

# Training with Explanations Alone

## Appendix A

### A.1 LRP fundamentals and tracing back bias

①

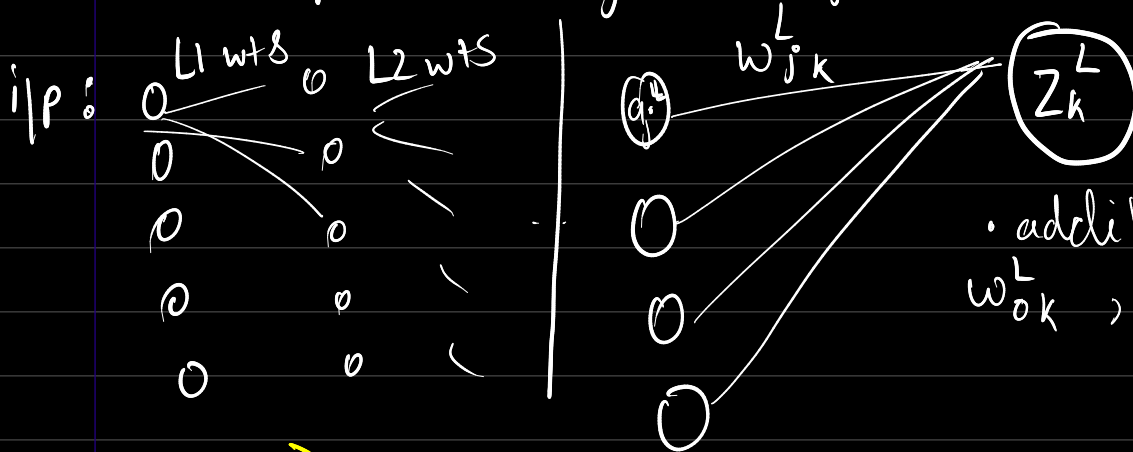
(A2)

$$z_k^L = \sum_j w_{jk}^L a_j^L$$

•  $z_k^L$  is the output of the  $k^{\text{th}}$  neuron in the  $L^{\text{th}}$  fully connected layer

BEFORE ReLU activation.

↳  $z_k^L$  is a function of the outputs of neurons  $j$  in the previous layer:  $a_j^L$



• additive bias is  $w_{0k}^L$ ,  $a_0^L = 1$ .

②

(A1)

$$R_j^{L-1} = \sum_k \frac{w_{jk}^L a_j^L}{z_k^L + \text{sign}(z_k^L) \epsilon} R_k^L$$

• LRP- $\epsilon$  relevance backprop rule used to propagate relevance from neuron  $k$  in layer  $L$  to neurons  $j$  in layer  $L-1$ .

•  $R_k^L$  denotes how relevant the neuron  $k$  in layer  $L$  is to the output of the neural network (logits?)

\* LRP back prop begins at one of the neurons in the final output layer say neuron  $c$  in layer  $N$ .  $R_c^N$  is the logit and  $z_c^N = R_c^N$

- We begin at the neuron associated with the class we want the LRP heatmap to explain ( $C$ ).  
 $\hookrightarrow$  relevances of other neurons  $R_k^N = 0$  and  $R_c^N = Z_c^N$ .
- Backprop ends at the input where LRP explanation heatmap (relevance of each input pixel to neural net's output) is formed.  $\bar{R}^0 = [R_k^0]$ .
- In the explanation heatmap  $\bar{R}^0 = [R_k^0]$ .
  - i) -ve values indicate areas that REDUCE the classifier's confidence for the explained class  $C$
  - ii) +ve values are areas that increased this confidence
  - iii) low absolute relevance indicates areas the classifier paid little attention to.
- eq<sup>n</sup> A1:  $\epsilon$  is a small ( $10^{-4}$ ,  $10^{-2}$ ) hyperparam to denoise the heatmap.
  - $\hookrightarrow$  it is added to  $Z_k^L (\sum_j W_{jk}^L a_j)$ .
  - $\therefore$  LRP prop. rule  $\div$  by  $Z_k^L$ , we add  $\epsilon$  to prevent DA from being too close to 0 and the relevance exploding
  - $\hookrightarrow$  also  $\therefore$  this method is based on Taylor series approximations,  $\epsilon$  represents the residue in the left out terms of the series?

$$R_j^{L-1} = \sum_k \frac{w_{jk}^L a_j^L}{z_k^L + \text{sign}(z_k^L) \epsilon} R_k^L$$

$$\text{where: } z_k^L = \sum_j w_{jk}^L a_j^L$$

$$R_k^N = \begin{cases} z_k^N & \text{if } k = c \\ 0 & \text{otherwise} \end{cases}$$

- the amount of relevance back-propagated from neuron  $k$  to neuron  $j$  depends on the relative influence of neuron  $j$  on the output of neuron  $k$  ( $z_k^L$ )

$$\hookrightarrow \frac{w_{jk}^L a_j^L}{z_k^L + \epsilon \frac{|z_k^L|}{z_k^L}}$$

- For the bias to influence the classifier's output its influence must flow from its position in the input image through all layers.

$\hookrightarrow \therefore \exists$  a corresponding backward flow from the o/p to the bias location in the input image.



$\rightarrow$  going forwards now?

- The relevance of a neuron  $k$  in layer  $L$  ( $R_k^L$ ) is a function of the outputs of neurons in layer  $L-1$  ( $\bar{a}^L = [a_j^L]$ )

Let  $\bar{a}^L = \tilde{a}^L + \tilde{h}$   $\rightarrow$  reference point  $x_0$  for Taylor's

then:

$$R_k^L(\bar{a}^L) = R_k^L(\tilde{a}^L) + \sum_j \underbrace{(\bar{a}_j^L - \tilde{a}_j^L)}_{\tilde{a}_j^L} \frac{\partial R_k^L}{\partial \tilde{a}_j^L} + \dots$$

$R_k^L(\bar{a}^L)$  is broken up

into one term for each neuron  $j$

each of these terms represents how much a change in the value of neuron  $j$  in layer  $L-1$  would affect the relevance  $R_k^L$  of neuron  $k$  in layer  $L$ .

$\hookrightarrow \therefore$  it is these terms LRP backpropagates from neuron  $k$  to each neuron  $j$ .

## Final Understanding :

- So basically, any bias to influence o/p's needs to go through every layer.
- Paper says we can see the relevance  $R_k^L$  as a function of o/p's of neurons in  $L-1$   
 $\therefore R_k^L = R_k^L(\bar{a}^L)$   
 $\hookrightarrow$  can do Taylor approximations etc... to see that  $R_k^L(\bar{a}^L)$  is split up into each neuron in  $\bar{a}^L$ 's indicating influence of each neuron on  $R_k^L$ .  
 $\hookrightarrow$  This is what LRP back propagates.

In summary, LRP performs approximate Taylor expansions at each neuron, using them to determine how much the neuron's relevance was influenced by the neurons in previous layers. Relevance is back-propagated proportionally to this influence. TEA trains the student classifier to match the heatmaps of a teacher classifier. Accordingly, features in the images will influence the student and the teacher in a similar manner. I.e., the reasons behind the TEA student decisions will match the reasons behind the teacher decisions. If the teacher pays no attention to bias, no LRP relevance flows to bias in its heatmaps. Similarly, no LRP will flow to bias in the student's heatmap, making it also pay no attention to bias.

\* Why training to match LRP o/p's guarantees correct outputs?

Under full conservativeness condition, LRP redistributes the relevance  $R_k^L$  at neuron  $k$  (layer  $L$ ) to the neurons  $j$  in layer  $L-1$  with no loss or gain of relevance. I.e., when being back-propagated, relevance is never destroyed or created, only redistributed across neurons. Since relevance starts as the value of a classifier logit (Eq. A3), the sum of all neuron's relevances in any layer matches the classifier output (logit) where the back-propagation started:

$$\sum_k R_k^L = z_c^N \forall L \in [0, \dots, N] \quad (\text{A11})$$

Therefore, under perfect conservativeness, LRP provides a perfect guarantee of correct outputs: if the LRP heatmaps (LRP heatmaps are  $\mathbf{R}^0 = [R_k^0]$  in Eq. A11) of two classifiers perfectly match, for all classes  $c$ , they will have the same logits  $z_c^N$ . However, setting  $\varepsilon$  to zero creates noisy heatmaps for deep neural networks and high-resolution images, making the convergence of the TEA less difficult (Tab. 5). Also

## \* Why TEA works better with LRP than GRAD-CAM

- GRAD-CAM has an issue called spurious-mapping: mapping input features to deep activations not aligned with these features.

↳ The deep layers GRAD-CAM is applied on have huge receptive fields, can even take up the whole image.