



**Maynooth
University**

National University
of Ireland Maynooth



Property-driven Machine Learning with Differentiable Logics

Thomas Flinkow

*Department of Computer Science
Maynooth University*

DAIR Course @ Heriot-Watt University
25th September 2024

About

Principles of Programming Research Group @ Maynooth University
<https://www.cs.nuim.ie/research/pop/>



Principles of Programming Research Group

Principal Researchers: Rosemary Monahan, Hao Wu,
Barak A. Pearlmutter, Kevin Casey

2 PostDocs: Medet Inkarbekov, Ali Bukhari

5 PhD Researchers: Oisín Sheridan, Dara MacConville,
Thomas Flinkow, Ankit Jha, Huan Zhang

Projects:

- MAIVV: Modular AI Verification and Visualisation (2021 – 2025, SFI)
- VerifAI: Writing Formal Requirements using AI (2024 – 2026, SFI ADAPT)
- VALU3S: Verification and Validation of Automated Systems' Safety and Security (2020 – 2023, EI and H2020 ESCEL)
- A Model Checker for Python (2022 – 2026, SFI CRT Foundations of Data Science)
- A Constructive Framework for Software Specification and Refinement (2012 – 2023, 2 × IRC)

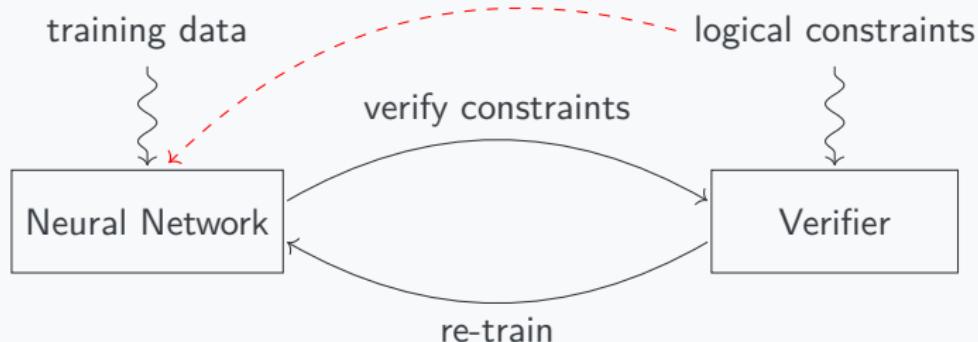
Background

Training with Logical Constraints

Task: train a neural network \mathcal{N} to satisfy constraint ϕ .

Train: given data, labels, and loss function, iteratively update network weights.

Verify: afterwards, α, β -CROWN, Marabou, NNV, ERAN, ...



Note

Training with constraints does not guarantee their satisfaction!

Property-driven Training with Differentiable Logics

Given data x_0 and label y , and constraint ϕ ,
obtain optimal network weights θ^+ by

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(x_0, y) + \beta \mathcal{L}_C(x_0, y, \phi).$$

Optimisation handled by Gradient Descent:

$$\theta_{i+1} := \theta_i - \eta \nabla \mathcal{L}(\theta_i)$$

Implications

Logical constraint ϕ must be differentiable!

- mapping $\llbracket \cdot \rrbracket_{\text{DL2}} : \Phi \rightarrow [0, \infty)$,
- $\llbracket \phi \rrbracket_{\text{DL2}} = 0$ iff ϕ is satisfied,
- $\llbracket \phi \rrbracket_{\text{DL2}}$ is differentiable almost everywhere.

Recursive definition of loss translation:

$$\llbracket x \leq y \rrbracket_{\text{DL2}} := \max\{x - y, 0\}$$

$$\llbracket \phi \wedge \psi \rrbracket_{\text{DL2}} := \llbracket \phi \rrbracket_{\text{DL2}} + \llbracket \psi \rrbracket_{\text{DL2}}$$

$$\llbracket \phi \vee \psi \rrbracket_{\text{DL2}} := \llbracket \phi \rrbracket_{\text{DL2}} \cdot \llbracket \psi \rrbracket_{\text{DL2}}.$$

¹Marc Fischer et al.: 'DL2: Training and Querying Neural Networks with Logic' 2019.

Fuzzy Logics^{2,3}

- mapping $\llbracket \cdot \rrbracket_L : \Phi \rightarrow [0, 1]$, where $\llbracket \top \rrbracket_L = 1$ and $\llbracket \perp \rrbracket_L = 0$,
- operators happen to be differentiable almost everywhere

Logic	Conjunction	Disjunction	Implication
Gödel	$\min\{x, y\}$	$\max\{x, y\}$	$\begin{cases} 1, & \text{if } x < y, \\ y, & \text{else.} \end{cases}$
Kleene-Dienes			
Łukasiewicz	$\max\{0, x + y - 1\}$	$\min\{1, x + y\}$	$S(N(x), y)$
Reichenbach			
Goguen	xy	$x + y - xy$	$\begin{cases} 1, & \text{if } x < y, \\ y^x, & \text{else.} \end{cases}$

²Natalia Ślusarz et al.: 'Logic of Differentiable Logics: Towards a Uniform Semantics of DL' 2023. DOI: [10.29007/c1nt](https://doi.org/10.29007/c1nt).

³Emile van Krieken et al.: 'Analyzing Differentiable Fuzzy Logic Operators' 2022. DOI: [10.1016/j.artint.2021.103602](https://doi.org/10.1016/j.artint.2021.103602).

Constraint Counterexamples

Given data \mathbf{x}_0 and label \mathbf{y} , and constraint ϕ ,
obtain optimal network weights θ^+ by

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(\mathbf{x}_0, \mathbf{y}) + \beta \mathcal{L}_C(\mathbf{x}_0, \mathbf{y}, \phi).$$

Insight from DL2 (Fischer et al. 2019)

Learning to satisfy $\forall x. x \models \phi$ by finding x^* such that $x^* \not\models \phi$.

1. Approximate counterexample *outside* of training set using PGD:

$$\mathbf{x}^* = \arg \max_{\mathbf{x} \in \|\mathbf{x} - \mathbf{x}_0\|_\infty \leq \epsilon} \mathcal{L}_C(\mathbf{x}_0, \mathbf{x}, \mathbf{y}, \phi)$$

2. Use this counterexample in training:

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(\mathbf{x}_0, \mathbf{y}) + \beta \mathcal{L}_C(\mathbf{x}_0, \mathbf{x}^*, \mathbf{y}, \phi).$$

Code

Full Code On DAIR Course GitHub⁴

The screenshot shows a GitHub repository page for 'DAIR-course / lab-exercises / lab4 / property-driven-training'. The left sidebar shows a tree view of files and folders. The main area displays the repository structure and requirements.

Repository Structure:

- README.md
- src
 - d12.py
 - fuzzy_logics.py
 - constraints.py
 - main.py
 - models.py
 - run.sh
 - util.py
 - reports
 - requirements.txt

Requirements:

The experiments were run on Python 3.12.3 (but should probably work with newer versions as well). The provided requirements.txt can be used to install the required packages using pip install -r requirements.txt.



⁴<https://github.com/KatyaKom/DAIR-course/tree/main/lab-exercises/lab4/property-driven-training/>

Integration into PyTorch

Standard PyTorch training + logical property

```
def train(..):
    for _, (inputs, labels) in enumerate(train_loader):
        outputs = NN(inputs)
        ce_loss = F.cross_entropy(outputs, labels)

        adv = pgd.attack(NN, inputs, labels, constraint)
        dl_loss = constraint.eval(NN, inputs, adv, labels)

        loss = alpha * ce_loss + beta * dl_loss

        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
```

Balancing Loss

Problem

It is *crucial* to find close to optimal values for α and β !

$$\theta^+ = \arg \min_{\theta} \alpha \mathcal{L}_{CE}(x_0, y) + \beta \mathcal{L}_C(x_0, x^*, y, \phi).$$

Adaptive Loss Balancing with GradNorm⁵

- Key point: $\alpha(t)$ and $\beta(t)$
- Better results than expensive grid search!

Loss balancing in code⁶

```
loss = grad_norm.weighted_loss(batch_index, ce_loss, dl_loss, alpha)
# loss = alpha * ce_loss + beta * dl_loss
```

⁵Zhao Chen et al.: 'GradNorm: Gradient Normalization for Adaptive Loss Balancing in Deep Multitask Networks' 2018.

⁶<https://github.com/KatyaKom/DAIR-course/blob/main/lab-exercises/lab4/property-driven-training/main.py#L56>

Basics – Logic⁷

```
class Logic(ABC):
    @abstractmethod
    def LEQ(self, x, y):
        pass

    @abstractmethod
    def NOT(self, x, y):
        pass

    @abstractmethod
    def AND(self, x, y):
        pass

    @abstractmethod
    def OR(self, x, y):
        pass

    def IMPL(self, x, y):
        return self.OR(self.NOT(x), y)

    def EQUIV(self, x, y):
        return self.AND(self.IMPL(x, y), self.IMPL(y, x))
```

Basics – Constraints⁸

```
class Constraint(ABC):
    def __init__(self, eps):
        self.eps = eps

    @abstractmethod
    def get_constraint(self, NN, inputs, adv, labels):
        pass

    # usage:    loss, sat = eval()
    def eval(self, NN, inputs, adv, labels, logic, train):
        ...
```

Remark

Why ϵ in the constructor? All constraints use PGD!

⁸<https://github.com/KatyaKom/DAIR-course/tree/main/lab-exercises/lab4/property-driven-training/constraints.py#L14>

Local Robustness Constraint



Figure 1: Adversarial attack (Goodfellow et al. 2015).

Definition

A neural network is **locally robust** in input x_0 , if

$$\underbrace{\forall x. \|x - x_0\|_\infty \leq \varepsilon}_{\text{all elements in the input space close to } x_0}$$

implies

$$\underbrace{\|\mathcal{N}(x) - \mathcal{N}(x_0)\|_\infty \leq \delta}_{\text{the classification is roughly the same}}$$

Local Robustness Constraint – Code⁹

Definition

A neural network is locally robust in input x_0 , if

$$\underbrace{\forall x. \|x - x_0\|_\infty \leq \varepsilon}_{\text{handled by PGD}} \quad \text{implies} \quad \underbrace{\|\mathcal{N}(x) - \mathcal{N}(x_0)\|_\infty \leq \delta}_{\text{the classification is roughly the same}}$$

```
class RobustnessConstraint(Constraint):
    def __init__(self, eps, delta):
        super().__init__(eps)
        self.delta = delta

    def get_constraint(self, NN, inputs, adv, _labels):
        diff = F.softmax(NN(adv)) - F.softmax(NN(inputs))
        norm = torch.linalg.vector_norm(diff, ord=float('inf'))

        return lambda l: l.LEQ(norm, self.delta)
```

⁹<https://github.com/KatyaKom/DAIR-course/tree/main/lab-exercises/lab4/property-driven-training/constraints.py#L45>

Group Constraint



(a) unique signs



(b) danger signs



(c) derestriction signs



(d) speed limit signs



(e) other prohibitory signs



(f) mandatory signs

Definition

$$\underbrace{\forall \mathbf{x} \in \|\mathbf{x} - \mathbf{x}_0\| \leq \epsilon}_{\text{handled by PGD}} \rightarrow \bigwedge_{G \in \mathcal{G}} p_G \leq \delta \vee p_G \geq 1 - \delta.$$

Group Constraint – Code¹⁰

Definition

$$\underbrace{\forall \mathbf{x} \in \|\mathbf{x} - \mathbf{x}_0\| \leq \epsilon}_{\text{handled by PGD}} \rightarrow \bigwedge_{G \in \mathcal{G}} p_G \leq \delta \vee p_G \geq 1 - \delta.$$

```
class GroupConstraint(Constraint):
    def __init__(self, eps, delta, groups):
        super().__init__(eps)
        self.delta = delta

    def get_constraint(self, NN, _inputs, adv, _labels):
        probs = F.softmax(NN(adv))
        sums = [torch.sum(probs[:, indices]) for indices in group_indices]

        return lambda l: reduce(l.AND, \
            [l.OR(l.LEQ(s, delta), l.LEQ(1. - delta, s)) for s in sums])
```

¹⁰<https://github.com/KatyaKom/DAIR-course/tree/main/lab-exercises/lab4/property-driven-training/constraints.py#L55>

Specifying Hypershapes¹¹

Defining ℓ_∞ hyperrectangles

```
hyper_rectangles = []

for data, _ in train_loader:
    data = data.to(device).view(data.size(0), -1)
    hyper_rectangles.append(torch.stack((data - eps, data + eps)))
```

¹¹<https://github.com/KatyaKom/DAIR-course/tree/main/lab-exercises/lab4/property-driven-training/main.py#L208>

**Thank you!
Any questions?**

References

- Chen, Zhao et al. (July 2018). 'GradNorm: Gradient Normalization for Adaptive Loss Balancing in Deep Multitask Networks'. In: *Proceedings of the 35th International Conference on Machine Learning*. PMLR, pp. 794–803.
- Fischer, Marc et al. (24th May 2019). 'DL2: Training and Querying Neural Networks with Logic'. In: *Proceedings of the 36th International Conference on Machine Learning*. International Conference on Machine Learning. PMLR, pp. 1931–1941.
- Goodfellow, Ian J. et al. (Mar. 2015). *Explaining and Harnessing Adversarial Examples*. DOI: [10.48550/arXiv.1412.6572](https://doi.org/10.48550/arXiv.1412.6572). arXiv: [1412.6572 \[cs, stat\]](https://arxiv.org/abs/1412.6572).

References (cont.)

- Ślusarz, Natalia et al. (3rd June 2023). 'Logic of Differentiable Logics: Towards a Uniform Semantics of DL'. In: *EPiC Series in Computing*. Proceedings of 24th International Conference on Logic for Programming, Artificial Intelligence and Reasoning. Vol. 94. EasyChair, pp. 473–493. DOI: [10.29007/cint](https://doi.org/10.29007/cint).
- Van Krieken, Emile et al. (Jan. 2022). 'Analyzing Differentiable Fuzzy Logic Operators'. In: *Artificial Intelligence* 302, p. 103602. ISSN: 00043702. DOI: [10.1016/j.artint.2021.103602](https://doi.org/10.1016/j.artint.2021.103602). arXiv: [2002.06100 \[cs\]](https://arxiv.org/abs/2002.06100).