

BlackboxNLP 2019

August 1, Florence, Italy, @ ACL 2019

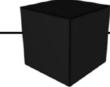
Evaluating Recurrent Neural Network Explanations

Leila Arras , Ahmed Osman, Klaus-Robert Müller and Wojciech Samek

Explaining predictions

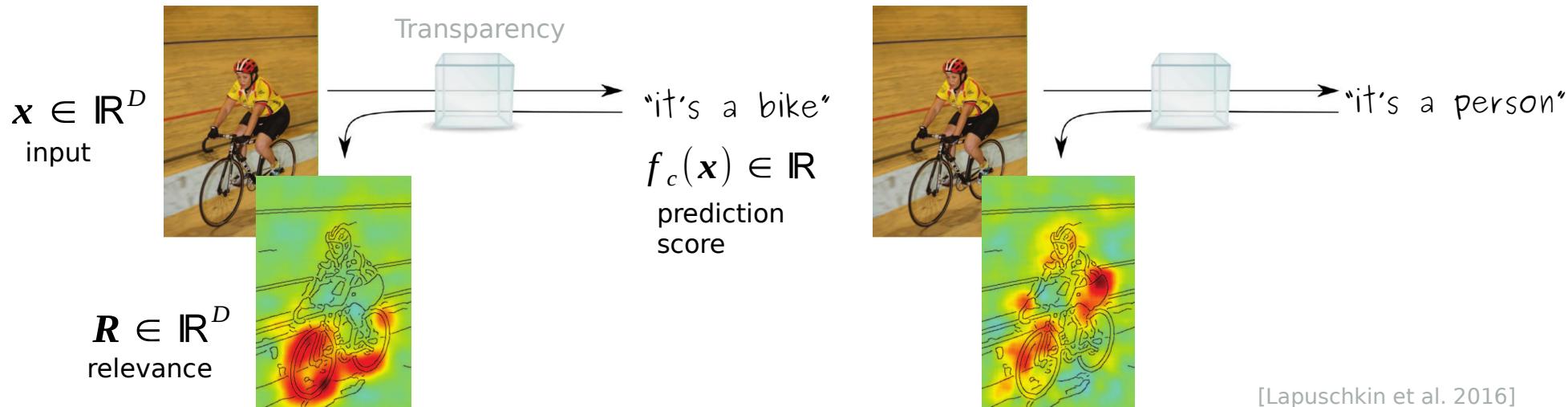


Blackbox



"it's a bike" **but why?**

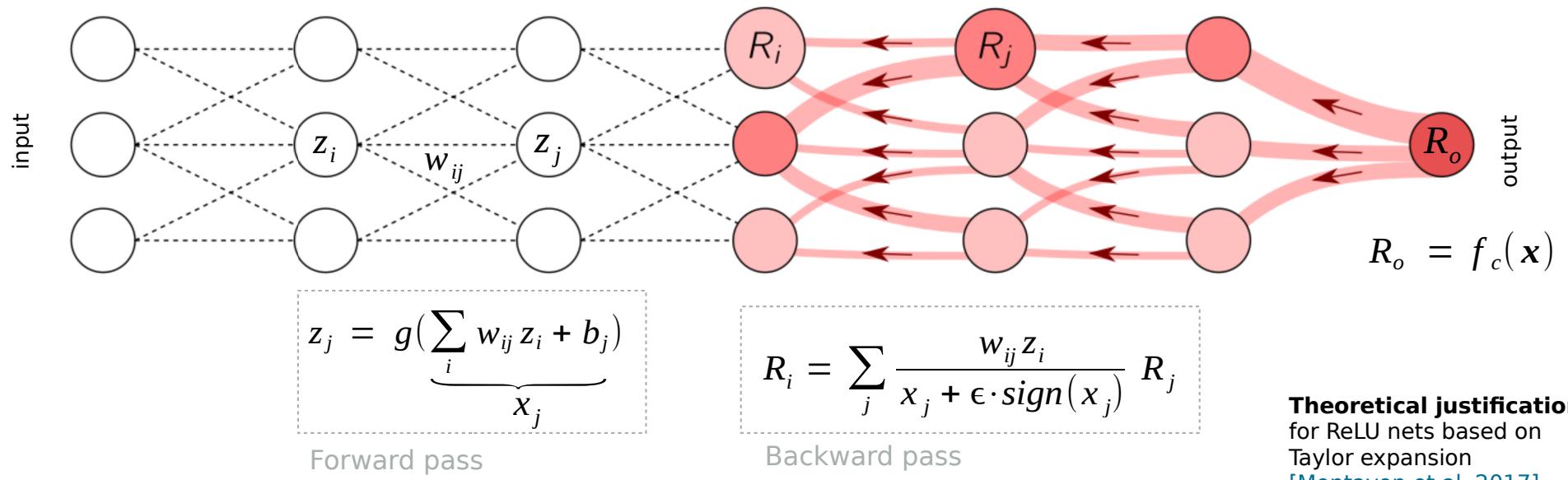
- Decision for right reason or biased?
- EU-GDPR 2016 compliant?



[Lapuschkin et al. 2016]

Layer-wise Relevance Propagation (LRP)

Idea: Decompose prediction function value $f_c(\mathbf{x})$. [Bach et al. 2015]

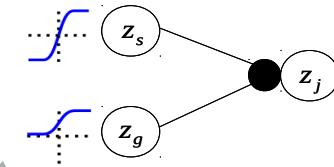


Extending LRP to recurrent networks

New: How to propagate relevance through products?

Forward pass

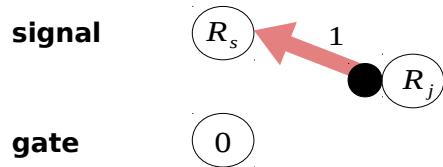
$$z_j = z_s \cdot z_g$$



signal (signed)
“is the information”

gate $\in (0,1)$
“controls the flow of information”

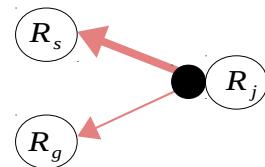
LRP-all



signal – take – all

Arras et al. 2017b

LRP-prop

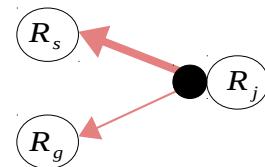


$$R_s = \frac{z_s}{z_s + z_g} R_j$$

$$R_g = \frac{z_g}{z_s + z_g} R_j$$

Ding et al. 2017
Arjona-Medina et al. 2018

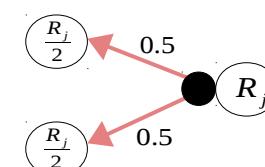
LRP-abs



$$R_s = \frac{|z_s|}{|z_s| + |z_g|} R_j$$

$$R_g = \frac{|z_g|}{|z_s| + |z_g|} R_j$$

LRP-half



equally

Arjona-Medina et al. 2018

Other methods for recurrent networks

- Contextual Decomposition (CD) [Murdoch et al. 2018]
- Gradient, Gradient x Input [Li et al. 2016, Denil et al. 2015, Gevrey et al. 2003]
- Occlusion relevance [Li et al. 2017]

see paper for details

now which one to use?

Evaluating Explanations

Look at a few heatmaps?

Predicted class: negative sentiment

this movie was actually neither that funny , nor super witty .

i hate the movie though the plot is interesting .

not worth the time

is n't a bad film (misclassified)

it never fails to engage us . (misclassified)

OK but this is not enough

Perturbation experiment

Idea:

Perturb the input according to **word relevance ordering**
Track the impact on the prediction [Arras et al. 2016]

akin **pixel-flipping/region perturbation** in Computer Vision
[Bach et al. 2015,
Samek et al. 2017]
or **ablation**

Step	Input	Predicted sentiment
0	i hate the movie though the plot is interesting .	--
1	i [] the movie though the plot is interesting .	+
2	i [] the movie though the [] is interesting .	+
3	i [] the movie [] the [] is interesting .	++

word embedding set to zero

Perturbation Results

LSTM sentiment analysis

Word deletion	Grad x Input	LRP-all	LRP-prop	LRP-abs	LRP-half	CD	occlusion
Impact on accuracy	✗	✓	✗	✗	✗	✓	✓

Model: Bidirectional LSTM from [Li et al. 2016]

Task: 5-class sentiment prediction (SST 5) [Socher et al. 2013]

Perturbation Results

LSTM sentiment analysis

Word deletion	Grad x Input	only one backward pass				multiple passes	
		LRP-all	LRP-prop	LRP-abs	LRP-half	CD	occlusion
Impact on accuracy	✗	✓	✗	✗	✗	✓	✓

Model: Bidirectional LSTM from [Li et al. 2016]

Task: 5-class sentiment prediction (SST 5) [Socher et al. 2013]

Toy Arithmetic Task

We propose a toy setup, with **Ground Truth** relevance

Input Sequence

$$\begin{bmatrix} 0 & 0 & 0 & n_a & 0 & 0 & 0 & n_b & 0 & 0 & 0 \\ \underbrace{n_1}_{x_1} & \dots & \underbrace{n_{a-1}}_{x_2} & \boxed{n_a \\ 0} & n_{a+1} & \dots & n_{b-1} & \boxed{n_b \\ 0} & n_{b+1} & \dots & \underbrace{n_T}_{x_T} \end{bmatrix}$$

Task 1

$$y_{target} = n_a + n_b$$

addition

$$n_t \in \mathbb{R}$$

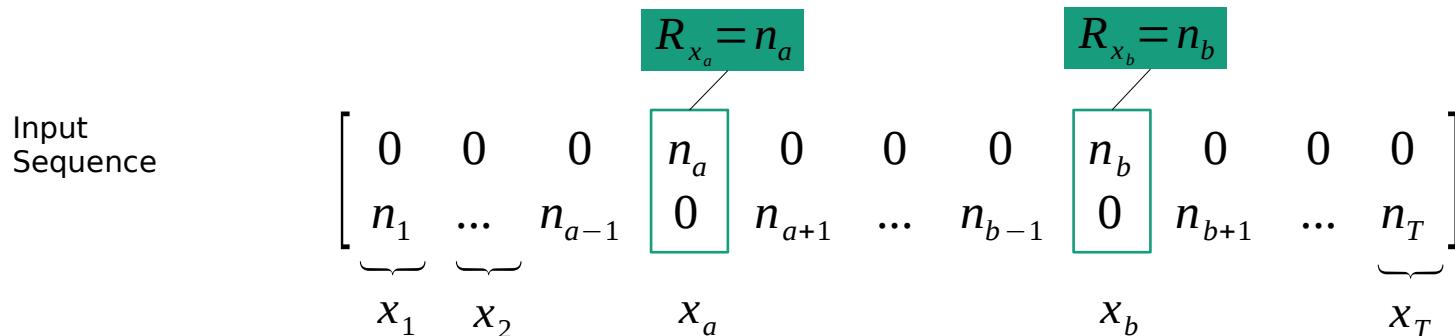
Model

LSTM with one cell

Similar to
adding problem
[Hochreiter and
Schmidhuber
1996]

Toy Arithmetic Task

We propose a toy setup, with **Ground Truth relevance**



Task 1

$$y_{target} = n_a + n_b$$

addition

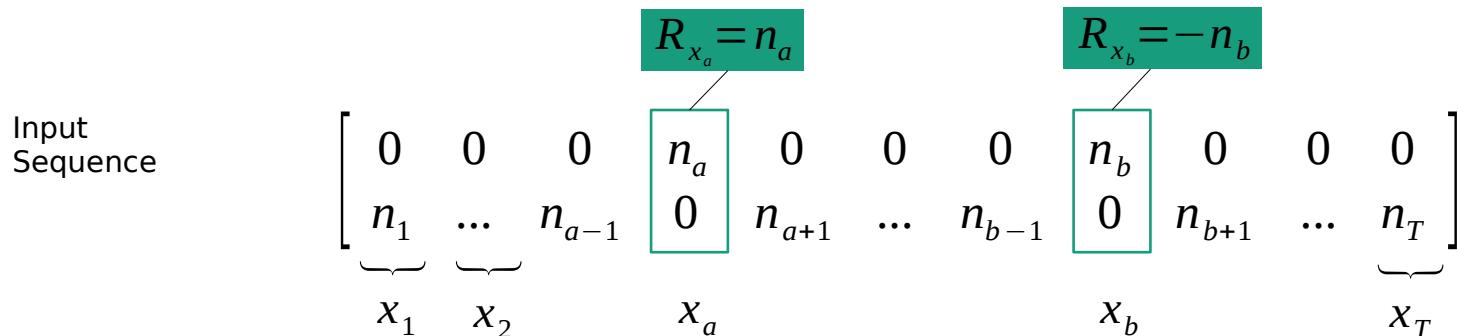
$$n_t \in \mathbb{R}$$

Model

LSTM with one cell

Toy Arithmetic Task

We propose a toy setup, with **Ground Truth relevance**



Task 2

$$y_{target} = n_a - n_b \quad \text{subtraction} \quad n_t \in \mathbb{R}^+$$

Model

LSTM with one cell

Results

Toy Task Addition	Correlation between relevance and input number	
	$\rho(R_{\alpha_a}, n_a)$ (in %)	$\rho(R_{\alpha_b}, n_b)$ (in %)
✓ Gradient × Input	99.960 (0.017)	99.954 (0.019)
✓ Occlusion	99.990 (0.004)	99.990 (0.004)
✗ LRP-prop	0.785 (3.619)	10.111 (12.362)
✗ LRP-abs	7.002 (6.224)	12.410 (17.440)
✗ LRP-half	29.035 (9.478)	51.460 (19.939)
✓ LRP-all	99.995 (0.002)	99.995 (0.002)
✓ CD	99.997 (0.002)	99.997 (0.002)

Results

	$\rho(R_{\mathbf{x}_a}, n_a)$ (in %)	$\rho(R_{\mathbf{x}_b}, n_b)$ (in %)		$\rho(R_{\mathbf{x}_a}, n_a)$ (in %)	$\rho(R_{\mathbf{x}_b}, n_b)$ (in %)
Toy Task Addition					
✓ Gradient×Input	99.960 (0.017)	99.954 (0.019)		✓ Gradient×Input	97.9 (1.6)
✓ Occlusion	99.990 (0.004)	99.990 (0.004)		✗ Occlusion	99.0 (2.0)
✗ LRP-prop	0.785 (3.619)	10.111 (12.362)		✗ LRP-prop	3.1 (4.8)
✗ LRP-abs	7.002 (6.224)	12.410 (17.440)		✗ LRP-abs	1.2 (7.6)
✗ LRP-half	29.035 (9.478)	51.460 (19.939)		✗ LRP-half	7.7 (15.3)
✓ LRP-all	99.995 (0.002)	99.995 (0.002)		✓ LRP-all	98.5 (3.5)
✓ CD	99.997 (0.002)	99.997 (0.002)		✗ CD	-25.9 (39.1)
Toy Task Subtraction					
✓ Gradient×Input				-98.8 (0.6)	
✗ Occlusion				-69.0 (19.1)	
✗ LRP-prop				-8.4 (18.9)	
✗ LRP-abs				-23.0 (11.1)	
✗ LRP-half				-28.9 (6.4)	
✓ LRP-all				-99.3 (1.3)	
✗ CD				-50.0 (29.2)	

Results

Why different results?

	$\rho(R_{\mathbf{x}_a}, n_a)$ (in %)	$\rho(R_{\mathbf{x}_b}, n_b)$ (in %)
Toy Task Addition		
✓ Gradient × Input	99.960 (0.017)	99.954 (0.019)
✓ Occlusion	99.990 (0.004)	99.990 (0.004)
✗ LRP-prop	0.785 (3.619)	10.111 (12.362)
✗ LRP-abs	7.002 (6.224)	12.410 (17.440)
✗ LRP-half	29.035 (9.478)	51.460 (19.939)
✓ LRP-all	99.995 (0.002)	99.995 (0.002)
✓ CD	99.997 (0.002)	99.997 (0.002)

Addition can be solved by a bag of words

	$\rho(R_{\mathbf{x}_a}, n_a)$ (in %)	$\rho(R_{\mathbf{x}_b}, n_b)$ (in %)
Toy Task Subtraction		
✓ Gradient × Input	97.9 (1.6)	-98.8 (0.6)
✗ Occlusion	99.0 (2.0)	-69.0 (19.1)
✗ LRP-prop	3.1 (4.8)	-8.4 (18.9)
✗ LRP-abs	1.2 (7.6)	-23.0 (11.1)
✗ LRP-half	7.7 (15.3)	-28.9 (6.4)
✓ LRP-all	98.5 (3.5)	-99.3 (1.3)
✗ CD	-25.9 (39.1)	-50.0 (29.2)

Subtraction is sequential (order matters!)

Summary

LSTM for sentiment analysis & toy tasks

Evaluation	LRP-all	LRP-prop	LRP-abs	LRP-half	CD	Grad x Input	occlusion
Perturbation	✓	✗	✗	✗	✓	✗	✓
Addition	✓	✗	✗	✗	✓	✓	✓
Subtraction	✓	✗	✗	✗	✗	✓	✗
Subject-verb agreement	✓	—	—	—	—	✗	✗

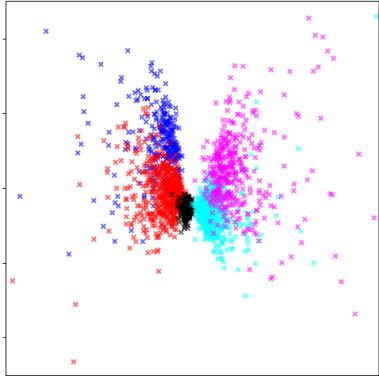
Evaluation by Poerner et al. 2018

Model: LSTM, GRU, Quasi-RNN

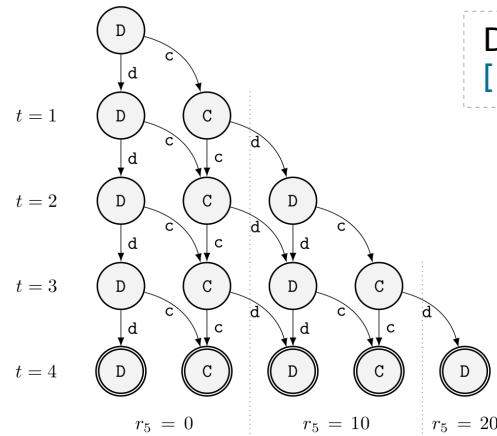
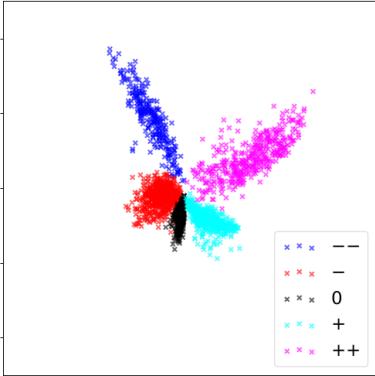
Task: verb number prediction [Linzen et al. 2016]

Applications of relevance

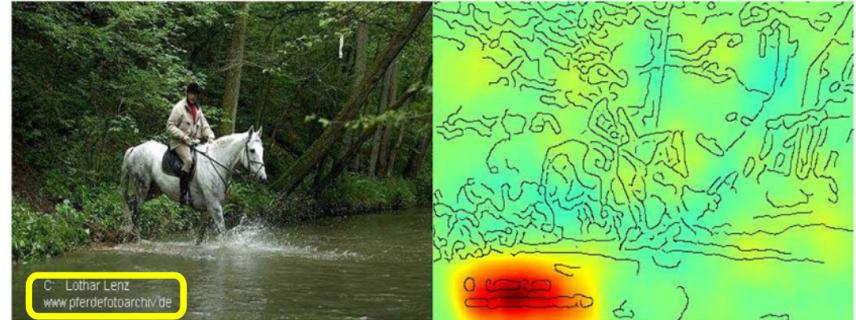
Applications



Sentence-level/Document-level vector representations
[present work, Arras et al. 2017a]



Detect model/dataset bias
[Lapuschkin et al. 2019]



Redistribute rewards in RL
[Arjona-Medina et al. 2018]

Conclusion

- LRP on LSTMs [Bach et al. 2015, Arras et al. 2017b] works well
→ but correct redistribution of relevance through product layer is crucial
- Some methods fail on a simple toy task
(CD [Murdoch et al. 2018], Occlusion [Li et al. 2017])
- Still, we need more evaluation tasks for relevance!
- Applications could also be extended.

Thanks for your attention!

More information on LRP (code, tutorials, ...):
www.heatmapping.org

Reference implementation of LRP for LSTMs:
github.com/ArrasL/LRP_for_LSTM

Tutorial paper:

Methods for interpreting and understanding deep neural networks
Grégoire Montavon, Wojciech Samek and Klaus-Robert Müller
Digital Signal Processing 73 (2018) 1-15

Upcoming book:

Chapter on “Explaining and Interpreting LSTMs”
Leila Arras*, Jose Arjona-Medina*, Michael Widrich, Grégoire Montavon, Michael Gillhofer, Klaus-Robert Müller, Sepp Hochreiter and Wojciech Samek (*equal contribution)

