Improving Out-of-Distribution Detection with Markov Logic Networks Konstantin Kirchheim, Frank Ortmeier

firstname.lastname@ovgu.de

University of Magdeburg, Department of Computer Science

Out-of-distribution (OOD) detection is essential for ensuring the reliability of deep learning models operating in open-world scenarios. Current OOD detectors mainly rely on statistical models to identify unusual patterns in the latent representations of a deep neural network. This work proposes to augment existing OOD detectors with probabilistic reasoning, utilizing Markov logic networks (MLNs). MLNs connect first-order logic with probabilistic reasoning to assign probabilities to inputs based on weighted logical constraints defined over human-understandable concepts, which offers improved explainability. Through extensive experiments on multiple datasets, we demonstrate that MLNs can significantly enhance the performance of a wide range of existing OOD detectors while maintaining computational efficiency. Furthermore, we introduce a simple algorithm for learning logical constraints for OOD detection from a dataset and showcase its effectiveness.

### **Background: OOD Detection with Logic**

#### Idea

- ► Train DNNs to predict some concepts, formulate constraints on plausible concepts
- Inputs that violate a constraint are marked as OOD

### Background: Markov Logic Network (MLN)

- Probabilistic Generalization of First-order Logic (FOL)
- Can be seen as Neuro-Symbolic or templates for large Markov Networks



#### Problem

- Strict logic too rigid for real-world applications where statistical associations dominate
- ► Rather, we want constraints that are violated often only have slight impact, etc.
- $\blacktriangleright$  Each FOL formula  $\varphi_i$  is associated with a weight  $w_i$
- For some input z, a MLN can predict:

$$P_{\mathcal{M}}(z) = \frac{1}{Z} \exp\left(\sum_{i} w_{i}\varphi_{i}(z)\right)$$

#### Standalone Markov Logic Network

- ► Train DNNs to approximate interpretation of FOL predicates  $\{\mathcal{P}_n\}_{n=1}^N$
- $\blacktriangleright$  Create constraint set  $\{\varphi_i\}_{i=1}^N$  with predicates
- $\blacktriangleright$  Train MLN weights  $w_i$  by maximizing likelihood on ID training set
- ► Inference time outlier score:

$$D_{\mathcal{M}}(\mathbf{x}) = -\sum_{i} w_{i} \varphi_{i}(\mathbf{x})$$

**Explainability**: we know exactly by what amount a violated rule changed the outlier score

#### **Combination with other Detectors**

- Normalizing outlier scores is necessary
- For detector  $D: \mathcal{X} \to \mathbb{R}$ , fit some distribution to outlier scores for ID data
- ▶ Then, use survival function  $p_D$ , which normalizes scores into the [0, 1] range
- Outlier score =  $p_D(\mathbf{x}) \times -\sum_i w_i \varphi_i(\mathbf{x})$



#### **Constraint Search**

#### Traffic Sign Classification (GTSRB)

- ► For some datasets, no prior knowledge available
- Idea: take dataset with ID and OOD examples and optimize set of rules

$$\max_{\varphi \in \mathscr{P}(\mathcal{T})} \underbrace{\mathbb{E}_{(\mathbf{x}_{\mathsf{ID}}, \mathbf{x}_{\mathsf{OOD}})}[J(\varphi, \mathbf{x}_{\mathsf{ID}}, \mathbf{x}_{\mathsf{OOD}})]}_{\mathsf{Performance}} - \lambda \underbrace{C(\varphi)}_{\mathsf{Complexity}}$$

- 1: Input: Training set  $\mathcal{D}_{train}$ , validation set  $\mathcal{D}_{val}$ , baseline  $J_0$ , rule set  $\mathcal{T}$
- 2: **Output:** Selected constraints  $\varphi$
- 3: Initialize  $\varphi \leftarrow \emptyset$
- 4: Initialize  $J \leftarrow J_0$
- 5: for all  $\varphi_i \in \mathcal{T}$  do
- $\varphi' \leftarrow \varphi \cup \{\varphi_i\}$
- Train detector with  $\varphi'$  on  $\mathcal{D}_{\text{train}}$ 7:
- $J' \leftarrow \mathsf{Evaluate} \text{ on } \mathcal{D}_{\mathsf{val}}$ 8:
- if  $J' > J + \delta_{\min}$  then 9:
- $J \leftarrow J'$ 10:
- $\varphi \leftarrow \varphi'$ 11:
- end if 12:
- 13: **end for**
- 14: return  $\varphi$

- ► We have 43 constraints over the predicates: class, shape and color
- Statistically significant performance gains, reduces FPR95 by 37% (relative)
- Combination with MLN improves performance of all tested detectors

## **Face Attribute Prediction (CelebA)**

Constraint search on CelebA yields the following result:

$orall \mathbf{x}$	young(x)	(1
$orall \mathbf{x}$	$\textbf{heavy\_makeup}(\mathbf{x}) \rightarrow \textbf{gray\_hair}(\mathbf{x})$	(2
$orall \mathbf{x}$	$wearing\_lipstick(\mathbf{x}) \rightarrow gray\_hair(\mathbf{x})$	(3
$orall \mathbf{x}$	$wearing\_lipstick(\mathbf{x}) \rightarrow no\_beard(\mathbf{x})$	(4
$orall \mathbf{x}$	$\neg male(\mathbf{x}) \rightarrow no\_beard(\mathbf{x})$	(5

- Since constraints are human-understandable, we can manually curate them
- Combination with MLN improves performance of all tested detectors
- ► e.g.: for MLN+Ensemble, FPR95 is reduced by 20% (relative)

## **Constraint Search Regularization**

Regularizing constraint optimization improves results

# **Omiting Rules (GTSRB)**

### **Constraint Search Dataset**

Rules depend on dataset used for optimization



Some rules are more important then others





Sufficient variability seems beneficial

