

A VULNERABILITY OF ATTRIBUTION METHODS USING PRE-SOFTMAX SCORES

MIGUEL LERMA¹ AND MIRTHA LUCAS²

ABSTRACT. We discuss a vulnerability involving a category of attribution methods used to provide explanations for the outputs of convolutional neural networks working as classifiers. It is known that this type of networks are vulnerable to adversarial attacks, in which imperceptible perturbations of the input may alter the outputs of the model. In contrast, here we focus on effects that small modifications in the model may cause on the attribution method without altering the model outputs.

1. INTRODUCTION

The black box nature of current artificial intelligence (AI) models is considered problematic in areas with low tolerance to errors, such as Computer Aided Diagnosis (CAD) and autonomous vehicles. To palliate the effect of mistakes and increase confidence in the model, explanation methods have been developed to justify the model outputs [2].

A class of explanation methods widely used on convolutional neural networks (CNN) take the form of attribution methods that determine how much different parts of the input of a model contribute to produce its final output. In general, the networks on which these methods are used consist of several convolutional layers that produce a vector of outputs $\mathbf{z} = (z_1, z_2, \dots, z_n)$, which is then transformed with a softmax function into a vector of probabilities $\mathbf{y} = (y_1, y_2, \dots, y_n)$, where n is the number of classes. (Figure 1). Each post-softmax output can be interpreted as the amount of confidence about the input sample belonging to each of the several classes $1, 2, \dots, n$. In classification tasks, the output with maximum value corresponds to the class to which the input sample is considered to belong.

Gradient-based attribution methods for convolutional networks work by computing the gradient $\nabla_{\mathbf{x}}S = (\partial S/\partial x_1, \dots, \partial S/\partial x_N)$ of an output or “score” S of the network respect to a set of inputs or unit activations $\mathbf{x} = (x_1, \dots, x_N)$, where N is the number of inputs or internal units, and S may represent either one of the pre-softmax outputs z_i , or one of the post-softmax outputs y_i . The assumption is that each derivative $\partial S/\partial x_i$ provides a measure of the impact of x_i on the score S . A few examples of attribution methods using this approach are Grad-CAM [11], Integrated Gradients (IG) [14], and RSI Grad-CAM [9].

In [8] there is a detailed analysis of the differences between using gradients of pre-softmax versus post-softmax outputs. In that paper it is argued that the post-softmax version of gradient-based methods is more robust and not affected by a vulnerability suffered by the pre-softmax version.

¹NORTHWESTERN UNIVERSITY, EVANSTON, USA

²DEPAUL UNIVERSITY, CHICAGO, USA

E-mail addresses: ¹mлерma@math.northwestern.edu, ²mlucas3@depaul.edu.

Date: October 25, 2023.

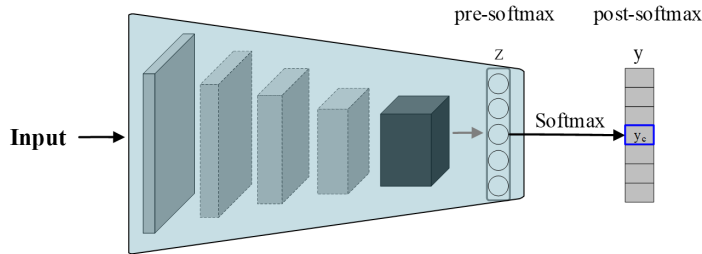


FIGURE 1. Structure of a typical classifier network. After a number of convolutional blocks this kind of network ends with a fully connected network producing a (pre-softmax) output \mathbf{z} , followed by a softmax activation function with (post-softmax) output \mathbf{y} .

Here we will provide a brief overview of the main argument leading to that conclusion, and a way in which the vulnerability could be exploited.

2. PREVIOUS WORK

The possibility of fooling a classification network with *adversarial attacks* by using slightly modified inputs is well known [1, 4]. On the other hand, the possibility of altering the output of an attribution method without modifying the model predictions has not been studied in the same extent, but there are also some findings in that direction. Since terminology may vary across works we must clarify that we use the term *attribution* method where other authors use *explanation* or *interpretation* method. We made this decision to stress the fact that an attribution method may not quite fulfill human expectations for an explanation, in particular Grad-CAM-like methods seem to do a good job in locating the parts of an input containing a sample of a class, i.e., it helps to determine *where* the object corresponding to the class predicted by the model is in the input image, but that does not necessarily explains *why* the output of the network is what it is. However, when citing a work we keep the authors terminology in this regard.

In [3] adversarial attacks against interpretation methods are tried and tested. They work in a similar way to adversarial attacks against network predictions, the main idea is to search for small perturbations of sample inputs that change the output of interpretation methods without altering the network predictions. The work is mainly experimental and requires extensive testing.

The works mentioned above focus on how perturbation of inputs can alter outputs of attribution methods. On the other hand, the authors of [6] study the possibility of fooling interpretation methods by adversarial model manipulation without perturbing model accuracy. Their approach consists of applying fine tuning to a given model with a loss term that includes the interpretation results in the penalty term of the objective function. So, rather than perturbing inputs the approach of the authors is to perturb the model itself. Again, the work is mainly empirical and requires extensive testing.

Concerned with the quality of explanation methods, the authors of [5] have built **Quantus**, a comprehensible tool for XAI evaluation, and they list a number of metrics that can be applied to explanation methods. The metric that is most closely related to our work is *robustness*, which (in their words) *measures to what extent explanations are stable when subject to slight perturbations in the input*,

assuming that the model output approximately stayed the same. As indicated, the metric is based on the effects of perturbations applied to input samples.

Before showing the details of our work we state how it differs from previous work in identifying possible adversarial attacks against attribution methods. First, our work does not require to perturb inputs. Second, our method does not require training or fine tuning a model. We just identify a vulnerability of Grad-CAM-like methods using pre-softmax scores, and show how the model can be modified to exploit the vulnerability. Going beyond the theory we show an specific modification that has the desired effect, and illustrate it with several examples as a proof of concept.

3. A VULNERABILITY OF ATTRIBUTION METHODS USING PRE-SOFTMAX SCORES.

In this section we examine a vulnerability that affects attribution methods for CNNs that work with pre-softmax scores, with a special emphasis on gradient-based methods, although many of the considerations can be easily extended to methods that work with finite differences rather than gradients, such as Layer-wise Relevance Propagation (LRP) [10] and DeepLIFT [12].

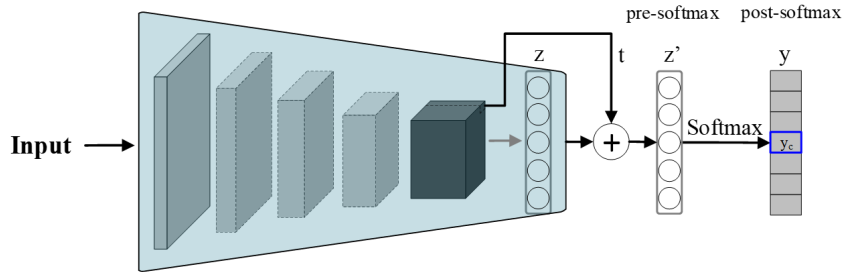


FIGURE 2. Example of alteration of a classifier network that changes attributions based on pre-softmax scores without changing post-softmax scores.

3.1. The softmax function. The output of the softmax function applied to a vector $\mathbf{z} = (z_1, z_2, \dots, z_n)$ is the vector $\mathbf{y} = (y_1, y_2, \dots, y_n)$ whose components are:

$$(1) \quad y_c = \frac{e^{z_c}}{\sum_{i=1}^n e^{z_i}}.$$

The outputs of the softmax verify $0 < y_c < 1$ for all classes $c = 1, \dots, n$, and $\sum_{c=1}^n y_c = 1$, so the y_c are usually interpreted as probabilities.

Note that adding an amount t independent of the class i to all the arguments of the softmax, $z'_i = z_i + t$, has no effect on its outputs:

$$(2) \quad \begin{aligned} y'_c &= \frac{e^{z'_c}}{\sum_{i=1}^n e^{z'_i}} = \frac{e^{z_c+t}}{\sum_{i=1}^n e^{z_i+t}} = \frac{e^t e^{z_c}}{\sum_{i=1}^n e^t e^{z_i}} \\ &= \frac{e^t e^{z_c}}{e^t \sum_{i=1}^n e^{z_i}} = \frac{e^{z_c}}{\sum_{i=1}^n e^{z_i}} = y_c. \end{aligned}$$

So, the change $z_i \mapsto z_i + t$ for every i does not change the network post-softmax outputs y_c . Note that t does not need to be a constant, all that is required is that t is independent of i .

Since adding t has no effect in the output of the softmax, the derivatives of the outputs of the softmax won't change after adding t to its arguments:

$$(3) \quad \frac{\partial y'_i}{\partial x} = \frac{\partial y_i}{\partial x},$$

however the derivatives of the pre-softmax z_i may change:

$$(4) \quad \frac{\partial z'_i}{\partial x} = \frac{\partial(z'_i + t)}{\partial x} = \frac{\partial z_i}{\partial x} + \frac{\partial t}{\partial x},$$

so that $\frac{\partial z'_i}{\partial x} \neq \frac{\partial z_i}{\partial x}$ if $\frac{\partial t}{\partial x} \neq 0$.

This theoretical result and its potential impact in gradient-based attribution methods are carefully examined in [8]. In the following section we will provide a proof of concept showing how this results can be used to radically modify a heatmap produced by an attribution method such as Grad-CAM.

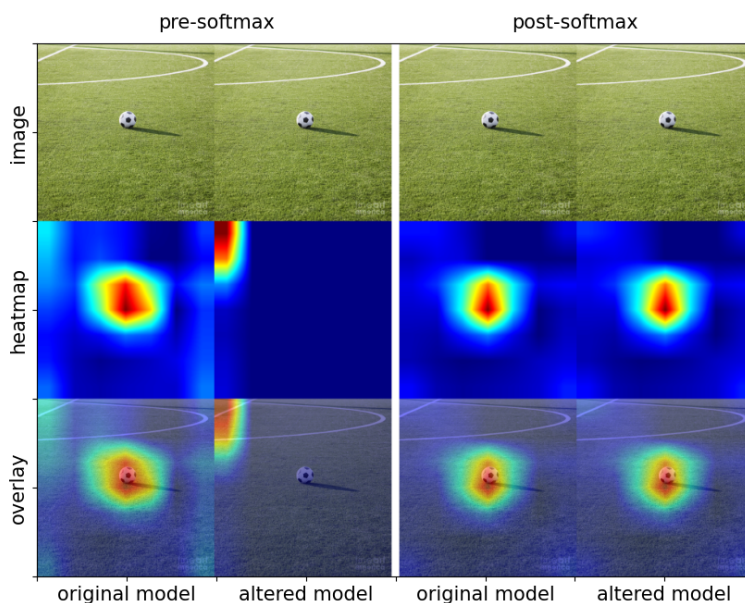


FIGURE 3. Heatmaps produced by Grad-CAM using pre-softmax and post-softmax outputs respectively, intended to locate the position of the soccer ball. The original model is a VGG19 network pretrained on ImageNet. The altered model is the same VGG19 network slightly modified, but still functionally equivalent (same final outputs) to the original network. The heatmaps are computed at the last convolutional layer of each model. Note that Grad-CAM working on pre-softmax outputs has been tricked to produce wrong heatmaps. The heatmaps obtained using post-softmax outputs remain unchanged.

3.2. A vulnerability of attribution methods using pre-softmax scores. Equation (2) shows that the softmax function has no unique inverse because we can add to its arguments z_1, \dots, z_n any scalar t independent of i without changing the output of the softmax.

In the example shown here (Figure 2) the network is a VGG19 pretrained on ImageNet [13]. Then, t is the result of adding the activations of the units placed in position $(0, 0)$ of the final pool layer

(block5_pool) across all its channels multiplied by a constant K . More specifically, if A_{ijk} presents the activation of unit in position (i, j) of channel k of the last pooling layer, then:

$$(5) \quad t = K \sum_k A_{00k},$$

where K is a constant—in our experiment we used $K = 10$.

After t is added to the original z_i pre-softmax scores of the network we get new pre-softmax scores $z'_i = z_i + t$. This makes the new pre-softmax scores strongly dependent on the units in position $(0, 0)$ of the final pool layer without altering the post-softmax scores of the network. Consequently, we expect that heatmaps produced by Grad-CAM to strongly highlight the upper left area of the image regardless of whether that part of the image is related to the network final output.

Figures 3-5 show that, for the altered model, the heatmaps produced using pre-softmax scores are strongly distorted, while the heatmaps produced using post-softmax scores remain unchanged.

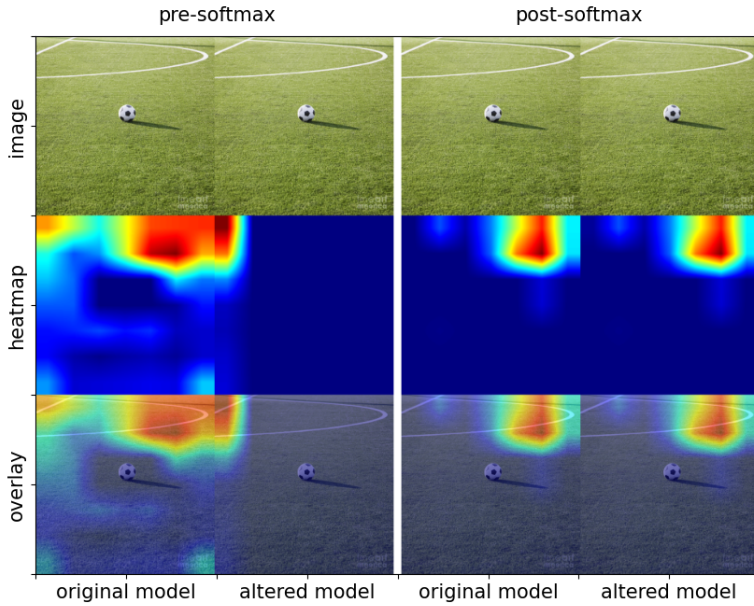


FIGURE 4. The altered model tends to produce the same heatmap regardless of the class assigned to the image. In this case Grad-CAM is used to locate a “maze” rather than a soccer ball in the image. The pre-softmax version of the heatmap on the altered model keeps highlighting the same upper left corner, while the other heatmaps focus on the lines drawn on the grass.

On the other hand, since the final (post-softmax) output of the network remains unchanged, the loss function used for training would sit on the same local minimum for both models (original and modified). Further training of the models won’t make a difference since the added connection cannot backpropagate error. More specifically, if E is the loss function used for training, then for the modified model we have (using multivariate chain rule):

$$(6) \quad \frac{\partial E}{\partial t} = \sum_{i=1}^n \frac{\partial E}{\partial y'_i} \frac{\partial y'_i}{\partial t} = 0$$

because $y'_i = y_i$, which does not depend on t , hence $\frac{\partial y'_i}{\partial t} = \frac{\partial y_i}{\partial t} = 0$ for all i . Consequently, the trainable parameters of both models would change in the same way, and if the error function E is

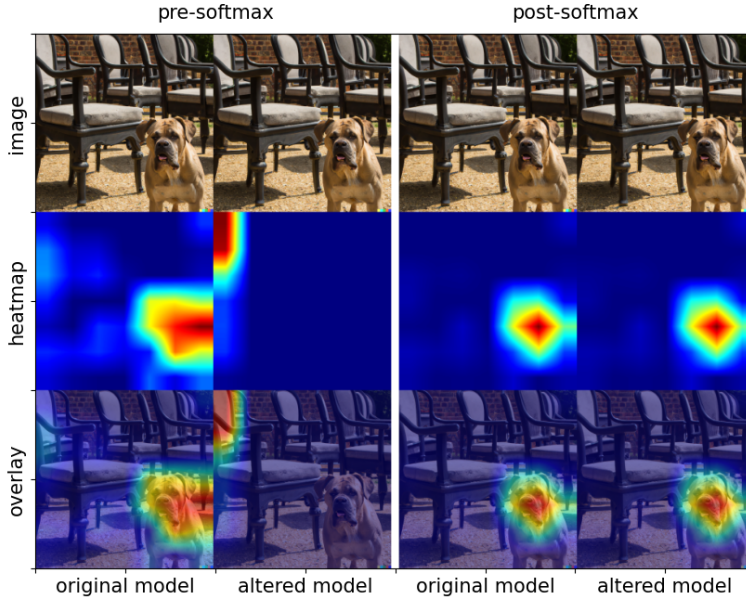


FIGURE 5. Another example showing the heatmap computed with pre-softmax outputs of the altered model concentrated in the upper left corner of the image. Heatmaps computed with post-softmax outputs remain unaltered highlighting the position of the dog.

at or near a minimum for the original model, the same would hold for the modified model. Also, if we trained the modified VGG19 network from scratch and with the same parameter initialization, the final trainable parameters would be the same as those of the original VGG19.

4. DISCUSSION

We note that the main property behind the vulnerability shown here is the possibility of altering pre-softmax scores of a classifier CNN without altering its post-softmax scores. One question could be whether this vulnerability can be exploited to deploy a malicious attack intended to undermine confidence in the model. This kind of attack would be available for anybody having access to model repositories. Since after modification the new model would be functionally equivalent to the original one (its outputs will not change) it would be hard to notice that it has been modified. Also, it is conceivable that the problem pointed out may manifest itself in an unintended way because, after training, both the original and modified model may end up at the same local minimum of the loss function used for training.

The phenomenon discussed may seem to have some similarities with *Clever Hans* effects [7], which also causes heatmaps to highlight wrong areas of the input. *Clever Hans* effects are due to the ability of a classifier to exploit spurious or artifactual correlations. For instance, in a dataset in which images of horses contain a watermark, the model may learn to correctly classify the image of a horse by paying attention only to the presence of the watermark rather than the horse. In that case, an appropriate attribution method would consistently highlight the area of the watermark in the images with horses, which is outside the actual area of interest. However, that would not happen because of a problem in the attribution method, which would be correctly revealing a problem with the model (trained with a biased dataset). On the contrary, the vulnerability discussed here tells

nothing about the ability of the model to extract the right information from the right parts of its inputs, it only depends on the fact that the gradients of the pre-softmax scores may not provide the right information to determine the impact of the inputs on the final (post-softmax) outputs.

5. CONCLUSIONS

We have shown that attribution methods using pre-softmax scores are vulnerable to a class of adversarial attacks that may modify the heatmaps produced without changing the model outputs. Post-softmax outputs are not vulnerable to this kind of attack. We have also noted that the vulnerability discussed here is not a Clever Hans effect. Future work can be used to determine in what extent the problem applies to a wider class of attribution methods.

REFERENCES

- [1] Naveed Akhtar and Ajmal Mian. Threat of adversarial attacks on deep learning in computer vision: A survey. *IEEE Access*, 6:14410–14430, 2018. [2](#)
- [2] Nadia Burkart and Marco F. Huber. A survey on the explainability of supervised machine learning. *J. Artif. Int. Res.*, 70:245–317, may 2021. [1](#)
- [3] Amirata Ghorbani, Abubakar Abid, and James Zou. Interpretation of neural networks is fragile. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):3681–3688, Jul. 2019. [2](#)
- [4] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *CoRR*, abs/1412.6572, 2014. [2](#)
- [5] Anna Hedström, Leander Weber, Daniel Krakowczyk, Dilyara Bareeva, Franz Motzkus, Wojciech Samek, Sebastian Lapuschkin, and Marina M.-C. Höhne. Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond. *Journal of Machine Learning Research*, 24(34):1–11, 2023. [2](#)
- [6] Juyeon Heo, Sunghwan Joo, and Taesup Moon. Fooling neural network interpretations via adversarial model manipulation. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. [2](#)
- [7] Sebastian Lapuschkin, Alexander Binder, Grégoire Montavon, Klaus-Robert Müller, and Wojciech Samek. Analyzing classifiers: Fisher vectors and deep neural networks. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2912–2920, 2016. [6](#)
- [8] Miguel Lerma and Mirtha Lucas. Pre or post-softmax scores in gradient-based attribution methods, what is best? In *2023 IEEE 13th International Conference on Pattern Recognition Systems (ICPRS)*, pages 1–4, 2023. [1](#), [4](#)
- [9] Mirtha Lucas, Miguel Lerma, Jacob Furst, and Daniela Raicu. Rsi-grad-cam: Visual explanations from deep networks via riemann-stieltjes integrated gradient-based localization. In *Advances in Visual Computing: 17th International Symposium, ISVC 2022, San Diego, CA, USA, October 3–5, 2022, Proceedings, Part I*, page 262–274, Berlin, Heidelberg, 2022. Springer-Verlag. [1](#)
- [10] Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, Wojciech Samek, and Klaus-Robert Müller. *Layer-Wise Relevance Propagation: An Overview*, pages 193–209. Springer International Publishing, Cham, 2019. [3](#)
- [11] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 618–626, 2017. [1](#)
- [12] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 3145–3153. JMLR.org, 2017. [3](#)
- [13] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. The 3rd International Conference on Learning Representations (ICLR2015). [4](#)
- [14] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 3319–3328. JMLR.org, 2017. [1](#)