for each example by means of the equivalent affine transformation. However, nor the nodes, nor the network can be logically interpreted.

the so-called  $\psi$ ,  $\mu$ , and ReLU out-of-the-box LENs that will be described in the following. Briefly, the  $\psi$  network, originally proposed in [Ciravegna et al., 2020b], is a fully interpretable model with limited learning capacity providing mediocre explanations; the  $\mu$  network is a neural model that can provide high-quality explanations, good learning capacity and modest interpretability; the ReLU network is a model enabling state-of-the-art learning capabilities and good explanations at the cost of very low interpretability. The characteristics of these three LENs are summarized in Table 2.

Table 2: Out-of-the-box LENs and their main properties.

| LEN                  | Pruning       | Activation | Learning  | Explanation | Interpretability |
|----------------------|---------------|------------|-----------|-------------|------------------|
| $\psi$ Net (Fig. 8a) | Node-level    | Sigmoid    | Low       | Low         | Very high        |
| $\mu$ Net (Fig. 8b)  | Network-level | Any        | High      | Very high   | High             |
| ReLU Net (Fig. 8c)   | Example-level | ReLU       | Very high | High        | Low              |

## 5.1 $\psi$ Network

A  $\psi$  network, originally proposed in [Ciravegna et al., 2020b, Ciravegna et al., 2020a], is the first model in the LENs family. A  $\psi$  network is based on three design principles: (i) all activation functions for all neurons (including hidden units) should be sigmoids; (ii) a strong  $L1$ -regularization is used in all the layers to shrink the less important weight values towards zero; (iii) a node-level pruning strategy is considered, where each neuron of the network must have the same number of incoming non-pruned weights (suggested values are between 2 and 9). This number is directly proportional to the number of terms involved in the explanations. In line with previous works [Ciravegna et al., 2020b, Ciravegna et al., 2020a], the rule generation mechanism involves the extraction of rules from each neuron of the network (also the hidden ones), and the node-level pruning favours simple rules on each neuron. Rules are then combined to get the final explanations. The number of hidden layers in the network needs to be small to avoid issues in the rule-merging procedure. In our implementation, pruning is performed after half of the training epochs has been completed, so that the pruned network can refine the value of the remaining weights during the last epochs.

We can cast the  $\psi$  network as a member of the LENs family. The computational pipeline presented in the original work corresponds to a Concept-Bottleneck pipeline. However, from the perspective of the rule-extraction mechanism, the  $\psi$  network can be seen as a Cascading LEN where each layer corresponds to a new level of LEN, each of them generating rules. However, the network is trained as a whole, since only the first layer gets human-understandable concepts. Neuron-specific explanations are then aggregated in order generate final explanations involving input concepts only. For instance, by recalling Example 1, a  $\psi$  network with a single hidden layer with two hidden nodes may learn the logic formula:

$$\varphi_1 = (\neg\bar{c}_1 \wedge \bar{c}_2) \vee (\bar{c}_1 \wedge \neg\bar{c}_2), \quad (17)$$

with hidden nodes learning the formulas  $\bar{h}_1 = \neg\bar{c}_1 \wedge \bar{c}_2$  and  $\bar{h}_2 = \bar{c}_1 \wedge \neg\bar{c}_2$  respectively, and the output node  $f_1$  learning  $\varphi_1 = \bar{h}_1 \vee \bar{h}_2$ .

We remark that the  $\psi$  network is specifically designed to be a fully interpretable neural network. As shown in Fig. 8a, this network provides a high level of interpretability as each neuron can be explained. On the contrary, the strong regularization and the restrictive pruning strategy lead to poor classification accuracy making  $\psi$  networks hard to train and not suitable for solving complex categorical learning problems. Besides, the provided explanations may be disappointing, due to the limited learning capabilities of the model.

## 5.2 $\mu$ Network

A  $\mu$  network is a LEN based on a multi-layer perceptron without the strong constraints on neurons' fan-in of the  $\psi$  net. In particular, a  $\mu$  network is based on two design principles: (i) a strong  $L1$ -regularization is used in the weights connecting the input to the first hidden layer of the network; (ii) a network-level pruning strategy, which prunes input neurons only. After the pruning step, the  $\mu$  network can be fine-tuned to better adapt it to the new configuration (we pruned the network after half of the training epochs has completed). No assumptions need to be made nor on hidden layers nor on activation functions. However, by applying a network-level pruning strategy, the same set of retained input concepts is used for the whole network. On one hand, this allows a simple logic interpretation of the network and, therefore, to provide concise explanations. On the other hand, this may represent a severe limitation for multi-class problems as each class may rely on its own input concepts, but these may be very different among the classes. However,  $\mu$  networks can be efficiently adapted to multi-class experiments (as the ones of Section 6) by splitting the multi-class problem into a set of binary classification problems (one per class), i.e. using a set of light binary  $\mu$  networks.

Still considering Example 1,  $\mu$  networks with a different number of hidden layers may easily learn the same logic formulas. The network-level pruning strategy in this specific example should not prune any of the input features as both  $c_1$  and  $c_2$  are relevant for the target class. However, if we assume the presence of additional redundant features, they would be likely discarded in favor of  $c_1$  and  $c_2$ , as shown in Fig. 8b. Assuming a successful training, the support set  $O_1$  will be composed by the second and third examples of the yellow-box in Fig. 8b, that will yield the following example-level explanations respectively

$$\neg c_1 \wedge c_2 \quad \text{and} \quad c_1 \wedge \neg c_2.$$

The two explanations can be then considered altogether (Eq. 2) to explain the whole class 1 (class-level explanation) by

$$\varphi_1 = \bigvee_{c \in O_1} \varphi_{1,c} = (\neg\bar{c}_1 \wedge \bar{c}_2) \vee (\bar{c}_1 \wedge \neg\bar{c}_2).$$

Thanks to the permissive pruning strategy, the learning capabilities of the network almost match an unconstrained network. As a consequence,  $\mu$  networks are suitable for solving and explaining more complex categorical learning

problems. As we will show in Section 6, the quality of the explanations provided by the  $\mu$  networks are among the highest of the proposed models. At last, a mild-level of interpretability is guaranteed as  $\mu$  nets can be logically interpreted as a whole, but hidden neurons are not as interpretable as in  $\psi$  networks.

### 5.3 ReLU network

The ReLU network is another member of the LENS family providing a different accuracy vs. interpretability trade-off. This model is based on three design principles: (i) all activation functions for all *hidden* neurons are rectified linear units; (ii) a mild L1-regularization is applied to all the weights associated to each layer of the network; (iii) an example-level pruning strategy is used, i.e. a specialized pruning strategy for each sample. Principle (iii) can be applied due to the presence of rectified linear units activation functions. The restriction to ReLU activations is not as limiting as it sounds since it is among the most widely used and efficient activation functions employed in deep learning [Glorot et al., 2011, Ramachandran et al., 2017]. What makes this LEN significantly different from both  $\psi$  and  $\mu$  nets, is that the pruning strategy does not alter the network structure at all. This is due to the fact that pruning is applied to the weights that belong to the single-affine instance of  $f(c)$  of Eq. 15, whose values are collected in  $\hat{W}^{(c)}$  and are only computed for rule-extraction purposes. This means that the original capacity of the model is fully preserved, eventually leading to state-of-the-art classification performances. However, this type of pruning does not provide general insights about the model behaviour, as it is only about the considered example  $c$ , and they may not always lead to optimal explanations.

Recalling again Example 1, a ReLU network with hidden layers can learn the correct logic formula by aggregating different example-level explanations, as we did in the case of the  $\mu$  network. For example, in Fig. 8c we show the case of the processing the second sample from the yellow table, that provides a portion of the explanation of the XOR function. However, when we restrict the connections to the arcs of  $\mathcal{G}_c$  for a certain sample  $c$ , we actually discard several weights of the original ReLU network, i.e. the ones of all the connections that are not needed for classifying the considered  $c$  correctly. This implies that the single-affine instance of  $f(c)$  of Eq. 15, being it function of  $\mathcal{G}_c$ , has a very localized dependence on the space region to which  $c$  belongs. Since Eq. 15 is the form of the LEN from which we extract the example-level explanation, there is the serious risk of obtaining an explanation that does not carry much information from the original structure  $\mathcal{G}$  and that does not globally applies to the whole class. As a matter of fact, there might be strong differences on the example level explanations of samples, even when belonging to the same class.

Summing up, this LEN has the capacity to provide the best performances in terms of classification accuracy thanks to the example-level pruning strategy that do not alter the original network. However, this comes at the cost of poor model interpretability and mild explanation capacity.

## 6 Benchmarking out-of-the-box LENS

In this section, we quantitatively assess the quality of the explanations and the performance of LENS, compared to state-of-the-art white-box models. We consider several tasks, covering multiple combinations of computational pipelines (Section 3.1) and objectives (Section 3.2). The summary of each task is reported in Table 3. A brief description of each task, the corresponding dataset and all the related experimental details are exposed in Section 6.1. In Section 6.2 six quantitative metrics are defined and used to compare LENS with the considered state-of-the-art methods.

Table 3: Summary of the experiments.

| Dataset                           | Description                                     | Pipeline  | Objective                             |
|-----------------------------------|---|-----------|---------------------------------------|
| MIMIC-II (Fig. 4, top)            | Predict patient survival from clinical data     | E2E       | Interpretable classification          |
| MNIST E/O (Fig. 5, top)           | Predict parity from digit images                | CB        | Interpretable classification          |
| CUB (Fig. 2)                      | Predict bird species from bird images           | CB        | Interpretable classification          |
| V-Dem (Fig. 6)                    | Predict electoral democracy from social indexes | Cascading | Interpretable classification          |
| MIMIC-II (EBB)                    | Predict patient survival from clinical data     | E2E       | Explaining black-box (Fig. 4, bottom) |
| MNIST E/O (ICLU) (Fig. 5, bottom) | Cluster digit properties from digit images      | CB        | Interpretable clustering              |

The Python code and the scripts used for the experiments, including parameter values and documentation, is freely available under Apache 2.0 Public License from a GitHub repository<sup>6</sup>. The code is based on our "Logic Explained Networks" library [Barbiero et al., 2021], designed to make out-of-the-box LENs accessible to researchers and neophytes by means of intuitive APIs requiring only a few lines of code to train and get explanations from a LEN, as we sketch in the code example of Listing 1.

```

1 import lens
2
3 # import train, validation and test data loaders
4 [...]
5
6 # instantiate a "psi network"
7 model = lens.models.PsiNN(n_classes=n_classes, n_features=n_features,
8                           hidden_neurons=[200], loss=torch.nn.CrossEntropyLoss(),
9                           l1_weight=0.001, fan_in=10)
10
11 # fit the model
12 model.fit(train_data, val_data, epochs=100, l_r=0.0001)
13
14 # get predictions on test samples
15 outputs, labels = model.predict(test_data)
16
17 # get first-order logic explanations for a specific target class
18 target_class = 1
19 formula = model.get_global_explanation(x_val, y_val, target_class)
20
21 # compute explanation accuracy
22 accuracy = lens.logic.test_explanation(formula, target_class, x_test, y_test)

```

Listing 1: Example on how to use the “Logic Explained Networks” library.

Further details about low-level APIs can be found in Appendix A.

## 6.1 Dataset and classification task details

We considered five categorical learning problems ranging from computer vision to medicine and democracy. Some datasets (e.g. CUB) were already annotated with high-level concepts (e.g. bird attributes) which can be leveraged to train a concept bottleneck pipeline. For datasets without annotations for high-level concepts, we transformed the input space into a predicate space (i.e.  $\mathbb{R}^9 \rightarrow [0, 1]^d$ ) to make it suitable for training LENs. According to the considered data type, we will evaluate and discuss the usage of different LEN pipelines and objectives, as already anticipated in Table 3. For the sake of consistency, in all the experiments based on supervised learning criteria (Section 4.2) LENs were trained by optimizing Eq. 6, therefore extracting IFF rules, that better cope with the ground truth of the datasets. However, as already discussed in Section 4.2, also Eq. 4 and Eq. 5 could be considered, yielding different target rule types. In the following we describe the details of each task.

**MIMIC-II - E2E, INTERPRETABLE CLASSIFICATION.** The Multiparameter Intelligent Monitoring in Intensive Care II (MIMIC-II, [Saeed et al., 2011, Goldberger et al., 2000]) is a public-access intensive care unit (ICU) database consisting of 32,536 subjects (with 40,426 ICU admissions) admitted to ICUs at a single tertiary care hospital. The dataset contains detailed description of a variety of clinical data classes: general, physiological, results of clinical laboratory tests, records of medications, fluid balance, and free text notes and reports of imaging studies (e.g. x-ray, CT, MRI, etc). In our experiments, we removed the text-based input features and we discretized blood pressure (BP) into five different categories (one-hot encoded): very low BP, low BP, normal BP, high BP, and very high BP. After such preprocessing step, we obtained an input space  $X$  composed of 90 key features. The task consists in training a classifier function to identify recovering or dying patients, after 28 days from ICU admission ( $Y = [0, 1]^2$ ). In particular, we considered the case of *interpretable classification* where an End-to-End LEN  $C \rightarrow E$  is employed to directly carry out the classification task, i.e. with  $C = X$  and  $E = Y$ .

**MIMIC-II (EBB) - E2E, EXPLAINING A BLACK-BOX MODEL.** On the same dataset, a second task is set up. A black-box model  $g$  is employed to solve the previously described classification task. The end-to-end LEN  $f$  instead is trained to mimic the behaviour of the black-box  $g$ , using the learning criterion of Eq. 11.

**MNIST E/O - CB, INTERPRETABLE CLASSIFICATION.** The Modified National Institute of Standards and Technology database (MNIST, [LeCun, 1998]) contains a large collection of images representing handwritten digits. The input space  $X \subset \mathbb{R}^{28 \times 28}$  is composed of 28x28 pixel images of digits from 0 to 9. However, the task we aim to solve is slightly different from the common digit-classification. We are interested in determining if a digit is either odd or even, and explain the assignment to one of these classes in terms of the specific digit categories from 0 to 9. In the notation of this

<sup>6</sup>[https://github.com/pietrobarbiero/logic\\_explainer\\_networks](https://github.com/pietrobarbiero/logic_explainer_networks)

paper, we can consider that each image comes with 12 attributes, where the first 10 ones are binary attributes about the specific digit type, i.e. from 0 to 9, while the last 2 ones are binary labels that encode whether the digit is even or odd ( $Y = [0, 1]^{12}$ ). We consider a concept-bottleneck pipeline that consists of a concept space  $C$  that is about the attributes of the specific digit type. We focus on the objective of *interpretable classification* of the odd/even classes, so that the mapping  $X \rightarrow C$ , with  $C = Y^{(1:10)}$  is learned by a ResNet10 classifier  $g$  [He et al., 2016] trained from scratch, while a LEN  $f$  is used to learn both the mapping and the explanation as a function  $C \rightarrow E$ , with  $E = Y^{(11:12)}$ .

**MNIST E/O (ICLU) - CB, INTERPRETABLE CLUSTERING.** The same dataset has been used with the objective of *interpretable clustering*. In this case, we considered a concept space  $C = Y = [0, 1]^{12}$  which comprise both the digits and the even/odd labels. The mapping  $X \rightarrow C$  is learned again by a Resnet10 model  $g$  trained as a multi-label classifier, while the LEN  $f$  performs clustering from  $C \rightarrow E$ , where  $E = [0, 1]^2$  space, in order to extract two clusters that are expected to group data due to two properties of the concept space that are not-defined in advance.

**CUB - CB, INTERPRETABLE CLASSIFICATION.** The Caltech-UCSD Birds-200-2011 dataset (CUB, [Wah et al., 2011]) is a fine-grained classification dataset. It includes 11,788 images representing 200 different bird species. Moreover, 312 lower-level binary attributes have been also attached to each image representing visual characteristics (color, pattern, shape) of particular parts (beak, wings, tail, etc.). Attributes annotations, however, are quite noisy. Similarly to [Koh et al., 2020], we denoised attributes by considering class-level annotations, i.e. a certain attribute is set as present only if it is also present in at least 50% of the images of the same class. Furthermore we only considered attributes present in at least 10 classes (species) after this refinement. In the end, a total of 108 lower-level attributes have been retained and paired with the higher-level attributes about the 200 species, so that each image is represented with binary targets in  $Y = [0, 1]^{200+108=308}$  (where the first 108 targets are for the lower-level attributes). In a concept-bottleneck pipeline with the LEN objective of *interpretable classification*, the mapping  $X \rightarrow C$  from images to lower-level concepts ( $C = Y^{(1:108)}$ ) is performed again with a ResNet10 model  $g$  trained from scratch while a LEN  $f$  learns the final function that classifies the bird specie,  $C \rightarrow E$ , with  $E = Y^{(109:308)}$ .

**V-DEM - CASCADING, INTERPRETABLE CLASSIFICATION.** Varieties of Democracy (V-Dem, [Pemstein et al., 2018, Coppedge et al., 2021]) is a dataset containing a collection of indicators of latent regime characteristics over 202 countries from 1789 to 2020. The database includes 483 low-level indicators (e.g. media bias, party ban, high-court independence, initiatives permitted, etc), 82 mid-level indices (e.g. freedom of expression, freedom of association, equality before the law, etc), and 5 high-level indices of democracy principles (i.e. electoral, liberal, participatory, deliberative, and egalitarian). In the experiments, we considered each example to be paired with low/mid level indices and the information on electoral/non-electoral democracies taken from high level indices. Keeping the same order of the previous description, each sample is then paired with binary targets in  $Y = [0, 1]^{483+82+2=567}$ . We considered the *interpretable classification* objective in the problem of classifying electoral/non-electoral democracies in a cascading LEN with two LEN components, where  $C_1 = Y^{(1:483)}$ ,  $E_1 = C_2 = Y^{(484:566)}$ , and  $E_2 = Y^{(566:567)}$ . In detail, cascading LENs are trained to learn the map  $C_1 \rightarrow E_1$  with  $E_1 = C_2$  and  $C_2 \rightarrow E_2$ , with  $E_2 = Y^{(566:567)}$ . We measured the quality of the rules extracted by the second LEN.

## 6.2 Metrics

Seven metrics are used to compare the performance of LENs with respect to state-of-the-art approaches. While measuring classification metrics is necessary for models to be viable in practice to perform interpretable classification tasks, assessing the quality of the explanations is required to justify their use for explainability. In contrast with other kinds of explanations, logic-based formulas can be evaluated quantitatively. Given a classification problem, first we extract a set of rules from a trained model and then we test/evaluate each explanation on an unseen set of samples. The results for each metric are reported in terms of the mean and standard deviation, computed over a 10-fold cross validation [Krzywinski and Altman, 2013]. Only in the CUB experiments a 5-fold cross validation is performed due to timing issues related to BRL (each fold required about 3 hours)—competitors were described in Section 2. In particular, for each experiment we consider the following metrics.

- *Model accuracy*: it measures how well the LEN or the competitor classifier correctly identifies the target classes, in the case of interpretable classification. When the LEN explains the predictions of a black-box classifier, this metric represents the accuracy of the model in mimicking the black-box classifier (Table 4).
- *Explanation accuracy*: it measures how well the extracted logic formula correctly identifies the target class (Table 5).
- *Complexity of an explanation*: it measures how hard would it be for a human being to understand the logic formula (Table 6). This is simulated by standardizing the explanations in disjunctive normal form and then by counting the number of terms of the standardized formula.

- *Fidelity of an explanation*: it measures how well the predictions obtained by applying the extracted explanations match the predictions obtained when simply using the classifier (Table 7). When the LEN is the classifier itself (i.e. interpretable classification), this metric represents the match between the extracted explanation and the LEN prediction. Instead, when the LEN is explaining the predictions of a black-box classifier, this metric represents the agreement between the extracted explanation and the prediction of black-box classifier.
- *Rule extraction time*: it measures the time required to obtain an explanation from scratch (Fig. 9). It is computed as the sum of the time required to train the model and the time required to extract the formula from a trained model. This is justified by the fact that for some models, like BRL, training and rule extraction consist of just one simultaneous process.
- *Consistency of an explanation*: it measures the similarity of the extracted explanations over different runs (Table 8). It is computed by counting how many times the same concepts appear in the logic formulas over different folds of a 5-fold cross-validation or over 5 different initialization seeds.

Table 4: Model accuracy (%). The two best results are in bold.

|                       | Tree             | BRL                                | $\psi$ net       | ReLU net                           | $\mu$ net                          |
|-----------------------|------------------|------------------------------------|------------------|------------------------------------|------------------------------------|
| <b>MIMIC-II</b>       | 77.53 $\pm$ 1.45 | 76.40 $\pm$ 1.22                   | 77.19 $\pm$ 1.09 | <b>80.11 <math>\pm</math> 1.87</b> | <b>80.00 <math>\pm</math> 0.95</b> |
| <b>vDem</b>           | 85.61 $\pm$ 0.57 | <b>91.23 <math>\pm</math> 0.75</b> | 89.78 $\pm$ 1.64 | <b>92.08 <math>\pm</math> 0.37</b> | 90.40 $\pm$ 0.51                   |
| <b>MNIST E/O</b>      | 99.75 $\pm$ 0.01 | 99.80 $\pm$ 0.02                   | 99.80 $\pm$ 0.03 | <b>99.88 <math>\pm</math> 0.02</b> | <b>99.83 <math>\pm</math> 0.01</b> |
| <b>CUB</b>            | 81.62 $\pm$ 1.17 | 90.79 $\pm$ 0.34                   | 91.92 $\pm$ 0.27 | <b>92.29 <math>\pm</math> 0.40</b> | <b>92.21 <math>\pm</math> 0.33</b> |
| <b>MIMIC-II (EBB)</b> | 77.53 $\pm$ 1.45 | 77.87 $\pm$ 1.24                   | 76.74 $\pm$ 1.52 | <b>80.00 <math>\pm</math> 0.95</b> | <b>79.44 <math>\pm</math> 0.97</b> |

Table 5: Explanation accuracy (%). The two best results are in bold (in case of ties in the average values, we highlighted all the involved models).

|                       | Tree                               | BRL                                | $\psi$ net       | ReLU net                           | $\mu$ net                          |
|-----------------------|------------------------------------|------------------------------------|------------------|------------------------------------|------------------------------------|
| <b>MIMIC-II</b>       | 69.15 $\pm$ 2.24                   | <b>70.59 <math>\pm</math> 2.17</b> | 49.51 $\pm$ 3.92 | 70.28 $\pm$ 1.67                   | <b>71.84 <math>\pm</math> 1.82</b> |
| <b>vDem</b>           | 85.45 $\pm$ 0.58                   | <b>91.21 <math>\pm</math> 0.75</b> | 67.08 $\pm$ 9.68 | <b>90.21 <math>\pm</math> 0.55</b> | 88.18 $\pm$ 1.07                   |
| <b>MNIST E/O</b>      | <b>99.74 <math>\pm</math> 0.01</b> | <b>99.79 <math>\pm</math> 0.02</b> | 65.64 $\pm$ 5.05 | <b>99.74 <math>\pm</math> 0.02</b> | <b>99.74 <math>\pm</math> 0.02</b> |
| <b>CUB</b>            | 89.36 $\pm$ 0.92                   | <b>96.02 <math>\pm</math> 0.17</b> | 76.10 $\pm$ 0.56 | 87.96 $\pm$ 2.81                   | <b>93.69 <math>\pm</math> 0.27</b> |
| <b>MIMIC-II (EBB)</b> | 69.15 $\pm$ 2.24                   | <b>71.68 <math>\pm</math> 2.21</b> | 51.71 $\pm$ 4.78 | 70.53 $\pm$ 1.56                   | <b>71.84 <math>\pm</math> 1.82</b> |

Table 6: Complexity of explanations. The two best results are in bold (the lower, the better).

|                       | Tree                               | BRL                               | $\psi$ net                         | ReLU net                           | $\mu$ net                          |
|-----------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|------------------------------------|
| <b>MIMIC-II</b>       | 66.60 $\pm$ 1.45                   | 57.70 $\pm$ 35.58                 | <b>20.60 <math>\pm</math> 5.36</b> | 39.50 $\pm$ 11.62                  | <b>15.80 <math>\pm</math> 1.37</b> |
| <b>vDem</b>           | 30.20 $\pm$ 1.20                   | 145.70 $\pm$ 57.93                | <b>5.40 <math>\pm</math> 2.70</b>  | 18.40 $\pm$ 2.17                   | <b>9.10 <math>\pm</math> 0.94</b>  |
| <b>MNIST E/O</b>      | <b>47.50 <math>\pm</math> 0.72</b> | 1352.30 $\pm$ 292.62              | 96.90 $\pm$ 10.01                  | <b>73.30 <math>\pm</math> 5.77</b> | 80.50 $\pm$ 4.85                   |
| <b>CUB</b>            | 45.92 $\pm$ 1.16                   | <b>8.87 <math>\pm</math> 0.11</b> | 15.96 $\pm$ 0.96                   | 60.57 $\pm$ 6.95                   | <b>14.65 <math>\pm</math> 0.16</b> |
| <b>MIMIC-II (EBB)</b> | 66.60 $\pm$ 1.45                   | 40.50 $\pm$ 32.46                 | <b>24.30 <math>\pm</math> 5.40</b> | 61.30 $\pm$ 13.61                  | <b>15.80 <math>\pm</math> 1.37</b> |

### 6.3 Results and discussion

Our results are organized in order to compare the behavior of the models in all the considered tasks jointly, with the exception of interpretable clustering, that will be discussed in a separate manner. Experiments of Table 4 show how all the considered out-of-the-box LENs generalize better than decision trees on complex Boolean functions (e.g. CUB) as expected [Tavares et al., 2020], and they usually outperform BRL as well. LENs based on ReLU networks

Table 7: Fidelity of explanations (%). Tree and BRL trivially get 100%. We highlighted in bold the best LEN model.

|                       | Tree              | BRL               | $\psi$ net        | ReLU net                           | $\mu$ net                          |
|-----------------------|-------------------|-------------------|-------------------|------------------------------------|------------------------------------|
| <b>MIMIC-II</b>       | 100.00 $\pm$ 0.00 | 100.00 $\pm$ 0.00 | 51.63 $\pm$ 6.68  | 75.62 $\pm$ 2.07                   | <b>88.37 <math>\pm</math> 2.73</b> |
| <b>vDem</b>           | 100.00 $\pm$ 0.00 | 100.00 $\pm$ 0.00 | 69.67 $\pm$ 10.43 | <b>96.36 <math>\pm</math> 0.64</b> | 94.73 $\pm$ 2.16                   |
| <b>MNIST E/O</b>      | 100.00 $\pm$ 0.00 | 100.00 $\pm$ 0.00 | 65.68 $\pm$ 5.05  | <b>99.85 <math>\pm</math> 0.02</b> | 99.83 $\pm$ 0.02                   |
| <b>CUB</b>            | 100.00 $\pm$ 0.00 | 100.00 $\pm$ 0.00 | 77.34 $\pm$ 0.52  | 89.28 $\pm$ 2.90                   | <b>95.21 <math>\pm</math> 0.25</b> |
| <b>MIMIC-II (EBB)</b> | 100.00 $\pm$ 0.00 | 100.00 $\pm$ 0.00 | 55.72 $\pm$ 8.55  | 74.14 $\pm$ 2.32                   | <b>88.37 <math>\pm</math> 2.73</b> |

Table 8: Rule consistency (%). The two best results are in bold.

|                       | Tree         | BRL           | $\psi$ net | ReLU net      | $\mu$ net     |
|-----------------------|--------------|---------------|------------|---------------|---------------|
| <b>MIMIC-II</b>       | 40.49        | 30.48         | 27.62      | <b>53.75</b>  | <b>71.43</b>  |
| <b>vDem</b>           | <b>72.00</b> | <b>73.33</b>  | 38.00      | 64.62         | 41.67         |
| <b>MNIST E/O</b>      | 41.67        | <b>100.00</b> | 96.00      | <b>100.00</b> | <b>100.00</b> |
| <b>CUB</b>            | 21.47        | 42.86         | 41.43      | <b>44.17</b>  | <b>76.92</b>  |
| <b>MIMIC-II (EBB)</b> | 40.49        | 40.00         | 25.71      | <b>58.67</b>  | <b>71.43</b>  |

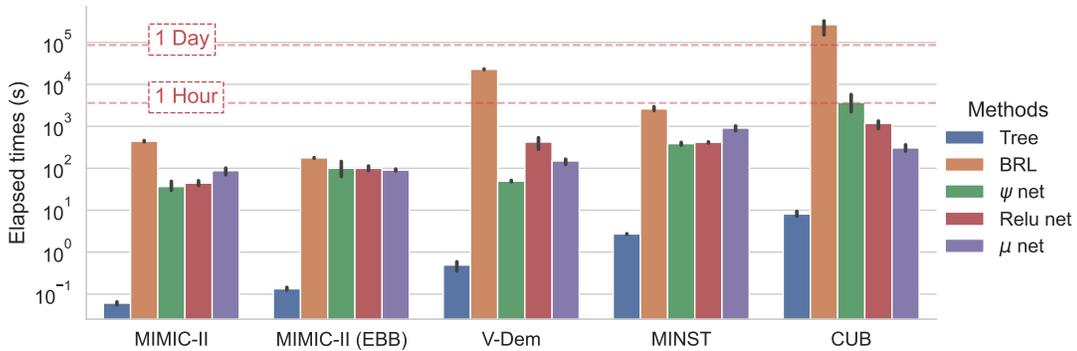


Figure 9: Rule extraction time (seconds). Time required to train models and to extract the explanations. Error bars show the 95% mean confidence interval.

are definitely the ones with better classification accuracy, confirming the intuitions reported in Section 5.3.  $\mu$  nets, however, are quite close in terms of classification accuracy. These results must be paired with the ones of Table 5, since having high classification performance and very low explanation quality would represent a major issue. For most experiments the formulas extracted from LENs are either better or almost as accurate as the formulas found by decision trees or mined by BRL, even if the top performance are reached by BRL. What makes LENs formulas preferable with respect to BRL is the significantly reduced complexity of the considered explanations, as shown in Table 6. Notice how less complex rules implies more readable formulas, that is a crucial feature in the context of Explainable AI. More specifically, the complexity of LENs explanations is usually lower than the complexity of the rules extracted both from a decision tree<sup>7</sup> or mined by BRL. Comparing the results of the different out-of-the-box LENs, we observe that  $\psi$  networks yield moderately complex explanations, sometimes with limited accuracy, while ReLUs, due to the variability of example-level rules, lead to more complex explanations. Overall, the case of  $\mu$  networks is the most attractive one, confirming the quality of their parsimony/pruning criterion.

Moving to more fine-grained details, Table 7 shows how white-box models, like decision trees and BRL, outperform LENs in terms of fidelity. This is due to the fact that such models make predictions based on the explanation directly. Therefore fidelity is (trivially) always 100%. However, we see how the fidelity of the formulas extracted by the  $\mu$

<sup>7</sup>Decision trees have been limited to a maximum of 5 decision levels in order to extract rules of comparable length with the other methods.

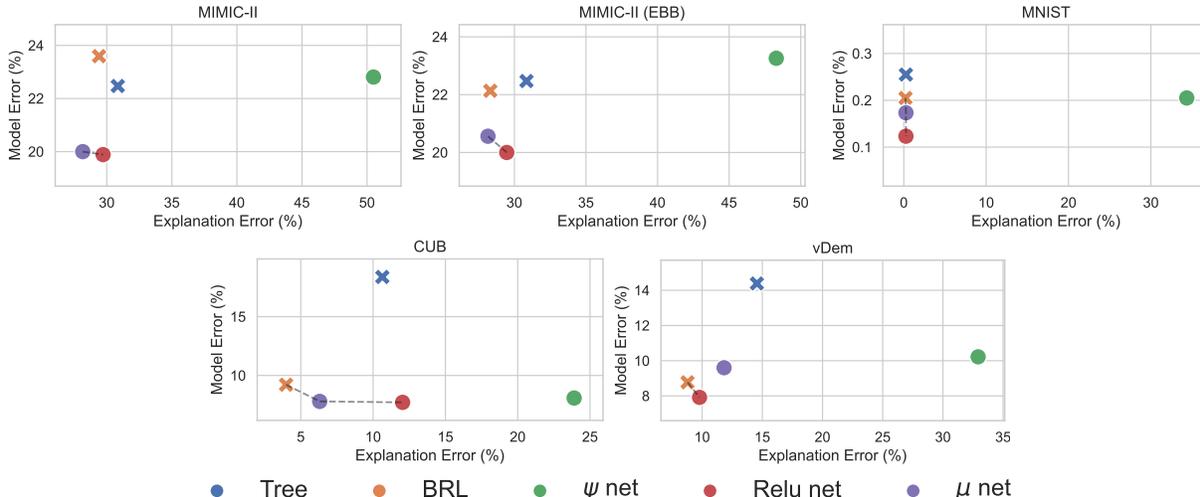


Figure 10: Pareto frontiers (dotted black line) in terms of average model test error and average explanation test error.

and the ReLU network is often higher than 90%. This means that almost any prediction has been correctly explained, making these networks very close to white boxes. Fig. 9 compares the rule extraction times. BRL is the slowest rule extractor over all the experiments, and LENs are faster by one to three orders of magnitude. In two cases BRL takes about 1 hour to extract an explanation and over 1 day in the case of CUB, making it unsuitable for supporting decision making. Decision trees are the fastest overall. In terms of consistency, LENs seem to provide more stable results, but overall there is not a dominant method (see Table 8). We can see, instead, how this result is clearly impacted by the considered dataset/task. Our intuition is that those datasets that are more coherently represented by the data in the different folds are expected to lead to more consistent behaviors.

Fig. 10 shows a combined view of two of the main metrics considered so far, reporting the Pareto frontiers [Marler and Arora, 2004] for each experiment in terms of the explanation and model error (100 minus the explanation/model accuracy). The figure provides a broad perspective showing how the LEN framework is flexible enough to instantiate different models having different compromises between explanation and classification accuracy. The limits of the  $\psi$  network are overcome by the two other out-of-the-box LENs we investigated. In all considered tasks, the  $\mu$  and ReLU LENs are always on the Pareto frontier, therefore providing, on average, the better accuracy vs. interpretability trade-offs.

Our experimental analysis quantitatively demonstrates the quality of LENs and of the FOL rules they extract. In order to provide the reader also a qualitative view on the type of rules that get extracted by the compared models, Table 9 reports a representative subset of the formulas over the different tasks. Comparing the LEN models, we can appreciate the compactness of rules coming from the  $\mu$  and  $\psi$  networks. All the LEN-rules are way more compact than the ones from Trees, and usually also the ones from BRL.

A separate investigation is dedicated to the interpretable clustering objective in MNIST E/O (ICLU). In this scenario, competitors are not involved as they only operate in a supervised manner, while LENs can also handle the unsupervised discovery of FOL formulas. LENs are instructed to learn two clusters considering the data represented in the concept space where the digit identity and the property of being even/odd are encoded. In this setting, we are interested in evaluating whether LENs can discover that data can be grouped into two balanced clusters of even and odd digits, and to provide explanations of them (such as, *odd and not even* and *even and not odd*). We remark that this task does not use any supervisions, and we recall that LENs do not know that odd/even are mutually exclusive properties, but they are just two input concepts out of 12. In Table 10, all the previously introduced metrics are reported for this task, where the model accuracy is computed considering how the system develops clusters that match to the even/odd classes. We can see that all methods are capable of reaching very high level of cluster accuracy, therefore correctly identifying the even and odd groups. By inspecting the extracted rules, we observed that while all the LENs correctly consider odd and even as predicates that participate to the FOL rules, only the  $\psi$  network consistently explains the two clusters in terms of only such important labels (even and odd), as we can see from the complexity of the rules (that is higher than 2 for the  $\mu$  network and ReLU net) and from the consistency of terms used in the explanation (100% for the  $\psi$  network). Other examples and applications of interpretable clustering have been performed in previous works employing the  $\psi$  network only [Ciravegna et al., 2020b, Ciravegna et al., 2020a].

Table 9: A selection of some of the best rules extracted for each experiment. Each rule involves concept names that are taken from the respective dataset. We dropped the argument ( $c$ ) from each predicate to make the notation more compact, and each rule is intended to hold  $\forall c \in C$ . Rule from decision trees are not shown since they involve a very large number of terms.

|                       | Model      | Sample Rule   |
|-----------------------|------------|---|
| <b>MIMIC-II</b>       | $\mu$ net  | Death $\leftrightarrow$ stroke $\wedge$ age_HIGH $\wedge$ $\neg$ atrial_fibrillation  |
|                       | ReLU net   | Death $\leftrightarrow$ stroke $\wedge$ age_HIGH $\wedge$ $\neg$ atrial_fibrillation $\wedge$ $\neg$ sapsi_first_LOW $\wedge$ $\neg$ sapsi_first_HIGH   |
|                       | $\psi$ net | Death $\leftrightarrow$ stroke $\vee$ bun_first_NORMAL $\vee$ sapsi_first_HIGH  |
|                       | Tree       | Death $\leftrightarrow$ *formula is too long  |
|                       | BRL        | Death $\leftrightarrow$ (stroke $\wedge$ resp $\wedge$ $\neg$ age_LOW $\wedge$ $\neg$ sapsi_first_LOW)<br>$\vee$ (stroke $\wedge$ age_HIGH $\wedge$ $\neg$ age_LOW $\wedge$ $\neg$ sapsi_first_LOW)   |
| <b>MIMIC-II (EBB)</b> | $\mu$ net  | Death $\leftrightarrow$ atrial_fibrillation $\wedge$ stroke $\wedge$ $\neg$ sapsi_first_HIGH $\wedge$ $\neg$ weight_first_NORMAL  |
|                       | ReLU net   | Death $\leftrightarrow$ stroke $\wedge$ $\neg$ sapsi_first_LOW $\wedge$ $\neg$ sapsi_first_HIGH $\wedge$ $\neg$ sodium_first_LOW  |
|                       | $\psi$ net | Death $\leftrightarrow$ stroke $\wedge$ age_HIGH  |
|                       | Tree       | Death $\leftrightarrow$ *formula is too long  |
|                       | BRL        | Death $\leftrightarrow$ stroke $\wedge$ age_HIGH $\wedge$ weight_first_LOW $\wedge$ $\neg$ sapsi_first_LOW)<br>$\vee$ (stroke $\wedge$ platelet_first_HIGH $\wedge$ weight_first_LOW $\wedge$ $\neg$ sapsi_first_LOW)   |
| <b>MINIST E/O</b>     | $\mu$ net  | Even $\leftrightarrow$ $\neg$ One $\wedge$ $\neg$ Three $\wedge$ $\neg$ Five $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Nine   |
|                       | ReLU net   | Even $\leftrightarrow$ (Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Two $\wedge$ $\neg$ Three $\wedge$ $\neg$ Four $\wedge$ $\neg$ Five $\wedge$ $\neg$ Six $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Eight $\wedge$ $\neg$ Nine)<br>$\vee$ (Two $\wedge$ $\neg$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Three $\wedge$ $\neg$ Four $\wedge$ $\neg$ Five $\wedge$ $\neg$ Six $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Eight $\wedge$ $\neg$ Nine)<br>$\vee$ (Four $\wedge$ $\neg$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Two $\wedge$ $\neg$ Three $\wedge$ $\neg$ Five $\wedge$ $\neg$ Six $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Eight $\wedge$ $\neg$ Nine)<br>$\vee$ (Six $\wedge$ $\neg$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Two $\wedge$ $\neg$ Three $\wedge$ $\neg$ Four $\wedge$ $\neg$ Five $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Eight $\wedge$ $\neg$ Nine)<br>$\vee$ (Eight $\wedge$ $\neg$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Two $\wedge$ $\neg$ Three $\wedge$ $\neg$ Four $\wedge$ $\neg$ Five $\wedge$ $\neg$ Six $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Nine)  |
|                       | $\psi$ net | Even $\leftrightarrow$ (Six $\wedge$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Seven) $\vee$ (Six $\wedge$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Three)<br>$\vee$ (Six $\wedge$ Zero $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Three) $\vee$ (Six $\wedge$ $\neg$ One $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Three)<br>$\vee$ (Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Three) $\vee$ (Six $\wedge$ Zero $\wedge$ $\neg$ One $\wedge$ $\neg$ Seven $\wedge$ $\neg$ Three)  |
|                       | Tree       | Even $\leftrightarrow$ *formula is too long   |
|                       | BRL        | Even $\leftrightarrow$ (Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Five $\wedge$ $\neg$ Two)) $\vee$ (Eight $\wedge$ $\neg$ Nine $\wedge$ $\neg$ Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Five $\wedge$ $\neg$ Two))<br>$\vee$ (Four $\wedge$ $\neg$ Nine $\wedge$ $\neg$ Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Eight $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Five $\wedge$ $\neg$ Two)<br>$\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Two)) $\vee$ (Two $\wedge$ $\neg$ One $\wedge$ $\neg$ Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Eight $\wedge$ $\neg$ Nine)<br>$\wedge$ $\neg$ (Five $\wedge$ $\neg$ Two) $\wedge$ $\neg$ (Four $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Nine)<br>$\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Two)) $\vee$ (Four $\wedge$ $\neg$ Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Eight $\wedge$ $\neg$ Nine)<br>$\wedge$ $\neg$ (Five $\wedge$ $\neg$ Two) $\wedge$ $\neg$ (Four $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Two) $\wedge$<br>$\neg$ (Two $\wedge$ $\neg$ One)) $\vee$ (Zero $\wedge$ $\neg$ Four $\wedge$ $\neg$ Nine $\wedge$ $\neg$ Six $\wedge$ $\neg$ Three $\wedge$ $\neg$ (Eight $\wedge$ $\neg$ Nine) $\wedge$<br>$\neg$ (Five $\wedge$ $\neg$ Two) $\wedge$ $\neg$ (Four $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Nine $\wedge$ $\neg$ Zero) $\wedge$ $\neg$ (One $\wedge$ $\neg$ Two)<br>$\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Nine) $\wedge$ $\neg$ (Seven $\wedge$ $\neg$ Two) $\wedge$ $\neg$ (Two $\wedge$ $\neg$ One)) |
| <b>CUB</b>            | $\mu$ net  | Black_foot_albatross $\leftrightarrow$ bill_shape_hooked_seabird $\wedge$ size_medium $\wedge$ wing_pattern_solid<br>$\wedge$ $\neg$ wing_color_black $\wedge$ $\neg$ underparts_color_white $\wedge$ $\neg$ upper_tail_color_grey<br>$\wedge$ $\neg$ breast_color_white $\wedge$ $\neg$ throat_color_white $\wedge$ $\neg$ tail_pattern_solid<br>$\wedge$ $\neg$ crown_color_white   |
|                       | ReLU net   | Black_foot_albatross $\leftrightarrow$ size_medium $\wedge$ $\neg$ bill_shape_allpurpose<br>$\wedge$ $\neg$ upperparts_color_black $\wedge$ $\neg$ head_pattern_plain<br>$\wedge$ $\neg$ under_tail_color_black $\wedge$ $\neg$ nape_color_buff<br>$\wedge$ $\neg$ wing_shape_roundedwings $\wedge$ $\neg$ shape_perchinglike<br>$\wedge$ $\neg$ leg_color_grey $\wedge$ $\neg$ leg_color_black $\wedge$ $\neg$ bill_color_grey<br>$\wedge$ $\neg$ bill_color_black $\wedge$ $\neg$ wing_pattern_multicolored   |
|                       | $\psi$ net | Black_foot_albatross $\leftrightarrow$ (bill_shape_hooked_seabird $\wedge$ tail_pattern_solid $\wedge$ $\neg$ underparts_color_white<br>$\wedge$ ( $\neg$ breast_color_white $\vee$ $\neg$ wing_color_grey))  |
|                       | Tree       | Black_foot_albatross $\leftrightarrow$ *formula is too long   |
|                       | BRL        | Black_foot_albatross $\leftrightarrow$ (bill_shape_hooked_seabird $\wedge$ forehead_color_blue $\wedge$ $\neg$ bill_color_black<br>$\wedge$ $\neg$ nape_color_white) $\vee$ (bill_shape_hooked_seabird $\wedge$ $\neg$ bill_color_black<br>$\wedge$ $\neg$ nape_color_white $\wedge$ $\neg$ tail_pattern_solid)   |
| <b>V-Dem</b>          | $\mu$ net  | Elect_Dem $\leftrightarrow$ v2x_freexp_alfinf $\wedge$ v2x_frassoc_thick $\wedge$ v2xel_frefair $\wedge$ v2x_elecoff $\wedge$ v2xcl_rol   |
|                       | ReLU net   | Elect_Dem $\leftrightarrow$ v2x_freexp_alfinf $\wedge$ v2xel_frefair $\wedge$ v2x_elecoff $\wedge$ v2xcl_rol  |
|                       | $\psi$ net | Elect_Dem $\leftrightarrow$ v2x_frassoc_thick $\wedge$ v2xeg_eqaccess   |
|                       | Tree       | Elect_Dem $\leftrightarrow$ *formula is too long  |
|                       | BRL        | Elect_Dem $\leftrightarrow$ v2x_freexp_alfinf $\wedge$ v2x_frassoc_thick $\wedge$ v2xel_frefair $\wedge$ v2x_elecoff $\wedge$ v2xcl_rol   |

Table 10: MNIST E/O (ICLU). All the metrics are reported.

| Method     | Model Acc        | Exp. Acc. (%)    | Complexity      | Extr. Time (sec) | Consistency (%) |
|------------|------------------|------------------|-----------------|------------------|-----------------|
| $\psi$ net | $99.90 \pm 0.01$ | $99.90 \pm 0.01$ | $2.00 \pm 0.00$ | $0.33 \pm 0.15$  | 100.00          |
| $\mu$ net  | $94.91 \pm 4.99$ | $95.41 \pm 4.50$ | $2.30 \pm 0.64$ | $2.36 \pm 0.75$  | 35.00           |
| ReLU net   | $96.87 \pm 3.04$ | $95.97 \pm 3.94$ | $2.15 \pm 0.39$ | $4.15 \pm 1.62$  | 45.00           |

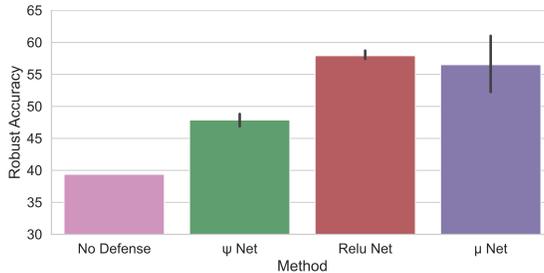


Figure 11: Quality of the explanations extracted by the LENSs in terms of robust accuracy when used to defend a model from adversarial examples. In pink, the accuracy of the model when facing a white box attack without any defense. Error bars show the 95% confidence interval of the mean.

#### 6.4 LENSs as Adversarial Defense

In nowadays machine learning literature, there is a serious concerns about the vulnerability of several machine learning models, such as neural networks, to the so called *adversarial examples*. These are input samples maliciously altered with a slight perturbation that may lead to wrong predictions, even if they do not look to be visually altered for a human (in the case of images). This kind of issue has received a lot of attention, and several work has been done to develop some techniques, *adversarial defenses*, to prevent an AI model from fraudulent *adversarial attacks* [Miller et al., 2020, Ozdag, 2018]. Recently, it has been shown how explicit domain knowledge, expressed by a set of FOL formulas, can be used to discriminate if a certain sample has to be considered as adversarial [Melacci et al., 2021], especially in case of multi-label classification. However this approach requires that the logic rules are already available for the considered task, so that a domain expert is expected to collect them.

Interestingly, LENSs can be applied to a set of clean data to automatically learn a set of FOL rules that may capture the important properties of the target domain. Then, these rules can be used to detect adversarial samples as in the framework introduced in [Melacci et al., 2021], since they are exactly of the same type of the LENS ones.

Without any attempts of being exhaustive on this topic, we briefly explored whether what we stated in the previous sentence would work in practice. We considered the bird-identification dataset (CUB 200), facing the classification problem in its original multi-label fashion, where a convolutional neural network  $g$  (ResNet18) is trained from scratch to classify both the 200 classes and the ones of the additional bird attributes. The FOL formulas automatically extracted by the LENSs in the already discussed experimental experience are directly plugged in the rejection mechanism of [Melacci et al., 2021], thus making  $g$  equipped with a rejection function. We fixed the adversarial perturbation upper-bound to  $\epsilon = 0.5$ , and we generated perturbed data attacking the 200 classes with the state-of-the-art attack APGD-CE [Croce and Hein, 2020], measuring the adversarial-aware accuracy described in [Melacci et al., 2021], that here we refer to as robust accuracy. The experimental results for the three out-of-the-box LENSs introduced in Section 5 are shown in Fig. 11. Interestingly, the accuracy of the attacked  $g$  classifier increases in a significant manner using the rules extracted by LENSs, confirming the validity of the proposed approach. The ReLU net leads to more detailed rules that better defend  $g$ , even if they introduce a larger computational burden in the defense mechanism, since they are usually more complex than the ones of the  $\mu$  net.

## 7 Conclusions and Future Work

In this paper we presented a family of neural networks called Logic Explained Networks (LENSs). We presented an in-depth study on the idea of using a neural model either to provide explanations for black-box or to solve a

classification problem in an interpretable manner. Explanations are provided by First-Order Logic formulas, whose type is strongly interconnected with the learning criterion that drives LENSs’ training. Our framework covers a large set of computational pipelines and different user objectives. We investigated and experimented the case of three out-of-the-box LENS models, showing that they represent a balanced trade-off among a set of important properties, comparing favourably with state-of-the-art white-box classifiers.

The extraction of a first-order logic explanation requires symbolic input and output spaces. This constraint is the main limitation of our framework, as it narrows the range of applications down to symbolic input/output problems. In some contexts, such as computer vision, the use of LENSs may require additional annotations and attribute labels to get a consistent symbolic layer of concepts. However, recent work may partially solve this issue leading to more cost-effective concept annotations [Ghorbani et al., 2019, Kazhdan et al., 2020]. Another area to focus on might be the improvement of out-of-the-box LENSs. The efficiency and the classification performances of fully interpretable LENSs, i.e.  $\psi$  network, is still quite limited due to the extreme pruning strategy. Even more flexible approaches, like  $\mu$  networks, are not as efficient as standard neural models when tackling multi-class or multi-label classification problems, as they require a bigger architecture to provide a disentangled explanation for each class. In our vision this framework would thrive and express its full potential by interacting with the external world and especially with humans in a dynamic way. In this case, logic formulas might be rewritten as sentences to allow for a more natural interaction with end users. Moreover, a dynamic interaction may call for an extended expressivity leading to higher-order or temporal logic explanations. Finally, in some contexts different neural architectures as graph neural networks might be worth exploring as they may be more suitable to solve the classification problem.

Current legislation in US and Europe binds AI to provide explanations especially when the economical, ethical, or financial impact is significant [EUGDPR, 2017, Law, 10]. This work contributes to a lawful and safer adoption of some of the most powerful AI technologies allowing deep neural networks to have a greater impact on society. The formal language of logic provides clear and synthetic explanations, suitable for laypeople, managers, and in general for decision makers outside the AI research field. The experiments show how this framework can be used to aid bias identification and to make black-boxes more robust to adversarial attacks. As out-of-the-box LENSs are easy to use even for neophytes and can be effectively used to understand the behavior of an existing algorithm, our approach might be used to reverse engineer competitor products, to find vulnerabilities, or to improve system design. From a scientific perspective, formal knowledge distillation from state-of-the-art networks may enable scientific discoveries or confirmation of existing theories.

## Acknowledgments

We thank Ben Day, Dobrik Georgiev, Dmitry Kazhdan, and Alberto Tonda for useful feedback and suggestions.

This work was partially supported by the GO-DS21 project funded by the European Union’s Horizon 2020 research and innovation programme under grant agreement No 848077 and by the PRIN 2017 project RexLearn (Reliable and Explainable Adversarial Machine Learning), funded by the Italian Ministry of Education, University and Research (grant no. 2017TWNMH2). This work was also partially supported by TAILOR, a project funded by EU Horizon 2020 research and innovation programme under GA No 952215.

## References

- [Adadi and Berrada, 2018] Adadi, A. and Berrada, M. (2018). Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160. 5
- [Agarwal et al., 2020] Agarwal, R., Frosst, N., Zhang, X., Caruana, R., and Hinton, G. E. (2020). Neural additive models: Interpretable machine learning with neural nets. *arXiv preprint arXiv:2004.13912*. 6
- [Ahmad et al., 2018] Ahmad, M. A., Eckert, C., and Teredesai, A. (2018). Interpretable machine learning in health-care. In *Proceedings of the 2018 ACM international conference on bioinformatics, computational biology, and health informatics*, pages 559–560. 2
- [Angelino et al., 2018] Angelino, E., Larus-Stone, N., Alabi, D., Seltzer, M., and Rudin, C. (2018). Learning certifiably optimal rule lists for categorical data. 5
- [Aristotle, nd] Aristotle (n.d.). Posterior analytics. 15
- [Barbiero et al., 2021] Barbiero, P., Ciravegna, G., Georgiev, D., and Giannini, F. (2021). Pytorch, explain! a python library for logic explained networks. *arXiv preprint arXiv:2105.11697*. 20, 31

- [Battaglia et al., 2018] Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., et al. (2018). Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*. 2
- [Breiman et al., 1984] Breiman, L., Friedman, J., Stone, C. J., and Olshen, R. A. (1984). *Classification and regression trees*. CRC press. 2, 5, 6
- [Brundage et al., 2020] Brundage, M., Avin, S., Wang, J., Belfield, H., Krueger, G., Hadfield, G., Khlaaf, H., Yang, J., Toner, H., Fong, R., et al. (2020). Toward trustworthy ai development: mechanisms for supporting verifiable claims. *arXiv preprint arXiv:2004.07213*. 1
- [Caruana et al., 2015] Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., and Elhadad, N. (2015). Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1721–1730. 6
- [Carvalho et al., 2019] Carvalho, D. V., Pereira, E. M., and Cardoso, J. S. (2019). Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8(8):832. 1, 2, 5
- [Chander et al., 2018] Chander, A., Srinivasan, R., Chelian, S., Wang, J., and Uchino, K. (2018). Working with beliefs: Ai transparency in the enterprise. In *IUI Workshops*. 1
- [Ciravegna et al., 2020a] Ciravegna, G., Giannini, F., Gori, M., Maggini, M., and Melacci, S. (2020a). Human-driven fol explanations of deep learning. In *Twenty-Ninth International Joint Conference on Artificial Intelligence and Seventeenth Pacific Rim International Conference on Artificial Intelligence {IJCAI-PRICAI-20}*, pages 2234–2240. International Joint Conferences on Artificial Intelligence Organization. 3, 10, 13, 14, 15, 18, 24
- [Ciravegna et al., 2020b] Ciravegna, G., Giannini, F., Melacci, S., Maggini, M., and Gori, M. (2020b). A constraint-based approach to learning and explanation. In *AAAI*, pages 3658–3665. 3, 13, 14, 15, 16, 17, 18, 24
- [Cohen, 1995] Cohen, W. W. (1995). Fast effective rule induction. In *Machine learning proceedings 1995*, pages 115–123. Elsevier. 6
- [Coppedge et al., 2021] Coppedge, M., Gerring, J., Knutsen, C. H., Lindberg, S. I., Teorell, J., Altman, D., Bernhard, M., Cornell, A., Fish, M. S., Gastaldi, L., et al. (2021). V-dem codebook v11. 21
- [Cowan, 2001] Cowan, N. (2001). The magical number 4 in short-term memory: A reconsideration of mental storage capacity. *Behavioral and brain sciences*, 24(1):87–114. 16
- [Cranmer et al., 2019] Cranmer, M. D., Xu, R., Battaglia, P., and Ho, S. (2019). Learning symbolic physics with graph networks. *arXiv preprint arXiv:1909.05862*. 2
- [Croce and Hein, 2020] Croce, F. and Hein, M. (2020). Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. 26
- [Das and Rad, 2020] Das, A. and Rad, P. (2020). Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371*. 1, 2
- [Devlin et al., 2018] Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*. 2
- [Doshi-Velez and Kim, 2017] Doshi-Velez, F. and Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*. 2
- [Doshi-Velez and Kim, 2018] Doshi-Velez, F. and Kim, B. (2018). Considerations for evaluation and generalization in interpretable machine learning. In *Explainable and interpretable models in computer vision and machine learning*, pages 3–17. Springer. 2
- [Dosovitskiy et al., 2020] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. 2
- [Erhan et al., 2010] Erhan, D., Courville, A., and Bengio, Y. (2010). Understanding representations learned in deep architectures. *Department dInformatique et Recherche Operationnelle, University of Montreal, QC, Canada, Tech. Rep*, 1355(1). 5, 6
- [EUGDPR, 2017] EUGDPR (2017). Gdpr. general data protection regulation. 1, 27
- [Frankle and Carbin, 2018] Frankle, J. and Carbin, M. (2018). The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*. 15
- [Ghorbani et al., 2019] Ghorbani, A., Wexler, J., Zou, J., and Kim, B. (2019). Towards automatic concept-based explanations. *arXiv preprint arXiv:1902.03129*. 2, 27

- [Gilpin et al., 2018] Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., and Kagal, L. (2018). Explaining explanations: An overview of interpretability of machine learning. In *2018 IEEE 5th International Conference on data science and advanced analytics (DSAA)*, pages 80–89. IEEE. 1
- [Glorot et al., 2011] Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier neural networks. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 315–323. JMLR Workshop and Conference Proceedings. 19
- [Gnecco et al., 2015] Gnecco, G., Gori, M., Melacci, S., and Sanguineti, M. (2015). Foundations of support constraint machines. *Neural computation*, 27(2):388–480. 13, 14
- [Goddard, 2017] Goddard, M. (2017). The eu general data protection regulation (gdpr): European regulation that has a global impact. *International Journal of Market Research*, 59(6):703–705. 1
- [Goldberger et al., 2000] Goldberger, A. L., Amaral, L. A., Glass, L., Hausdorff, J. M., Ivanov, P. C., Mark, R. G., Mietus, J. E., Moody, G. B., Peng, C.-K., and Stanley, H. E. (2000). Physiobank, physiotookit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220. 20
- [Guidotti et al., 2018] Guidotti, R., Monreale, A., Ruggieri, S., Pedreschi, D., Turini, F., and Giannotti, F. (2018). Local rule-based explanations of black box decision systems. *arXiv preprint arXiv:1805.10820*. 6
- [Gunning, 2017] Gunning, D. (2017). Explainable artificial intelligence (xai). *Defense Advanced Research Projects Agency (DARPA), nd Web*, 2(2). 1
- [Hahnloser et al., 2000] Hahnloser, R. H., Sarpeshkar, R., Mahowald, M. A., Douglas, R. J., and Seung, H. S. (2000). Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. *Nature*, 405(6789):947–951. 16
- [Hassibi and Stork, 1993] Hassibi, B. and Stork, D. G. (1993). *Second order derivatives for network pruning: Optimal brain surgeon*. Morgan Kaufmann. 15
- [Hastie and Tibshirani, 1987] Hastie, T. and Tibshirani, R. (1987). Generalized additive models: some applications. *Journal of the American Statistical Association*, 82(398):371–386. 6
- [He et al., 2016] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778. 21
- [Kazhdan et al., 2020] Kazhdan, D., Dimanov, B., Jamnik, M., Liò, P., and Weller, A. (2020). Now you see me (cme): Concept-based model extraction. *arXiv preprint arXiv:2010.13233*. 27
- [Kim et al., 2018] Kim, B., Gilmer, J., Wattenberg, M., and Viégas, F. (2018). Tcav: Relative concept importance testing with linear concept activation vectors. 2, 7
- [Koh et al., 2020] Koh, P. W., Nguyen, T., Tang, Y. S., Mussmann, S., Pierson, E., Kim, B., and Liang, P. (2020). Concept bottleneck models. In *International Conference on Machine Learning*, pages 5338–5348. PMLR. 2, 3, 21
- [Krzywinski and Altman, 2013] Krzywinski, M. and Altman, N. (2013). Error bars: the meaning of error bars is often misinterpreted, as is the statistical significance of their overlap. *Nature methods*, 10(10):921–923. 21
- [Kukačka et al., 2017] Kukačka, J., Golkov, V., and Cremers, D. (2017). Regularization for deep learning: A taxonomy. *arXiv preprint arXiv:1710.10686*. 15
- [Law, 10] Law, P. A. (10). Code of federal regulations. *Wash.: Gov. print. off.* 1, 27
- [LeCun, 1998] LeCun, Y. (1998). The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>. 20
- [LeCun et al., 1989] LeCun, Y., Denker, J. S., Solla, S. A., Howard, R. E., and Jackel, L. D. (1989). Optimal brain damage. In *NIPs*, volume 2, pages 598–605. Citeseer. 15
- [Letham et al., 2015] Letham, B., Rudin, C., McCormick, T. H., Madigan, D., et al. (2015). Interpretable classifiers using rules and bayesian analysis: Building a better stroke prediction model. *Annals of Applied Statistics*, 9(3):1350–1371. 2, 5, 6
- [Lipton, 2018] Lipton, Z. C. (2018). The mythos of model interpretability: In machine learning, the concept of interpretability is both important and slippery. *Queue*, 16(3):31–57. 1, 2
- [Lundberg and Lee, 2017] Lundberg, S. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. *arXiv preprint arXiv:1705.07874*. 5, 6
- [Ma et al., 2014] Ma, W. J., Husain, M., and Bays, P. M. (2014). Changing concepts of working memory. *Nature neuroscience*, 17(3):347. 16
- [MacKay and Mac Kay, 2003] MacKay, D. J. and Mac Kay, D. J. (2003). *Information theory, inference and learning algorithms*. Cambridge university press. 15

- [Marcinkevičs and Vogt, 2020] Marcinkevičs, R. and Vogt, J. E. (2020). Interpretability and explainability: A machine learning zoo mini-tour. *arXiv preprint arXiv:2012.01805*. 1, 5
- [Marler and Arora, 2004] Marler, R. T. and Arora, J. S. (2004). Survey of multi-objective optimization methods for engineering. *Structural and multidisciplinary optimization*, 26(6):369–395. 24
- [McCluskey, 1956] McCluskey, E. J. (1956). Minimization of boolean functions. *The Bell System Technical Journal*, 35(6):1417–1444. 2
- [McColl, 1878] McColl, H. (1878). The calculus of equivalent statements (third paper). *Proceedings of the London Mathematical Society*, 1(1):16–28. 2
- [Melacci et al., 2021] Melacci, S., Ciravegna, G., Sotgiu, A., Demontis, A., Biggio, B., Gori, M., and Roli, F. (2021). Domain knowledge alleviates adversarial attacks in multi-label classifiers. *arXiv*, 2006.03833. 26
- [Melacci and Gori, 2012] Melacci, S. and Gori, M. (2012). Unsupervised learning by minimal entropy encoding. *IEEE transactions on neural networks and learning systems*, 23(12):1849–1861. 14
- [Mendelson, 2009] Mendelson, E. (2009). *Introduction to mathematical logic*. CRC press. 12
- [Miller et al., 2020] Miller, D. J., Xiang, Z., and Kesidis, G. (2020). Adversarial learning targeting deep neural network classification: A comprehensive review of defenses against attacks. *Proceedings of the IEEE*, 108(3):402–433. 26
- [Miller, 1956] Miller, G. A. (1956). The magical number seven, plus or minus two: Some limits on our capacity for processing information. *Psychological review*, 63:81–97. 16
- [Miller, 2019] Miller, T. (2019). Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38. 2
- [Molnar, 2020] Molnar, C. (2020). *Interpretable machine learning*. Lulu. com. 2, 5
- [Ozdag, 2018] Ozdag, M. (2018). Adversarial attacks and defenses against deep neural networks: a survey. *Procedia Computer Science*, 140:152–161. 26
- [Pemstein et al., 2018] Pemstein, D., Marquardt, K. L., Tzelgov, E., Wang, Y.-t., Krusell, J., and Miri, F. (2018). The v-dem measurement model: latent variable analysis for cross-national and cross-temporal expert-coded data. *V-Dem Working Paper*, 21. 21
- [Quine, 1952] Quine, W. V. (1952). The problem of simplifying truth functions. *The American mathematical monthly*, 59(8):521–531. 2
- [Quinlan, 1987] Quinlan, J. R. (1987). Simplifying decision trees. *International journal of man-machine studies*, 27(3):221–234. 6
- [Ramachandran et al., 2017] Ramachandran, P., Zoph, B., and Le, Q. V. (2017). Searching for activation functions. *arXiv preprint arXiv:1710.05941*. 19
- [Ribeiro et al., 2016a] Ribeiro, M. T., Singh, S., and Guestrin, C. (2016a). "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. 5, 6
- [Ribeiro et al., 2016b] Ribeiro, M. T., Singh, S., and Guestrin, C. (2016b). Model-agnostic interpretability of machine learning. *arXiv preprint arXiv:1606.05386*. 5
- [Ribeiro et al., 2018] Ribeiro, M. T., Singh, S., and Guestrin, C. (2018). Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32. 6
- [Rudin, 2019] Rudin, C. (2019). Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215. 1, 2
- [Rudin et al., 2021] Rudin, C., Chen, C., Chen, Z., Huang, H., Semenova, L., and Zhong, C. (2021). Interpretable machine learning: Fundamental principles and 10 grand challenges. *arXiv preprint arXiv:2103.11251*. 2
- [Russell and Norvig, 2016] Russell, S. J. and Norvig, P. (2016). *Artificial intelligence: a modern approach*. Malaysia; Pearson Education Limited,. 13
- [Saeed et al., 2011] Saeed, M., Villarroel, M., Reisner, A. T., Clifford, G., Lehman, L.-W., Moody, G., Heldt, T., Kyaw, T. H., Moody, B., and Mark, R. G. (2011). Multiparameter intelligent monitoring in intensive care ii (mimic-ii): a public-access intensive care unit database. *Critical care medicine*, 39(5):952. 20
- [Samek et al., 2020] Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J., and Müller, K.-R. (2020). Toward interpretable machine learning: Transparent deep neural networks and beyond. *arXiv preprint arXiv:2003.07631*. 2

- [Santosa and Symes, 1986] Santosa, F. and Symes, W. W. (1986). Linear inversion of band-limited reflection seismograms. *SIAM Journal on Scientific and Statistical Computing*, 7(4):1307–1330. 15
- [Sato and Tsukimoto, 2001] Sato, M. and Tsukimoto, H. (2001). Rule extraction from neural networks via decision tree induction. In *IJCNN'01. International Joint Conference on Neural Networks. Proceedings (Cat. No. 01CH37222)*, volume 3, pages 1870–1875. IEEE. 6
- [Schmidt and Lipson, 2009] Schmidt, M. and Lipson, H. (2009). Distilling free-form natural laws from experimental data. *science*, 324(5923):81–85. 2
- [Selvaraju et al., 2017] Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., and Batra, D. (2017). Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626. 5
- [Simon, 1956] Simon, H. A. (1956). Rational choice and the structure of the environment. *Psychological review*, 63(2):129. 13
- [Simon, 1979] Simon, H. A. (1979). Rational decision making in business organizations. *The American economic review*, 69(4):493–513. 2
- [Simonyan et al., 2013] Simonyan, K., Vedaldi, A., and Zisserman, A. (2013). Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034*. 5, 6
- [Srinivasan and Chander, 2020] Srinivasan, R. and Chander, A. (2020). Explanation perspectives from the cognitive sciences—a survey. In *29th International Joint Conference on Artificial Intelligence*, pages 4812–4818. 2
- [Tavares et al., 2020] Tavares, A. R., Avelar, P., Flach, J. M., Nicolau, M., Lamb, L. C., and Vardi, M. (2020). Understanding boolean function learnability on deep neural networks. *arXiv preprint arXiv:2009.05908*. 4, 5, 22
- [Tibshirani, 1996] Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288. 15
- [Wah et al., 2011] Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. (2011). The caltech-ucsd birds-200-2011 dataset. 21
- [Wilson, 2020] Wilson, A. G. (2020). The case for bayesian deep learning. *arXiv preprint arXiv:2001.10995*. 15
- [Xie et al., 2020] Xie, Q., Luong, M.-T., Hovy, E., and Le, Q. V. (2020). Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10687–10698. 2
- [Zeiler and Fergus, 2014] Zeiler, M. D. and Fergus, R. (2014). Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer. 5
- [Zilke et al., 2016] Zilke, J. R., Loza Mencía, E., and Janssen, F. (2016). Deepred – rule extraction from deep neural networks. In Calders, T., Ceci, M., and Malerba, D., editors, *Discovery Science*, pages 457–473, Cham. Springer International Publishing. 6

## A Python APIs for LENS

In order to make LENS paradigms accessible to the whole community, we released "PyTorch, Explain!" [Barbiero et al., 2021], a Python package<sup>8</sup> with an extensive documentation on methods and low-level APIs. Low levels APIs allow the design of custom LENS as illustrated in the example of Listing 2.

<sup>8</sup><https://pypi.org/project/torch-explain/>.

```

1 import torch
2 from torch.nn.functional import one_hot
3 import torch_explain as te
4 from torch_explain.nn.functional import l1_loss
5 from torch_explain.logic.nn import psi
6 from torch_explain.logic.metrics import test_explanation, complexity
7
8 # train data
9 x_train = torch.tensor([
10     [0, 0],
11     [0, 1],
12     [1, 0],
13     [1, 1],
14 ], dtype=torch.float)
15 y_train = torch.tensor([0, 1, 1, 0], dtype=torch.float).unsqueeze(1)
16
17 # instantiate a "psi network"
18 layers = [
19     torch.nn.Linear(x_train.shape[1], 10),
20     torch.nn.Sigmoid(),
21     torch.nn.Linear(10, 5),
22     torch.nn.Sigmoid(),
23     torch.nn.Linear(5, 1),
24     torch.nn.Sigmoid(),
25 ]
26 model = torch.nn.Sequential(*layers)
27
28 # fit (and prune) the model
29 optimizer = torch.optim.AdamW(model.parameters(), lr=0.01)
30 loss_form = torch.nn.BCELoss()
31 model.train()
32 for epoch in range(6001):
33     optimizer.zero_grad()
34     y_pred = model(x_train)
35     loss = loss_form(y_pred, y_train) + 0.000001 * l1_loss(model)
36     loss.backward()
37     optimizer.step()
38
39     model = prune_equal_fanin(model, epoch, prune_epoch=1000, k=2)
40
41 # get first-order logic explanations for a specific target class
42 y1h = one_hot(y_train.squeeze().long())
43 explanation = psi.explain_class(model, x_train)
44
45 # compute explanation accuracy and complexity
46 accuracy, preds = test_explanation(explanation, x_train, y1h, target_class=1)
47 explanation_complexity = complexity(explanation)

```

Listing 2: Example on how to use the "PyTorch Explain!" library to solve the XOR problem.