

Extending Layerwise Relevance Propagation using Semiring Annotations

Antoine Groudiev
L3, ENS Ulm

Silviu Maniu – Supervisor
SLIDE Team, LIG

Tuesday, July 9th

Plan

Introduction

- Problem statement

- Layerwise Relevance Propagation

- Semiring-based provenance annotations

Extending LRP

Applications

- Image mask computation

- Network pruning using LRP ranking

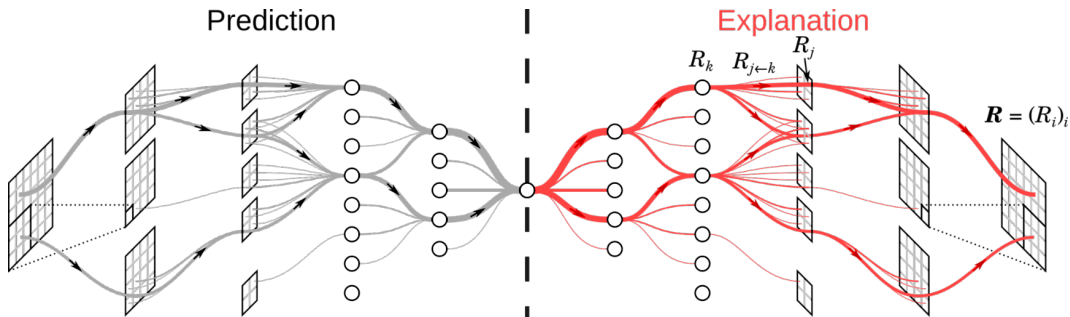
- Comparison to image perturbation

Conclusion



Problem statement

Layerwise Relevance Propagation



Layerwise Relevance Propagation

Propagation rules

Initialization:

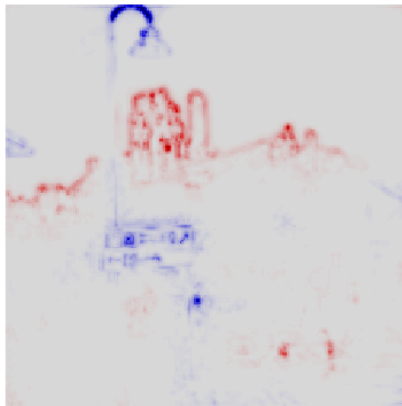
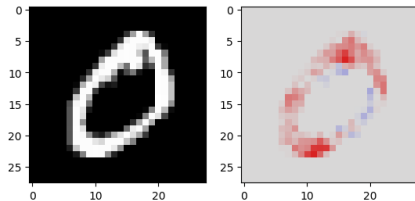
$$R_i^{(L)} = \begin{cases} a_i^{(L)} & \text{if } i = y \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

LRP-0 rule:

$$R_j^{(l)} = \sum_k \frac{a_j^{(l)} w_{j,k}}{\sum_{j'} a_{j'}^{(l)} w_{j',k}} R_k^{(l+1)} \quad (2)$$

Layerwise Relevance Propagation

Results visualization



Semiring-based provenance annotations

Definition (Semiring)

A semiring $(\mathbb{K}, \oplus, \otimes, \mathbf{0}, \mathbf{1})$ is composed of a set \mathbb{K} , binary operators \oplus and \otimes such that \otimes distributes over \oplus , verifying the following properties:

- $(\mathbb{K}, \oplus, \mathbf{0})$ is a commutative monoid
- $(\mathbb{K}, \otimes, \mathbf{1})$ is a monoid such that $\mathbf{0}$ is absorbing

Example

The following structures are semirings:

- Real semiring: $(\mathbb{R}, +, \times, 0, 1)$
- Boolean semiring: $(\{\perp, \top\}, \vee, \wedge, \perp, \top)$
- Counting semiring: $(\mathbb{N}, +, \times, 0, 1)$
- Viterbi semiring: $([0, 1], \max, \times, 0, 1)$

Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Conclusion

Semiring generalization of the LRP rule

Conversion functions for activations, weights:

$$\Theta_a : \mathbb{R} \longrightarrow \mathbb{K}$$

$$\Theta_w : \mathbb{R} \longrightarrow \mathbb{K}$$

Initialization:

$$R_i^{(L)} = \begin{cases} \Theta_a \left(a_i^{(L)} \right) & \text{if } i = y \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (3)$$

Propagation rule:

$$R_j^{(l)} = \bigoplus_k \Theta_a \left(a_j^{(l)} \right) \otimes \Theta_w \left(w_{j,k}^{(l)} \right) \otimes R_k^{(l+1)} \quad (4)$$

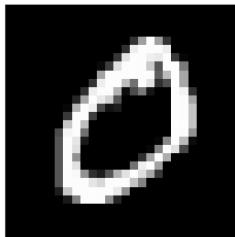
Boolean Semiring

$(\{\perp, \top\}, \vee, \wedge, \perp, \top)$

$$\Theta_a = a \mapsto \begin{cases} \top & \text{if } a \geq \theta_a \\ \perp & \text{otherwise} \end{cases}$$

$$\Theta_w = w \mapsto \begin{cases} \top & \text{if } w \geq \theta_w \\ \perp & \text{otherwise} \end{cases}$$

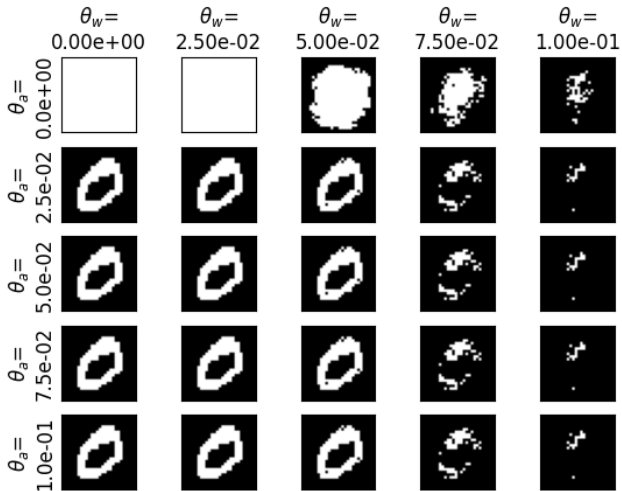
Reference



LRP-0 without z^B
Boolean Semiring



Influence of the thresholds

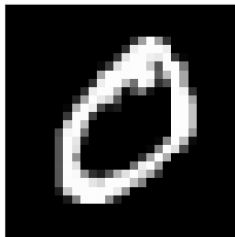


$(\mathbb{N}, +, \times, 0, 1)$

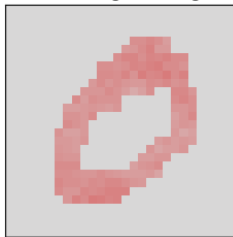
$$\Theta_a = a \mapsto \begin{cases} 1 & \text{if } a \geq \theta_a \\ 0 & \text{otherwise} \end{cases}$$

$$\Theta_w = w \mapsto \begin{cases} 1 & \text{if } w \geq \theta_w \\ 0 & \text{otherwise} \end{cases}$$

Reference



LRP-0 without z^B
Counting semiring



Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Conclusion

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Accuracies per attack zone
Kernel size: 4 — Step: 1

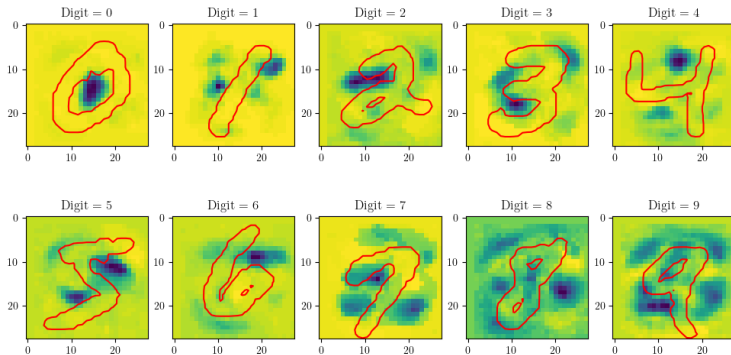


Figure: Accuracies per attack zone

Plan

Introduction

Problem statement

Layerwise Relevance Propagation

Semiring-based provenance annotations

Extending LRP

Applications

Image mask computation

Network pruning using LRP ranking

Comparison to image perturbation

Conclusion

- [1] Sebastian Bach et al. “On Pixel-Wise Explanations for Non-Linear Classifier Decisions by Layer-Wise Relevance Propagation”. In: *PLOS ONE* (2015), pp. 1–46. DOI: 10.1371/journal.pone.0130140. URL: <https://doi.org/10.1371/journal.pone.0130140>.
- [2] Ruth C Fong and Andrea Vedaldi. “Interpretable explanations of black boxes by meaningful perturbation”. In: *Proceedings of the IEEE international conference on computer vision*. 2017, pp. 3429–3437. URL: <https://arxiv.org/abs/1704.03296>.
- [3] Robert Geirhos et al. “Shortcut learning in deep neural networks”. In: *Nature Machine Intelligence* 2 (2020), pp. 665–673.
- [4] Todd J Green, Grigoris Karvounarakis, and Val Tannen. “Provenance semirings”. In: *Proceedings of the twenty-sixth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*. 2007, pp. 31–40.
- [5] Grégoire Montavon et al. “Layer-Wise Relevance Propagation: An Overview”. In: *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*. Springer International Publishing, 2019, pp. 193–209. URL: https://doi.org/10.1007/978-3-030-28954-6_10.
- [6] Yann Ramusat, Silviu Maniu, and Pierre Senellart. “Provenance-Based Algorithms for Rich Queries over Graph Databases”. In: *EDBT 2021 - 24th International Conference on Extending Database Technology*. 2021. URL: <https://inria.hal.science/hal-03140067>.