

Explanations can be manipulated and geometry is to blame

Ann-Kathrin Dombrowski¹, Maximilian Alber¹, Christopher J. Anders¹,
Marcel Ackermann², Klaus-Robert Müller^{1,3,4}, Pan Kessel¹

¹Machine Learning Group, EE & Computer Science Faculty, TU-Berlin

²Department of Video Coding & Analytics, Fraunhofer Heinrich-Hertz-Institute

³Max Planck Institute for Informatics

⁴Department of Brain and Cognitive Engineering, Korea University

{klaus-robert.mueller, pan.kessel}@tu-berlin.de

Abstract

Explanation methods aim to make neural networks more trustworthy and interpretable. In this paper, we demonstrate a property of explanation methods which is disconcerting for both of these purposes. Namely, we show that explanations can be manipulated *arbitrarily* by applying visually hardly perceptible perturbations to the input that keep the network's output approximately constant. We establish theoretically that this phenomenon can be related to certain geometrical properties of neural networks. This allows us to derive an upper bound on the susceptibility of explanations to manipulations. Based on this result, we propose effective mechanisms to enhance the robustness of explanations.



Figure 1: Original image with corresponding explanation map on the left. Manipulated image with its explanation on the right. The chosen target explanation was an image with a text stating "this explanation was manipulated".

1 Introduction

Explanation methods have attracted significant attention over the last years due to their promise to open the black box of deep neural networks. Interpretability is crucial for scientific understanding and safety critical applications.

Explanations can be provided in terms of explanation maps [1–19] that visualize the relevance attributed to each input feature for the overall classification result. In this work, we establish that these explanation maps can be changed to an *arbitrary target map*. This is done by applying a visually hardly perceptible perturbation to the input. We refer to Figure 1 for an example. This perturbation does not change the output of the neural network, i.e. in addition to the classification result also the vector of all class probabilities is (approximately) the same.

This finding is clearly problematic if a user, say a medical doctor, is expecting a robustly interpretable explanation map to rely on in the clinical decision making process.

Motivated by this unexpected observation, we provide a theoretical analysis that establishes a relation of this phenomenon to the geometry of the neural network’s output manifold. This novel understanding allows us to derive a bound on the degree of possible manipulation of the explanation map. This bound is proportional to two differential geometric quantities: the principle curvatures and the geodesic distance between the original input and its manipulated counterpart. Given this theoretical insight, we propose efficient ways to limit possible manipulations and thus enhance resilience of explanation methods.

In summary, this work provides the following key contributions:

- We propose an algorithm which allows to manipulate an image with a hardly perceptible perturbation such that the explanation matches an arbitrary target map. We demonstrate its effectiveness for six different explanation methods and on four network architectures as well as two datasets.
- We provide a theoretical understanding of this phenomenon for gradient-based methods in terms of differential geometry. We derive a bound on the principle curvatures of the hypersurface of equal network output. This implies a constraint on the maximal change of the explanation map due to small perturbations.
- Using these insights, we propose methods to undo the manipulations and increase the robustness of explanation maps by smoothing the explanation method. We demonstrate experimentally that smoothing leads to increased robustness not only for gradient but also for propagation-based methods.

1.1 Related work

In [20], it was demonstrated that explanation maps can be sensitive to small perturbations in the image. Their results may be thought of as untargeted manipulations, i.e. perturbations to the image which lead to an unstructured change in the explanation map. Our work focuses on targeted manipulations instead, i.e. to reproduce a given target map. Another approach [21] adds a constant shift to the input image, which is then eliminated by changing the bias of the first layer. For some methods, this leads to a change in the explanation map. Contrary to our approach, this requires to change the network’s biases. In [22], explanation maps are changed by randomization of (some of) the network weights. This is different from our method as it does not aim to change the explanation in a targeted manner and modifies the weights of the network.

2 Manipulating explanations

We consider a neural network $g : \mathbb{R}^d \rightarrow \mathbb{R}^K$ with relu non-linearities which classifies an image $x \in \mathbb{R}^d$ in K categories with the predicted class given by $k = \arg \max_i g(x)_i$. The

explanation map is denoted by $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and associates an image with a vector of the same dimension whose components encode the relevance score of each pixel for the neural network’s prediction. For a given explanation method and specified target $h^t \in \mathbb{R}^d$, a manipulated image $x_{adv} = x + \delta x$ has the following properties:

1. The output of the network stays approximately constant, i.e. $g(x_{adv}) \approx g(x)$.
2. The explanation is close to the target map, i.e. $h(x_{adv}) \approx h^t$.
3. The norm of the perturbation δx added to the input image is small, i.e. $\|\delta x\| = \|x_{adv} - x\| \ll 1$ and therefore not perceptible.

Throughout this paper, we will use the following explanation methods:

- **Gradient:** The map $h(x) = \frac{\partial g}{\partial x}(x)$ is used and quantifies how infinitesimal perturbations in each pixel change the prediction $g(x)$ [1, 2].
- **Gradient \times Input:** This method uses the map $h(x) = x \odot \frac{\partial g}{\partial x}(x)$ [14]. For linear models, this measure gives the exact contribution of each pixel to the prediction.
- **Integrated Gradients:** This method defines $h(x) = (x - \bar{x}) \odot \int_0^1 \frac{\partial g(\bar{x} + t(x - \bar{x}))}{\partial x} dt$ where \bar{x} is a suitable baseline. See the original reference [13] for more details.
- **Guided Backpropagation (GBP):** This method is a variation of the gradient explanation for which negative components of the gradient are set to zero while backpropagating through the non-linearities [4].
- **Layer-wise Relevance Propagation (LRP):** This method [5, 16] propagates relevance backwards through the network. For the output layer, relevance is defined by¹

$$R_i^L = \delta_{i,k}, \tag{1}$$

which is then propagated backwards through all layers but the first using the z^+ rule

$$R_i^l = \sum_j \frac{x_i^l (W^l)_{ji}^+}{\sum_i x_i^l (W^l)_{ji}^+} R_j^{l+1}, \tag{2}$$

where $(W^l)^+$ denotes the positive weights of the l -th layer and x^l is the activation vector of the l -th layer. For the first layer, we use the z^B rule to account for the bounded input domain

$$R_i^0 = \sum_j \frac{x_j^0 W_{ji}^0 - l_j (W^0)_{ji}^+ - h_j (W^0)_{ji}^-}{\sum_i (x_j^0 W_{ji}^0 - l_j (W^0)_{ji}^+ - h_j (W^0)_{ji}^-)} R_j^1, \tag{3}$$

where l_i and h_i are the lower and upper bounds of the input domain respectively.

- **Pattern Attribution (PA):** This method is equivalent to standard backpropagation upon element-wise multiplication of the weights W^l with learned patterns A^l . We refer to the original publication for more details [17].

These methods cover two classes of attribution methods, namely *gradient-based* and *propagation-based* explanations, and are frequently used in practice [23, 24].

¹Here we use the Kronecker symbol $\delta_{i,k} = \begin{cases} 1, & \text{for } i = k \\ 0, & \text{for } i \neq k \end{cases}$.

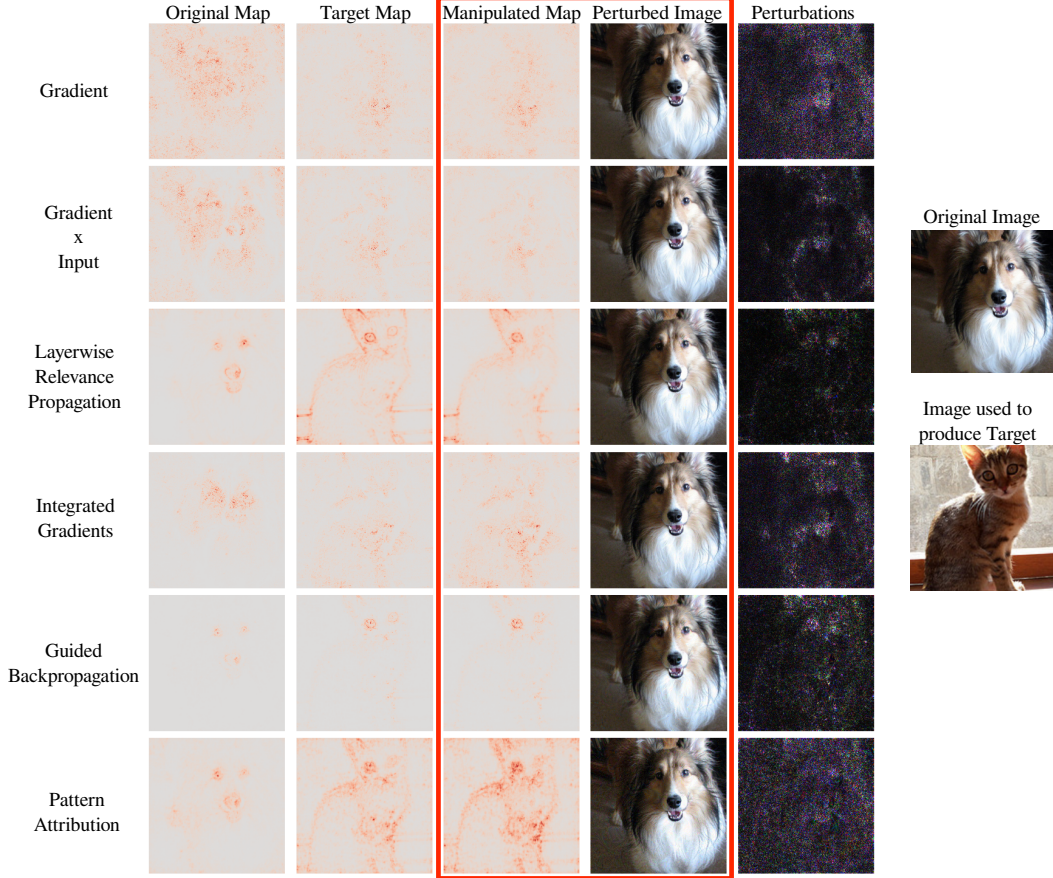


Figure 2: The explanation map of the cat is used as the target and the image of the dog is perturbed. The red box contains the manipulated images and the corresponding explanations. The first column corresponds to the original explanations of the unperturbed dog image. The target map, shown in the second column, is generated with the cat image. The last column visualizes the perturbations.

2.1 Manipulation Method

Let $h^t \in \mathbb{R}^d$ be a given target explanation map and $x \in \mathbb{R}^d$ an input image. As explained previously, we want to construct a manipulated image $x_{adv} = x + \delta x$ such that it has an explanation very similar to the target h^t but the output of the network stays approximately constant, i.e. $g(x_{adv}) \approx g(x)$. We obtain such manipulations by optimizing the loss function

$$\mathcal{L} = \|h(x_{adv}) - h^t\|^2 + \gamma \|g(x_{adv}) - g(x)\|^2, \quad (4)$$

with respect to x_{adv} using gradient descent. We clamp x_{adv} after each iteration so that it is a valid image. The first term in the loss function (4) ensures that the manipulated explanation map is close to the target while the second term encourages the network to have the same output. The relative weighting of these two summands is controlled by the hyperparameter $\gamma \in \mathbb{R}_+$.

The gradient with respect to the input $\nabla h(x)$ of the explanation often depends on the vanishing second derivative of the relu non-linearities. This causes problems during optimization of the loss (4). As an example, the gradient method leads to

$$\partial_{x_{adv}} \|h(x_{adv}) - h^t\|^2 \propto \frac{\partial h}{\partial x_{adv}} = \frac{\partial^2 g}{\partial x_{adv}^2} \propto \text{relu}'' = 0.$$

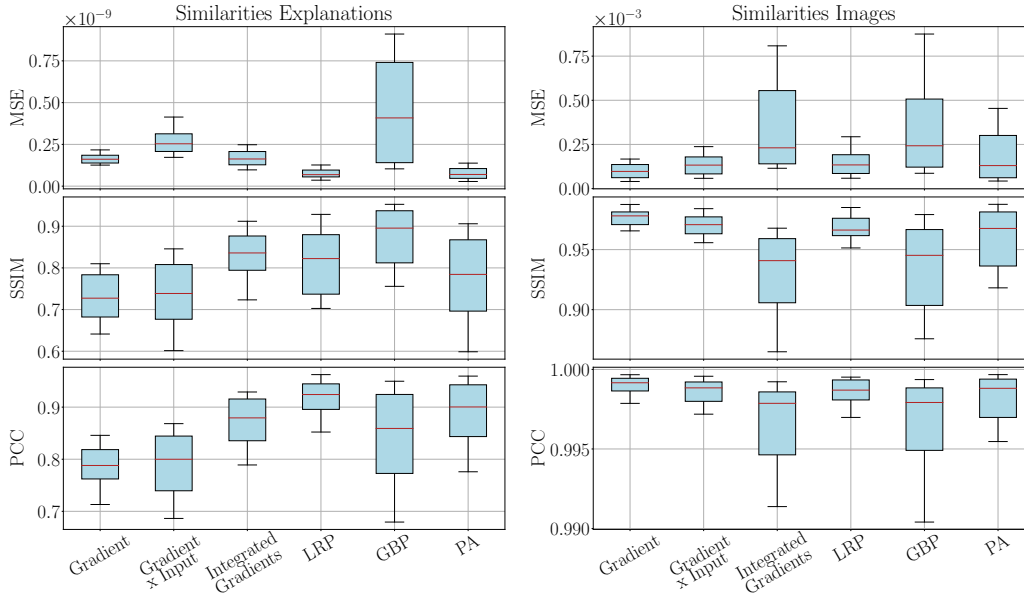


Figure 3: Left: Similarity measures between target h^t and manipulated explanation map $h(x_{adv})$. Right: Similarity measures between original x and perturbed image x_{adv} . For SSIM and PCC large values indicate high similarity while for MSE small values correspond to similar images.

We therefore replace the relu by softplus non-linearities

$$\text{softplus}_\beta(x) = \frac{1}{\beta} \log(1 + e^{\beta x}). \quad (5)$$

For large β values, the softplus approximates the relu closely but has a well-defined second derivative. After optimization is complete, we test the manipulated image with the original relu network.

Similarity metrics: In our analysis, we assess the similarity between both images and explanation maps. To this end, we use three metrics following [22]: the structural similarity index (SSIM), the Pearson correlation coefficient (PCC) and the mean squared error (MSE). SSIM and PCC are relative similarity measures with values in $[0, 1]$, where larger values indicate high similarity. The MSE is an absolute error measure for which values close to zero indicate high similarity. We normalize the sum of the explanation maps to be one and the images to have values between 0 and 1.

2.2 Experiments

To evaluate our approach, we apply our algorithm to 100 randomly selected images for each explanation method. We use a pre-trained VGG-16 network [25] and the ImageNet dataset [26]. For each run, we randomly select two images from the test set. One of the two images is used to generate a target explanation map h^t . The other image is perturbed by our algorithm with the goal of replicating the target h^t using a few thousand iterations of gradient descent. We sum over the absolute values of the channels of the explanation map to get the relevance per pixel. Further details about the experiments are summarized in Supplement A.

Qualitative analysis: Our method is illustrated in Figure 2 in which a dog image is manipulated in order to have an explanation of a cat. For all explanation methods, the target is closely emulated and the perturbation of the dog image is small. More examples can be found in the supplement.

Quantitative analysis: Figure 3 shows similarity measures between the target h^t and the manipulated explanation map $h(x_{adv})$ as well as between the original image x and perturbed image x_{adv} .² All considered metrics show that the perturbed images have an explanation closely resembling the targets. At the same time, the perturbed images are very similar to the corresponding original images. We also verified by visual inspection that the results look very similar. We have uploaded the results of all runs so that interested readers can assess their similarity themselves³ and will provide code to reproduce them. In addition, the output of the neural network is approximately unchanged by the perturbations, i.e. the classification of all examples is unchanged and the median of $\|g(x_{adv}) - g(x)\|$ is of the order of magnitude 10^{-3} for all methods. See Supplement B for further details.

Other architectures and datasets: We checked that comparable results are obtained for ResNet-18 [27], AlexNet [28] and Densenet-121 [29]. Moreover, we also successfully tested our algorithm on the CIFAR-10 dataset [30]. We refer to the Supplement C for further details.

3 Theoretical considerations

In this section, we analyze the vulnerability of explanations theoretically. We argue that this phenomenon can be related to the large curvature of the output manifold of the neural network. We focus on the gradient method starting with an intuitive discussion before developing mathematically precise statements.

We have demonstrated that one can drastically change the explanation map while keeping the output of the neural network constant

$$g(x + \delta x) = g(x) = c \tag{6}$$

using only a small perturbation in the input δx . The perturbed image $x_{adv} = x + \delta x$ therefore lies on the hypersurface of constant network output $S = \{p \in \mathbb{R}^d | g(p) = c\}$.⁴ We can exclusively consider the winning class output, i.e. $g(x) := g(x)_k$ with $k = \arg \max_i g(x)_i$ because the gradient method only depends on this component of the output. Therefore, the hyperplane S is of co-dimension one. The gradient ∇g for every $p \in S$ is normal to this hypersurface. The fact that the normal vector ∇g can be drastically changed by slightly perturbing the input along the hypersurface S suggests that the curvature of S is large.

While the latter statement may seem intuitive, it requires non-trivial concepts of differential geometry to make it precise, in particular the notion of the second fundamental form. We will briefly summarize these concepts in the following (see e.g. [31] for a standard textbook). To this end, it is advantageous to consider a normalized version of the gradient method

$$n(x) = \frac{\nabla g(x)}{\|\nabla g(x)\|}. \tag{7}$$

This normalization is merely conventional as it does not change the relative importance of any pixel with respect to the others. For any point $p \in S$, we define the tangent space $T_p S$ as the vector space spanned by the tangent vectors $\dot{\gamma}(0) = \frac{d}{dt} \gamma(t)|_{t=0}$ of all possible curves $\gamma : \mathbb{R} \rightarrow S$ with $\gamma(0) = p$. For $u, v \in T_p S$, we denote their inner product by $\langle u, v \rangle$. For any $u \in T_p S$, the *directional derivative* is uniquely defined for any choice of γ by

$$D_u f(p) = \left. \frac{d}{dt} f(\gamma(t)) \right|_{t=0} \quad \text{with} \quad \gamma(0) = p \quad \text{and} \quad \dot{\gamma}(0) = u. \tag{8}$$

²Throughout this paper, boxes denote 25th and 75th percentiles, whiskers denote 10th and 90th percentiles, solid lines show the medians and outliers are depicted by circles.

³<https://drive.google.com/drive/folders/1TZewngoevHRuIw6gb5CZDIRrc7EWf5yb?usp=sharing>

⁴It is sufficient to consider the hypersurface S in a neighbourhood of the unperturbed input x .

We then define the *Weingarten map* as⁵

$$L : \begin{cases} T_p S & \rightarrow T_p S \\ u & \mapsto -D_u n(p), \end{cases}$$

where the unit normal $n(p)$ can be written as (7). This map quantifies how much the unit normal changes as we infinitesimally move away from p in the direction u . The *second fundamental form* is then given by

$$\mathcal{L} : \begin{cases} T_p S \times T_p S & \rightarrow \mathbb{R} \\ u, v & \mapsto -\langle v, L(u) \rangle = -\langle v, D_u n(p) \rangle. \end{cases}$$

It can be shown that the second fundamental form is bilinear and symmetric $\mathcal{L}(u, v) = \mathcal{L}(v, u)$. It is therefore diagonalizable with real eigenvalues $\lambda_1, \dots, \lambda_{d-1}$ which are called *principle curvatures*.

We have therefore established the remarkable fact that the sensitivity of the gradient map (7) is described by the principle curvatures, a key concept of differential geometry.

In particular, this allows us to derive an upper bound on the maximal change of the gradient map $h(x) = n(x)$ as we move slightly on S . To this end, we define the *geodesic distance* $d_g(p, q)$ of two points $p, q \in S$ as the length of the shortest curve on S connecting p and q . In the supplement, we show that:

Theorem 1 *Let $g : \mathbb{R}^d \rightarrow \mathbb{R}$ be a network with softplus $_{\beta}$ non-linearities and $\mathcal{U}_{\epsilon}(p) = \{x \in \mathbb{R}^d; \|x - p\| < \epsilon\}$ an environment of a point $p \in S$ such that $\mathcal{U}_{\epsilon}(p) \cap S$ is fully connected. Let g have bounded derivatives $\|\nabla g(x)\| \geq c$ for all $x \in \mathcal{U}_{\epsilon}(p) \cap S$. It then follows for all $p_0 \in \mathcal{U}_{\epsilon}(p) \cap S$ that*

$$\|h(p) - h(p_0)\| \leq |\lambda_{max}| d_g(p, p_0) \leq \beta C d_g(p, p_0), \quad (9)$$

where λ_{max} is the principle curvature with the largest absolute value for any point in $\mathcal{U}_{\epsilon}(p) \cap S$ and the constant $C > 0$ depends on the weights of the neural network.

This theorem can intuitively be motivated as follows: for relu non-linearities, the lines of equal network output are piece-wise linear and therefore have kinks, i.e. points of divergent curvature. These relu non-linearities are well approximated by softplus non-linearities (5) with large β . Reducing β smoothes out the kinks and therefore leads to reduced maximal curvature, i.e. $|\lambda_{max}| \leq \beta C$. For each point on the geodesic curve connecting p and p_0 , the normal can at worst be affected by the maximal curvature, i.e. the change in explanation is bounded by $|\lambda_{max}| d_g(p, p_0)$.

There are two important lessons to be learned from this theorem: the geodesic distance can be substantially greater than the Euclidean distance for curved manifolds. In this case, inputs which are very similar to each other, i.e. the Euclidean distance is small, can have explanations that are drastically different. Secondly, the upper bound is proportional to the β parameter of the softplus non-linearity. Therefore, smaller values of β provably result in increased robustness with respect to manipulations.

4 Robust explanations

Using the fact that the upper bound of the last section is proportional to the β parameter of the softplus non-linearities, we propose β -smoothing of explanations. This method calculates an explanation using a network for which the relu non-linearities are replaced by softplus

⁵The fact that $D_u n(p) \in T_p S$ follows by taking the directional derivative with respect to u on both sides of $\langle n, n \rangle = 1$.

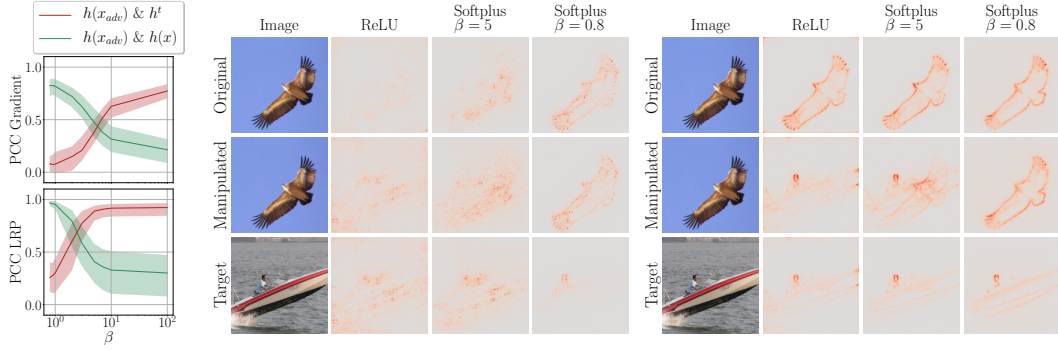


Figure 4: Left: β dependence for the correlations of the manipulated explanation (here Gradient and LRP) with the target and original explanation. Lines denote the medians, 10th and 90th percentiles are shown in semitransparent colour. Center and Right: network input and the respective explanation maps as β is decreased for Gradient (center) and LRP (right).

with a small β parameter to smooth the principle curvatures. The precise value of β is a hyperparameter of the method, but we find that a value around one works well in practice.

As shown in the supplement, a relation between SmoothGrad [12] and β -smoothing can be proven for a one-layer neural network:

Theorem 2 For a one-layer neural network $g(x) = \text{relu}(w^T x)$ and its β -smoothed counterpart $g_\beta(x) = \text{softplus}_\beta(w^T x)$, it holds that

$$\mathbb{E}_{\epsilon \sim p_\beta} [\nabla g(x - \epsilon)] = \nabla g_{\frac{\beta}{\|w\|}}(x),$$

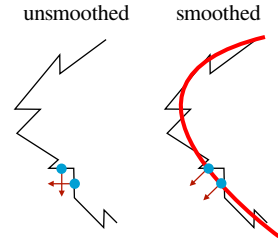
where $p_\beta(\epsilon) = \frac{\beta}{(e^{\beta\epsilon/2} + e^{-\beta\epsilon/2})^2}$.

Since $p_\beta(x)$ closely resembles a normal distribution with variance $\sigma = \log(2) \frac{\sqrt{2\pi}}{\beta}$, β -smoothing can be understood as $N \rightarrow \infty$ limit of SmoothGrad $h(x) = \frac{1}{N} \sum_{i=1}^N \nabla g(x - \epsilon_i)$ where $\epsilon_i \sim g_\beta \approx \mathcal{N}(0, \sigma)$. We emphasize that the theorem only holds for a one-layer neural network, but for deeper networks we empirically observe that both lead to visually similar maps as they are considerably less noisy than the gradient map. The theorem therefore suggests that SmoothGrad can similarly be used to smooth the curvatures and can thereby make explanations more robust.⁶

Experiments: Figure 4 demonstrates that β -smoothing allows us to recover the original explanation map by lowering the value of the β parameter. We stress that this works for all considered methods. We also note that the same effect can be observed using SmoothGrad by successively increasing the standard deviation σ of the noise distribution. This further underlines the similarity between the two smoothing methods.

If an attacker knew that smoothing was used to undo the manipulation, they could try to attack the smoothed method directly. However, both β -smoothing and SmoothGrad are substantially more robust than their non-smoothed counterparts, see Figure 5. It is important to note that β -smoothing achieves this at considerably lower computational cost: β -smoothing only requires a single forward and backward pass, while SmoothGrad requires as many as the number of noise samples (typically between 10 to 50).

We refer to Supplement D for more details on these experiments.



⁶For explanation methods $h(x)$ other than gradient, SmoothGrad needs to be used in a slightly generalized form, i.e. $\frac{1}{N} \sum_{i=1}^N h(x - \epsilon_i)$.

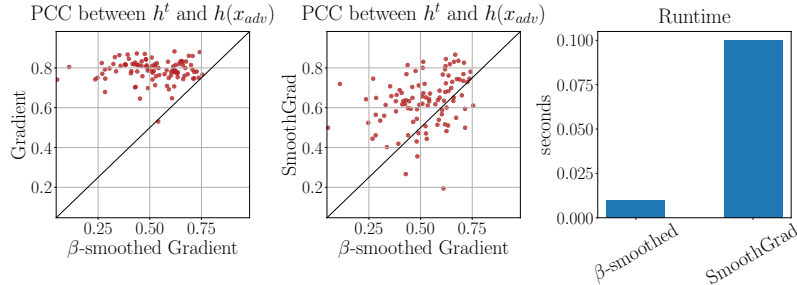


Figure 5: Left: markers are clearly left of the diagonal, i.e. explanations are more robust to manipulations when β -smoothing is used. Center: SmoothGrad has comparable results to β -smoothing, i.e. markers are distributed around the diagonal. Right: β -smoothing has significantly lower computational cost than SmoothGrad.

5 Conclusion

Explanation methods have recently become increasingly popular among practitioners. In this contribution we show that dedicated imperceptible manipulations of the input data can yield arbitrary and drastic changes of the explanation map. We demonstrate both qualitatively and quantitatively that explanation maps of many popular explanation methods can be arbitrarily manipulated. Crucially, this can be achieved while keeping the model’s output constant. A novel theoretical analysis reveals that in fact the large curvature of the network’s decision function is one important culprit for this unexpected vulnerability. Using this theoretical insight, we can profoundly increase the resilience to manipulations by smoothing *only* the explanation process while leaving the model itself unchanged.

Future work will investigate possibilities to modify the training process of neural networks itself such that they can become less vulnerable to manipulations of explanations. Another interesting future direction is to generalize our theoretical analysis from gradient to propagation-based methods. This seems particularly promising because our experiments strongly suggest that similar theoretical findings should also hold for these explanation methods.

References

- [1] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11(Jun):1803–1831, 2010.
- [2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings*, 2014.
- [3] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, pages 818–833, 2014.
- [4] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving for simplicity: The all convolutional net. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Workshop Track Proceedings*, 2015.

- [5] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS one*, 10(7):e0130140, 2015.
- [6] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. Grad-cam: Why did you say that? visual explanations from deep networks via gradient-based localization. *CoRR*, abs/1610.02391, 2016.
- [7] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144. ACM, 2016.
- [8] Luisa M. Zintgraf, Taco S. Cohen, Tameem Adel, and Max Welling. Visualizing deep neural network decisions: Prediction difference analysis. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- [9] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3145–3153, 2017.
- [10] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 4765–4774. Curran Associates, Inc., 2017.
- [11] Piotr Dabkowski and Yarín Gal. Real time image saliency for black box classifiers. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 6967–6976. Curran Associates, Inc., 2017.
- [12] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. *CoRR*, abs/1706.03825, 2017.
- [13] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3319–3328, 2017.
- [14] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 3145–3153, 2017.
- [15] Ruth C Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *2017 IEEE international conference on computer vision (ICCV)*, pages 3449–3457. IEEE, 2017.
- [16] Grégoire Montavon, Sebastian Lapuschkin, Alexander Binder, Wojciech Samek, and Klaus-Robert Müller. Explaining nonlinear classification decisions with deep Taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
- [17] Pieter-Jan Kindermans, Kristof T Schütt, Maximilian Alber, Klaus-Robert Müller, Dumitru Erhan, Been Kim, and Sven Dähne. Learning how to explain neural networks: Patternnet and patternattribution. *International Conference on Learning Representations*, 2018. <https://openreview.net/forum?id=Hkn7CBaTW>.

- [18] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, and Rory Sayres. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). *arXiv preprint arXiv:1711.11279*, 2017.
- [19] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1096, 2019.
- [20] Amirata Ghorbani, Abubakar Abid, and James Y. Zou. Interpretation of neural networks is fragile. *CoRR*, abs/1710.10547, 2017.
- [21] Pieter-Jan Kindermans, Sara Hooker, Julius Adebayo, Maximilian Alber, Kristof T. Schütt, Sven Dähne, Dumitru Erhan, and Been Kim. The (un)reliability of saliency methods. *CoRR*, abs/1711.00867, 2017.
- [22] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. *CoRR*, abs/1810.03292, 2018.
- [23] Maximilian Alber, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T. Schütt, Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter-Jan Kindermans. investigate neural networks! *Journal of Machine Learning Research* **20**, 2019.
- [24] Marco Ancona, Enea Ceolini, Cengiz Oztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *6th International Conference on Learning Representations (ICLR 2018)*, 2018.
- [25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*, 2014.
- [26] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015.
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [28] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS'12*, pages 1097–1105, USA, 2012. Curran Associates Inc.
- [29] Gao Huang, Zhuang Liu, and Kilian Q. Weinberger. Densely connected convolutional networks. *CoRR*, abs/1608.06993, 2016.
- [30] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009.
- [31] Loring W Tu. *Differential geometry: connections, curvature, and characteristic classes*, volume 275. Springer, 2017.

Supplement

Contents

A	Details on experiments	12
A.1	Beta growth	13
B	Difference in network output	14
C	Generalization over architectures and data sets	15
C.1	Additional architectures	15
C.2	Additional datasets	16
D	Smoothing explanation methods	20
E	Proofs	25
E.1	Theorem 1	25
E.2	Theorem 2	27
F	Additional examples for VGG	28

A Details on experiments

We provide a *run_attack.py* file in our reference implementation which allows one to produce manipulated images. The hyperparameter choices used in our experiments are summarized in Table 1. We set $\beta_0 = 10$ and $\beta_e = 100$ for beta growth (see section below for a description). The column 'factors' summarizes the weighting of the mean squared error of the heatmaps and the images respectively.

method	iterations	lr	factors
Gradient	1500	10^{-3}	$10^{11}, 10^6$
Grad x Input	1500	10^{-3}	$10^{11}, 10^6$
IntGrad	500	5×10^{-3}	$10^{11}, 10^6$
LRP	1500	2×10^{-4}	$10^{11}, 10^6$
GBP	1500	10^{-3}	$10^{11}, 10^6$
PA	1500	2×10^{-3}	$10^{11}, 10^6$

Table 1: Hyperparameters used in our analysis.

The patterns for explanation method PA are trained on a subset of the ImageNet training set. The baseline \bar{x} for explanation method IG was set to zero. To approximate the integral, we use 30 steps for which we verified that the attributions approximately adds up to the score at the input.

A.1 Beta growth

In practise, we observe that we get slightly better results by increasing the value of β of the softplus $sp(x) = \frac{1}{\beta} \ln(1 + e^{\beta x})$ during training a start value β_0 to a final value β_e using

$$\beta(t) = \beta_0 \left(\frac{\beta_e}{\beta_0} \right)^{t/T}, \quad (10)$$

where t is the current optimization step and T denotes the total number of steps. Figure 6 shows the MSE for images and explanation maps during training with and without β -growth. This strategy is however not essential for our results.

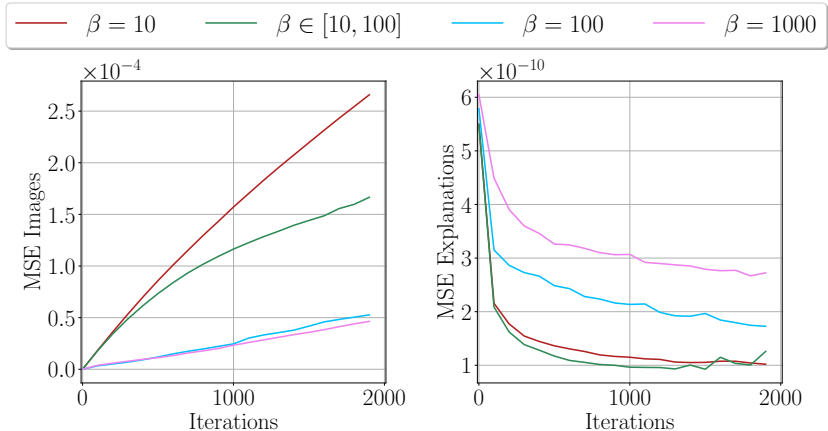


Figure 6: MSE between x and x_{adv} (left) and between h^t and $h(x_{adv})$ (right) for various values for β .

We use beta growth for all methods except LRP for which we do not find any speed-up in the optimization as the LRP rules do not explicitly depend on the second derivative of the relu activations. Figure 7 demonstrates that for large beta values the softplus networks approximate the relu network well. Figure 8 and Figure 9 show this for an example for the gradient and the LRP explanation method. We also note that for small beta the gradient explanation maps become more similar to LRP/GPB/PA explanation maps.

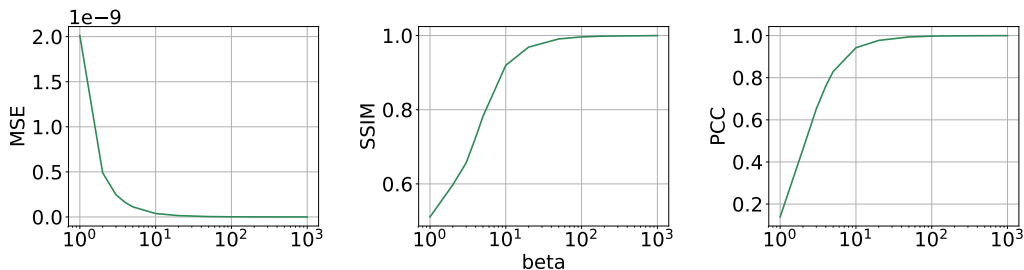


Figure 7: Error measures between the gradient explanation map produced with the original network and explanation maps produced with a network with softplus activation functions using various values for β .

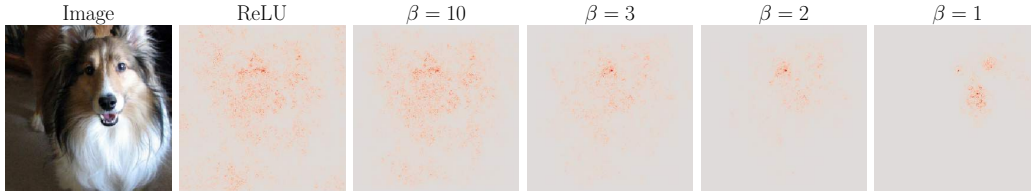


Figure 8: Gradient explanation map produced with the original network and a network with softplus activation functions using various values for β .

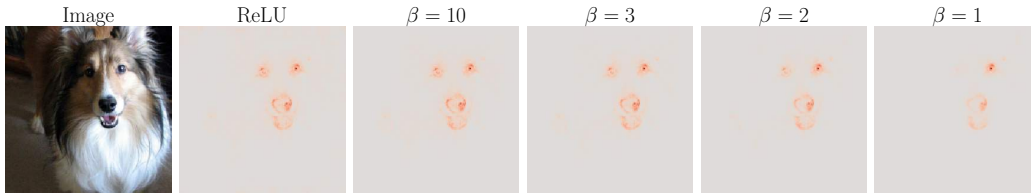


Figure 9: LRP explanation map produced with the original network and a network with softplus activation functions using various values for β .

B Difference in network output

Figure 10 summarizes the change in the output of the network due to the manipulation. We note that all images have the same classification result as the originals. Furthermore, we note that the change in confidence is small. Last but not least, norm of the vector of all class probabilities is also very small.

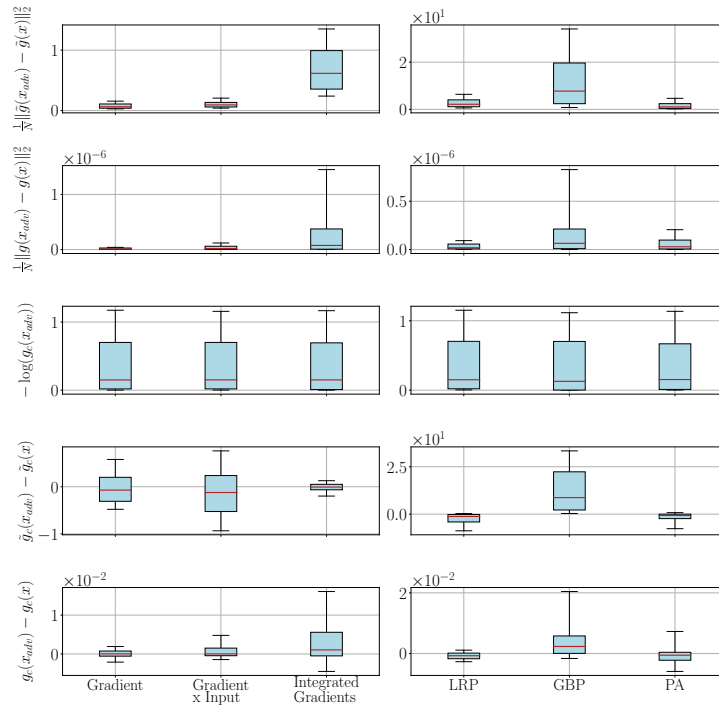


Figure 10: Error analysis of Network output. $\tilde{g}(x)$ denotes pre-activation of the last layer. $g(x)$ is the network output after applying the softmax function to the pre-activation $\tilde{g}(x)$.

C Generalization over architectures and data sets

Manipulable explanations are not only a property of the VGG-16 network. In this section, we show that our algorithm to manipulate explanations can also be applied to other architectures and data sets. For the experiments, we optimize the loss function given in the main text. We keep the pre-activation for all network architectures approximately constant, which also leads to approximately constant activation.

C.1 Additional architectures

In addition to the VGG architecture we also analyzed the explanation’s susceptibility to manipulations for the AlexNet, Densenet and ResNet architectures. The hyperparameter choices used in our experiments are summarized in Table 2. We set $\beta_0 = 10$ and $\beta_e = 100$ for beta growth. Only for Densenet we set $\beta_0 = 30$ and $\beta_e = 300$ as for smaller beta values the explanation map produced with softplus does not resemble the explanation map produced with relu. Figure 12 and 11 show that the similarity measures are comparable for all network architectures for the gradient method.

Figure 13, 15, 16 and 14 show one example image for each architecture.

network	iterations	lr	factors
VGG16	1500	10^{-3}	1e11, 10
AlexNet	4000	10^{-3}	1e11, 10
Densenet-121	2000	5×10^{-4}	1e11, 10
ResNet-18	2000	10^{-3}	1e11, 10

Table 2: Hyperparameters used in our analysis for all networks.

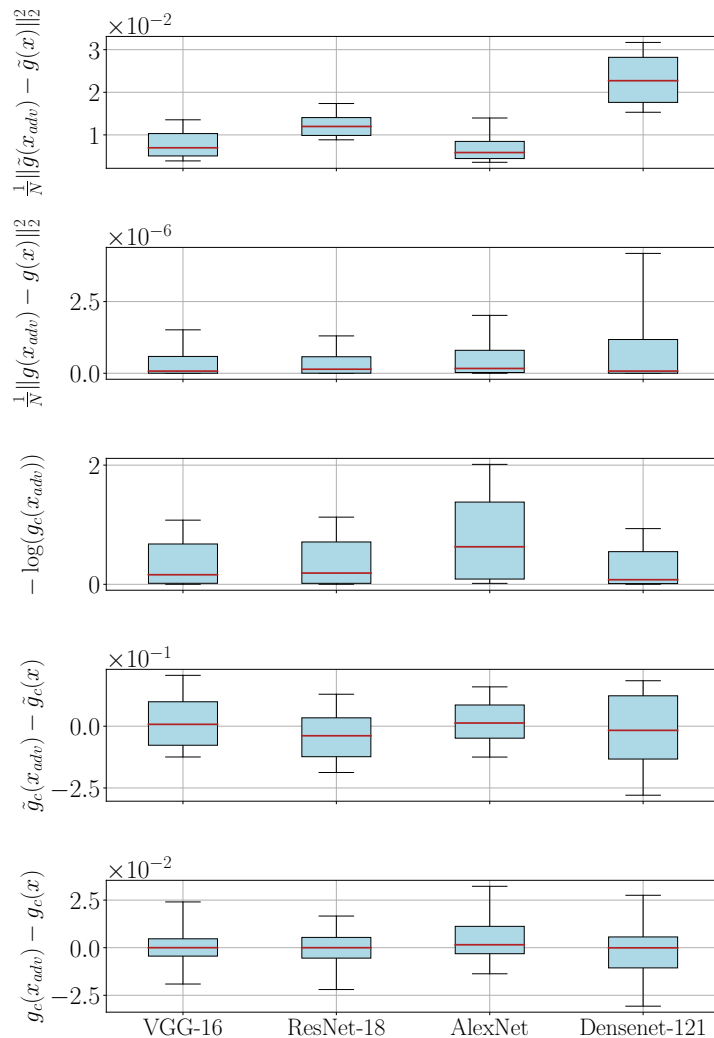


Figure 11: . Change in output for various architectures.

C.2 Additional datasets

We trained the VGG-16 architecture on the CIFAR-10 dataset¹. The test accuracy is approximately 92%. We then used our algorithm to manipulate the explanations for the LRP method. The hyperparameters are summarized in Table 3. Two example images can be seen in Figure 17.

method	iterations	lr	factors
LRP	1500	2×10^{-4}	$10^7, 10^2$

Table 3: Hyperparameters used in our analysis for the CIFAR-10 Dataset.

¹code for training VGG on CIFAR-10 from <https://github.com/chengyangfu/pytorch-vgg-cifar10>

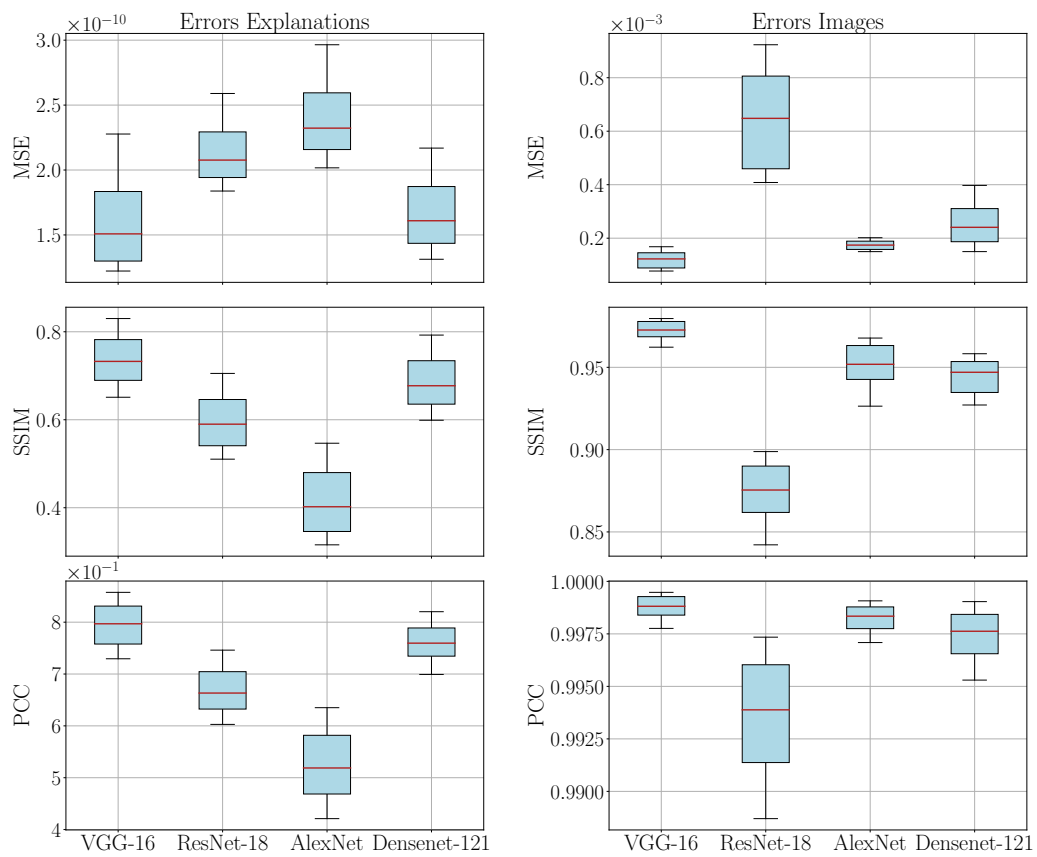


Figure 12: Similarity measures for gradient method for various architectures.

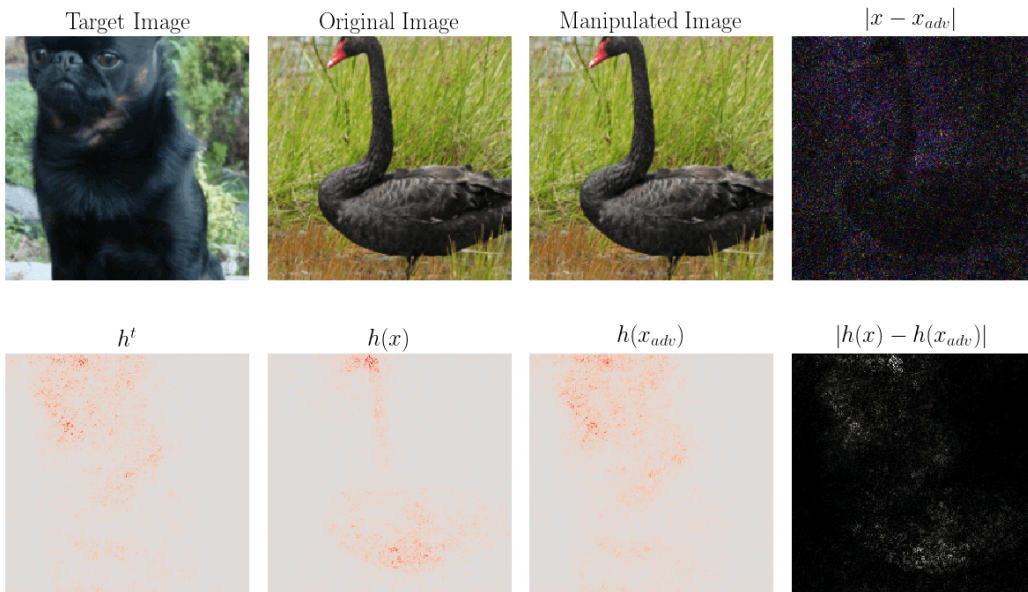


Figure 13: Gradient explanation maps produced with VGG-16 model.

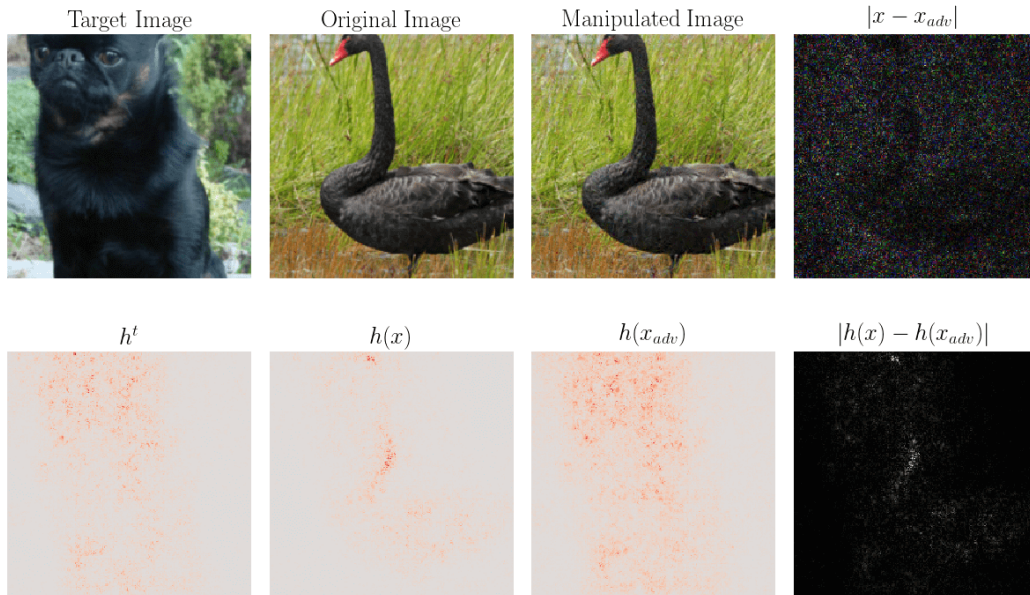


Figure 14: Gradient explanation maps produced with ResNet-18 model.

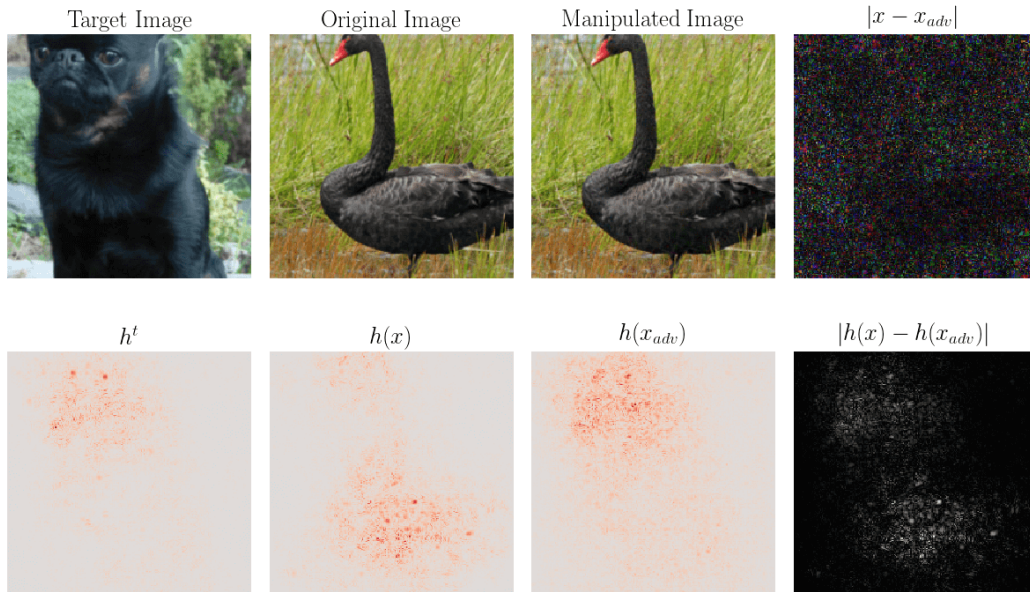


Figure 15: Gradient explanation maps produced with AlexNet model.

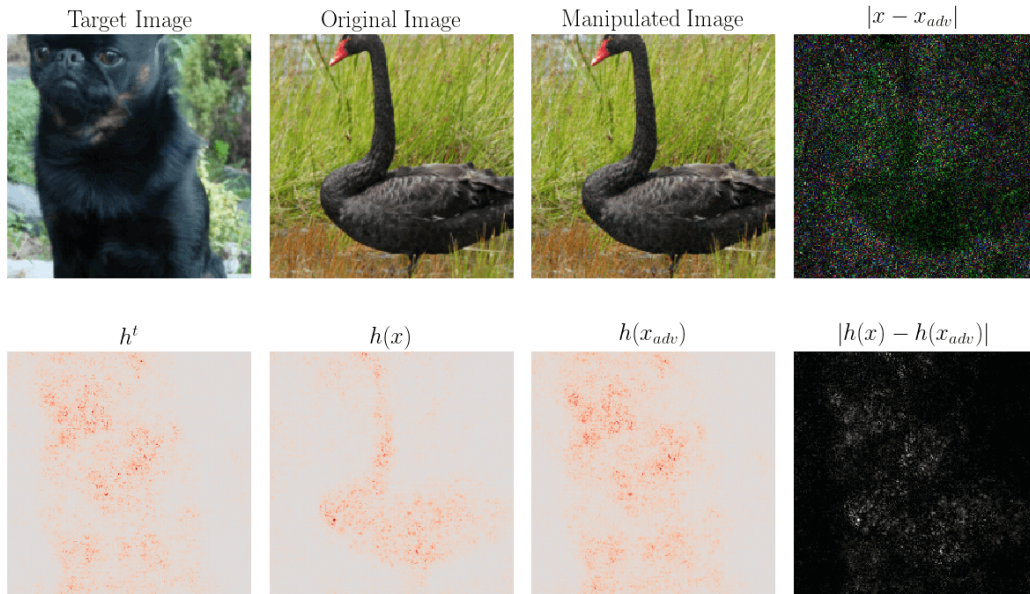


Figure 16: Gradient explanation maps produced with Densenet-121 model.

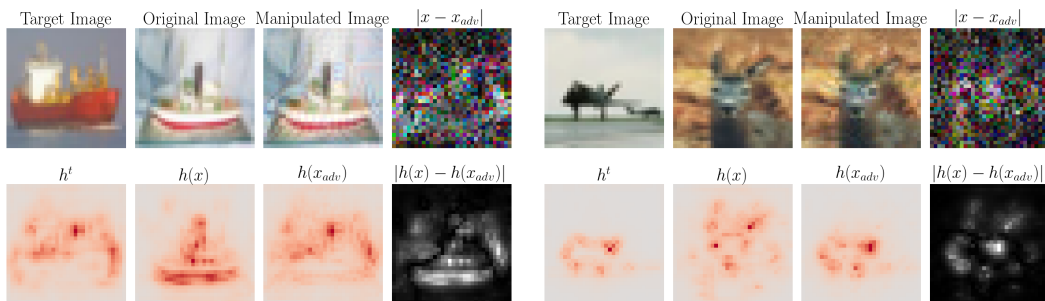


Figure 17: LRP Method on CIFAR-10 dataset

D Smoothing explanation methods

One can achieve a smoothing effect when substituting the relu activations for softplus $_{\beta}$ activations and then applying the usual rules for the different explanation methods.

A smoothing effect can also be achieved by applying the smoothgrad explanation method, see Figure 18. That is adding random perturbation to the image and then averaging over the resulting explanation maps. We average over 10 perturbed images with different values for the standard deviation σ of the Gaussian noise. The noise level n is related to σ as $\sigma = n \cdot (x_{max} - x_{min})$, where x_{max} and x_{min} are the maximum and minimum values the input image can have.

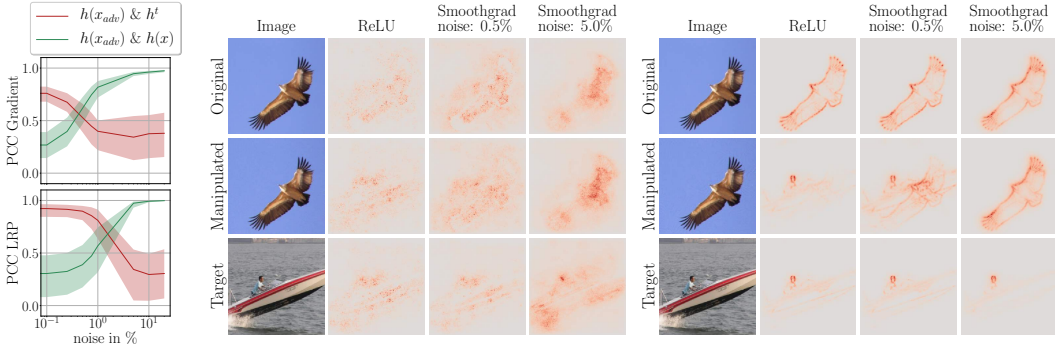


Figure 18: Recovering the original explanation map with SmoothGrad. Left: β dependence for the correlations of the manipulated explanation (here Gradient and LRP) with the target and original explanation. Line denotes median, 10^{th} and 90^{th} percentile are shown in semitransparent colour. Center and Right: network input and the respective explanation maps as β is lowered for Gradient (center) and LRP (right).

The β -smoothing or SmoothGrad explanation maps are more robust with respect to manipulations. Figure 19, 20 and 21 show results (MSE, SSIM and PCC) for 100 targeted attacks on the original explanation, the SmoothGrad explanation and the β -smoothed explanation for explanation methods Gradient and LRP.

For manipulation of SmoothGrad we use beta growth with $\beta_0 = 10$ and $\beta_e = 100$. For manipulation of β -Smoothing we set $\beta = 0.8$ for all runs. The hyperparameters for SmoothGrad and β -Smoothing are summarized in Table 4 and Table 5.

method	iterations	lr	factors
Gradient	1500	3×10^{-3}	$10^{11}, 10^6$
LRP	1500	3×10^{-4}	$10^{11}, 10^6$

Table 4: Hyperparameters used in our analysis for SmoothGrad.

method	iterations	lr	factors
Gradient	500	2.5×10^{-4}	$10^{11}, 10^6$
Grad x Input	500	2.5×10^{-4}	$10^{11}, 10^6$
IntGrad	200	2.5×10^{-3}	$10^{11}, 10^6$
LRP	1500	2.0×10^{-4}	$10^{11}, 10^6$
GBP	500	5.0×10^{-4}	$10^{11}, 10^6$
PA	500	5.0×10^{-4}	$10^{11}, 10^6$

Table 5: Hyperparameters used in our analysis for β -smoothing.

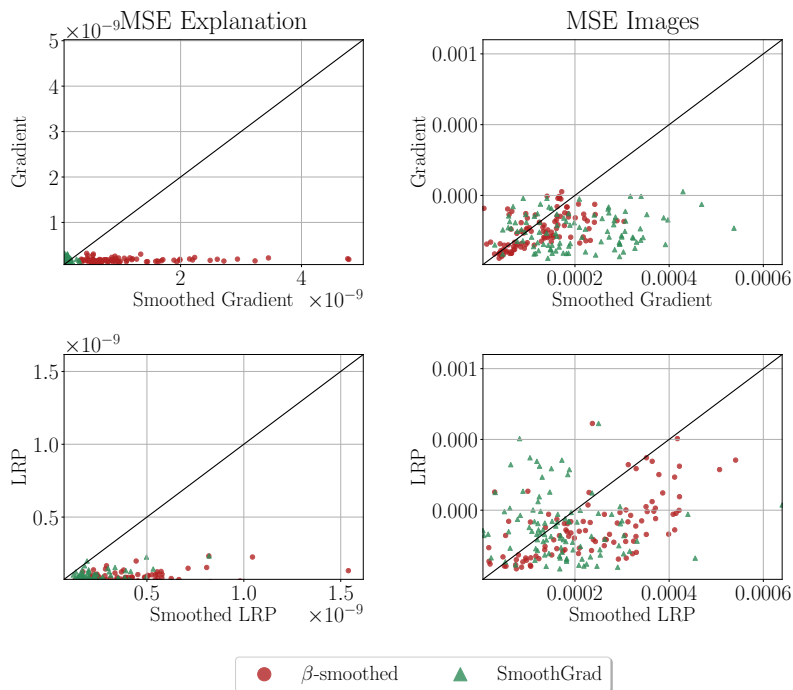


Figure 19: Left: Similarities between explanations. Markers are mostly right of the diagonal, i.e. the MSE for the smoothed explanations is higher than for the unsmoothed explanations which means the manipulated smoothed explanation map does not closely resemble the target h^t . Right: Similarities between Images. The MSE for the smoothed methods is higher (right of the diagonal) or comparable (on the diagonal), i.e. bigger or comparable perturbations in the manipulated Images when using smoothed explanation methods.

In Figure 22 and Figure 23, we directly compare the original explanation methods with the β -smoothed explanation methods. An increase in robustness can be seen for all methods: explanation maps for β -smoothed explanations have higher MSE and lower SSIM and PCC than explanation maps for the original methods. The similarity measures for the manipulated images are of comparable magnitude.

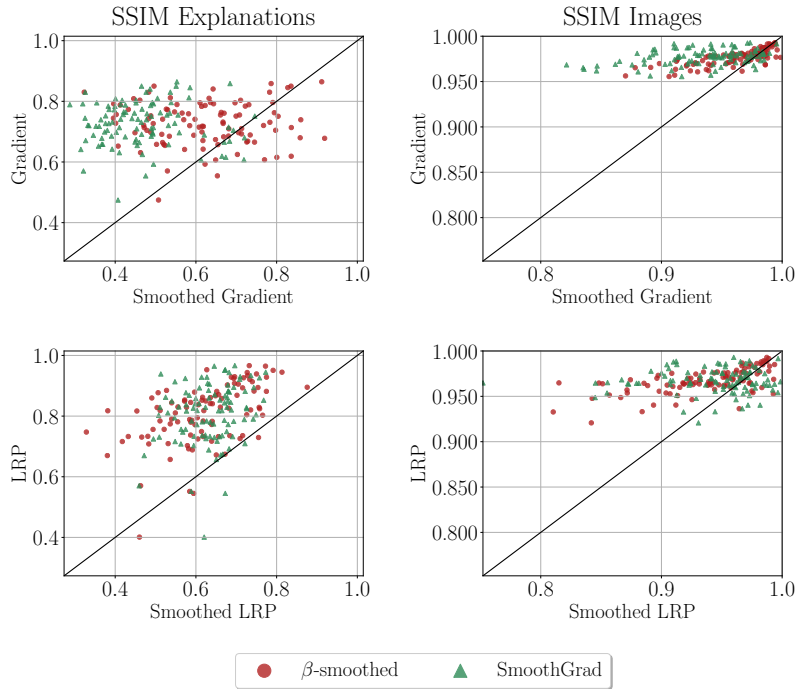


Figure 20: Left: Similarities between explanations. Markers are mostly left of the diagonal, i.e. the SSIM for the smoothed explanations is lower than for the unsmoothed explanations which means the manipulated smoothed explanation map does not closely resemble the target h^t . Right: Similarities between Images. The SSIM for the smoothed methods is lower (left of the diagonal) or comparable (on the diagonal), i.e. bigger or comparable perturbations in the manipulated Images when using smoothed explanation methods.

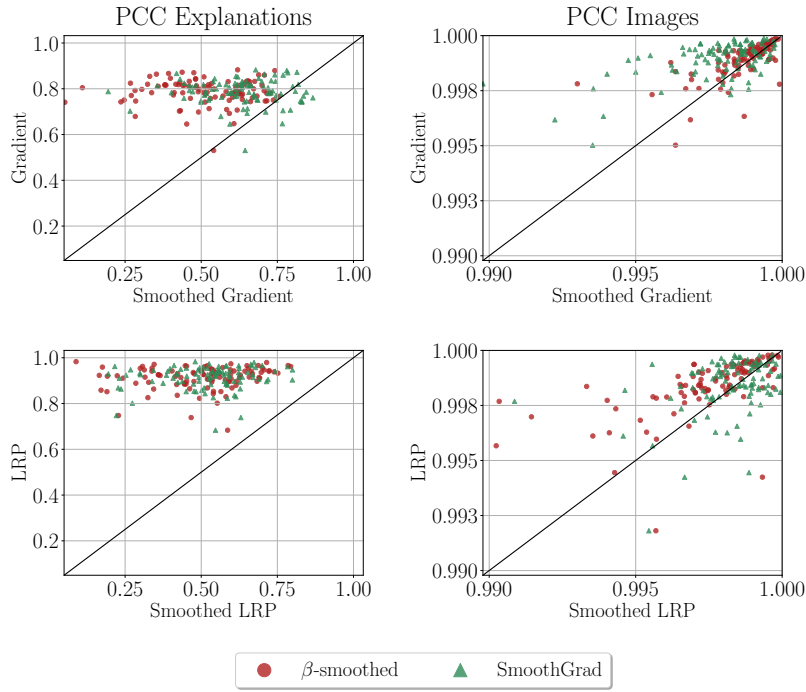


Figure 21: Left: Similarities between explanations. Markers are mostly left of the diagonal, i.e. the PCC for the smoothed explanations is lower than for the unsmoothed explanations which means that the manipulated smoothed explanation map does not closely resemble the target h^t . Right: Similarities between Images. The PCC for the smoothed methods is lower (left of the diagonal) or comparable (on the diagonal), i.e. bigger or comparable perturbations in the manipulated Images when using smoothed explanation methods.

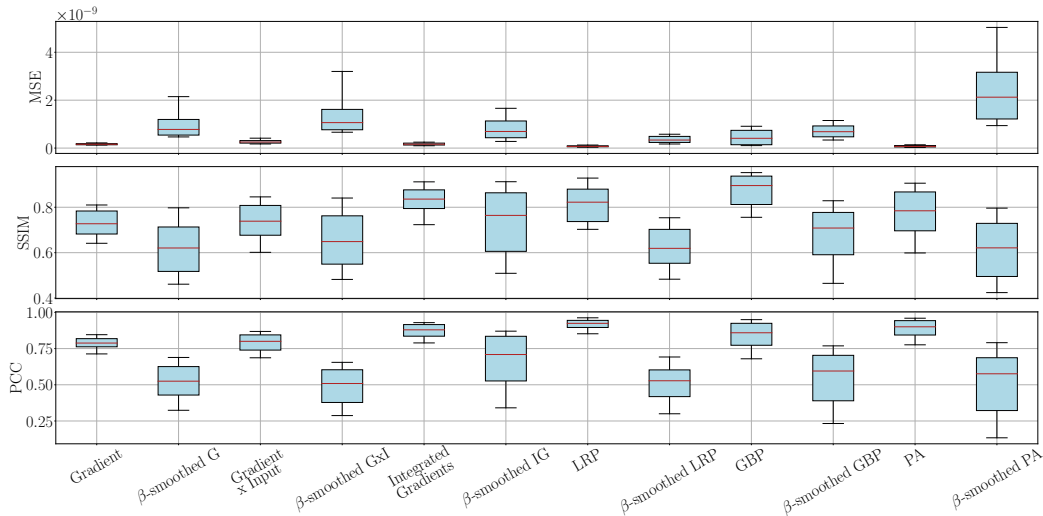


Figure 22: Comparison of Similarities of Explanation Maps for the original Explanation Methods and the β -smoothed Explanation Methods. Targeted attacks do not work very well on β -smoothed explanations, i.e. MSE is higher and SSIM and PCC are lower for the β -smoothed explanations than for the original explanations.

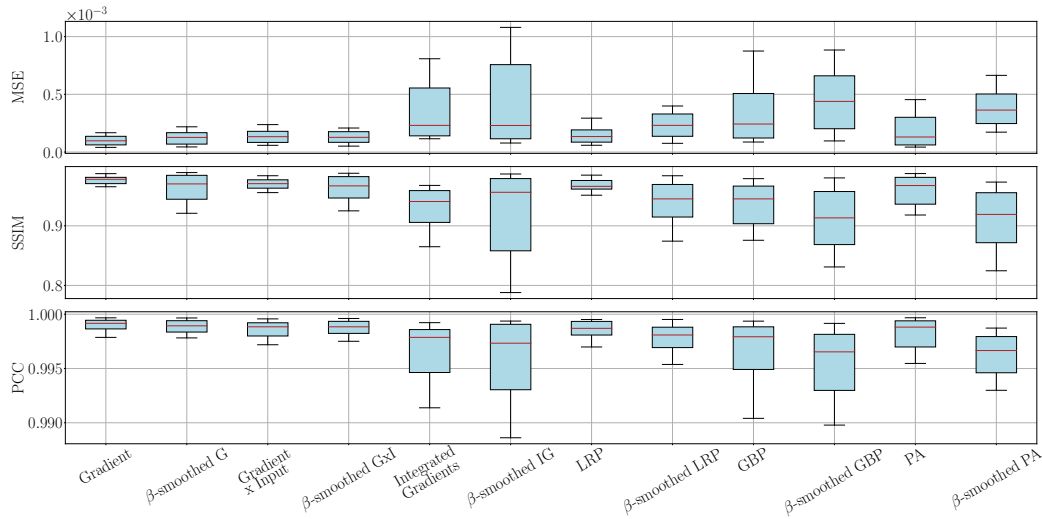


Figure 23: Comparison of Similarities between original and manipulated images. The similarity measures for images for the β -smoothed explanation methods are of comparable size or slightly worse (higher MSE, lower SSIM and lower PCC) than for the original explanation method, i.e. the manipulations are more visible for the β -smoothed explanation methods.

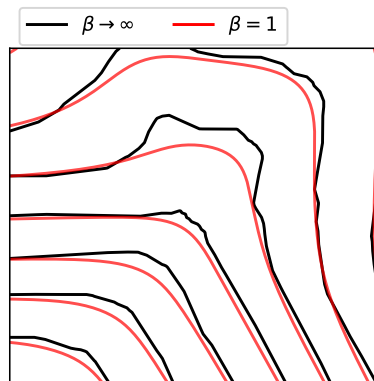


Figure 24: Contour plot of a 2-Layer Neural Network $f(x) = V^\top \text{softplus}(W^\top x)$ with $x \in [-1, 1]^2$, $W \in \mathbb{R}^{2 \times 50}$, $V \in \mathbb{R}^{50}$ and $V_i, W_{ij} \sim \mathcal{U}(-1, 1)$. Using a softplus activation with $\beta = 1$ visibly reduces curvature compared to a ReLU activation with $\beta \rightarrow \infty$.

E Proofs

In this section, we collect the proofs of the theorems stated in the main text.

E.1 Theorem 1

Theorem 3 *Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a network with softplus_β non-linearities and $\mathcal{U}_\epsilon(p) = \{x \in \mathbb{R}^d; \|x - p\| < \epsilon\}$ an environment of a point $p \in S$ such that $\mathcal{U}_\epsilon(p) \cap S$ is fully connected. Let f have bounded derivatives $\|\nabla f(x)\| \geq c$ for all $x \in \mathcal{U}_\epsilon(p) \cap S$. It then follows for all $p_0 \in \mathcal{U}_\epsilon(p) \cap S$ that*

$$\|h(p) - h(p_0)\| \leq |\lambda_{max}| d_g(p, p_0) \leq \beta C d_g(p, p_0), \quad (11)$$

where λ_{max} is the principle curvatures with the largest absolute value for any point in $\mathcal{U}_\epsilon(p) \cap S$ and the constant $C > 0$ depends on the weights of the neural network.

Proof: This proof will proceed in four steps. We will first bound the Frobenius norm of the Hessian of the network f . From this, we will deduce an upper bound on the Frobenius norm of the second fundamental form. This in turn will allow us to bound the largest principle curvature $|\lambda_{max}| = \max\{|\lambda_1| \dots |\lambda_{d-1}|\}$. Finally, we will bound the maximal and minimal change in explanation.

Step 1: Let $\text{softplus}^{(l)}(x) = \text{softplus}(W^{(l)}x)$ where $W^{(l)}$ are the weights of layer l .² We note that

$$\partial_k \text{softplus}\left(\sum_j W_{ij}x_j\right) = W_{ik} \sigma\left(\sum_j W_{ij}x_j\right) \quad (12)$$

$$\partial_l \sigma\left(\sum_j W_{ij}x_j\right) = \beta W_{il} g\left(\sum_j W_{ij}x_j\right) \quad (13)$$

where

$$\sigma(x) = \frac{1}{(1 + e^{-\beta x})}, \quad g(x) = \frac{1}{(e^{\beta x/2} + e^{-\beta x/2})^2}. \quad (14)$$

The activation at layer L is then given by

$$a^{(L)}(x) = (\text{softplus}^{(L)} \circ \dots \circ \text{softplus}^{(1)})(x) \quad (15)$$

Its derivative is given by

$$\begin{aligned} \partial_k a_i^{(L)} &= \sum_{s_2 \dots s_L} W_{is_L}^{(L)} \sigma\left(\sum_j W_{ij}^{(L)} a_j^{(L-1)}\right) W_{s_L s_{L-1}}^{(L-1)} \sigma\left(\sum_j W_{s_L j}^{(L-1)} a_j^{(L-2)}\right) \\ &\quad \dots W_{s_2 k}^{(1)} \sigma\left(\sum_j W_{s_2 j}^{(1)} x_j\right) \end{aligned}$$

We therefore obtain

$$\|\nabla a^{(L)}\| \leq \prod_{l=1}^L \|W^{(l)}\|_F \quad (16)$$

²We do not make the dependence of softplus on its β parameter explicit to ease notation.

Deriving the expression for $\partial_k a_i^{(L)}$ again, we obtain

$$\begin{aligned} \partial_l \partial_k a_i^{(L)} &= \sum_m \sum_{s_2 \dots s_L} \{ \\ &W_{i s_L}^{(L)} \sigma \left(\sum_j W_{ij}^{(L)} a_j^{(L-1)} \right) W_{s_L s_{L-1}}^{(L-1)} \sigma \left(\sum_j W_{ij}^{(L-1)} a_j^{(L-2)} \right) \\ &\dots \beta \sum_{\hat{s}_m} W_{s_{m+1} \hat{s}_m}^{(m)} W_{s_{m+1} s_m}^{(m)} g \left(\sum_j W_{s_{m+1} j}^{(m)} a_j^{(m-1)}(x) \right) \partial_l a_{\hat{s}_m}^{(m-1)}(x) \\ &\dots W_{s_2 k}^{(1)} \sigma \left(\sum_j W_{s_2 j}^{(1)} x_j \right) \} \end{aligned}$$

We now restrict to the case for which the index i only takes a single value and use $|\sigma(\bullet)| \leq 1$. The Hessian $H_{ij} = \partial_i \partial_j a^L(x)$ is then bounded by

$$\|H\|_F \leq \beta \tilde{C} \quad (17)$$

where the constant is given by

$$\tilde{C} = \sum_m \left\| W^{(L)} \right\|_F \left\| W^{(L-1)} \right\|_F \dots \left\| W^{(m)} \right\|_F^2 \dots \left\| W^{(1)} \right\|_F^2. \quad (18)$$

Step 2: Let $e_1 \dots e_{d-1}$ be a basis of the tangent space $T_p S$. Then the second fundamental form for the hypersurface $f(x) = c$ at point p is given by

$$\mathcal{L}(e_i, e_j) = -\langle D_{e_i} n(p), e_j \rangle \quad (19)$$

$$= -\langle D_{e_i} \frac{\nabla f(p)}{\|\nabla f(p)\|}, e_j \rangle = -\frac{1}{\|\nabla f(p)\|} \langle H[f] e_i, e_j \rangle + (\dots) \langle \nabla f(p), e_j \rangle \quad (20)$$

We now use the fact that $\langle \nabla f(p), e_j \rangle = 0$, i.e. the gradient of f is normal to the tangent space. This property was explained in the main text. This allows us to deduce that

$$\mathcal{L}(e_i, e_j) = -\frac{1}{\|\nabla f(p)\|} H[f]_{ij}. \quad (21)$$

Step 3: The Frobenius norm of the second fundamental form (considered as a matrix in the sense of step 2) can be written as

$$\|\mathcal{L}\|_F = \sqrt{\lambda_1^2 + \dots + \lambda_{d-1}^2}, \quad (22)$$

where λ_i are the principle curvatures. This property follows from the fact that the second fundamental form is symmetric and can therefore be diagonalized with real eigenvectors, e.g. the principle curvatures. Using the fact that the derivative of the network is bounded from below, $\|\nabla f(p)\| \geq c$, we obtain

$$|\lambda_{max}| \leq \beta \frac{\tilde{C}}{c}. \quad (23)$$

Step 4: For $p, p_0 \in \mathcal{U}_\epsilon(p) \cap S$, we choose a curve γ with $\gamma(t_0) = p_0$ and $\gamma(t) = p$. Furthermore, we use the notation $u(t) = \dot{\gamma}(t)$. It then follows that

$$n(p) - n(p_0) = \int_{t_0}^t \frac{d}{dt} (n(\gamma(t))) dt = \int_{t_0}^t D_{u(t)} n(\gamma(t)) dt \quad (24)$$

Using the fact that $D_{u(t)}n(\gamma(t)) \in T_{\gamma(t)}S$ and choosing an orthonormal basis $e_i(t)$ for the tangent spaces, we obtain

$$\int_{t_0}^t D_{u(t)}n(\gamma(t)) dt = \int_{t_0}^t \sum_j \langle e_j(t), D_{u(t)}n(\gamma(t)) \rangle e_j(t) dt \quad (25)$$

$$= \int_{t_0}^t \sum_j \mathcal{L}(e_j(t), u(t)) e_j(t) dt \quad (26)$$

The second fundamental form \mathcal{L} is bilinear and therefore

$$\int_{t_0}^t \sum_i \mathcal{L}(e_j(t), u(t)) e_j(t) dt = \int_{t_0}^t \sum_{i,j} \mathcal{L}(e_j(t), e_i(t)) u^i(t) e_j(t) dt \quad (27)$$

We now use the notation $\mathcal{L}_{ij}(t) = \mathcal{L}(e_j(t), e_i(t))$ and choose its eigenbasis for $e_i(t)$. We then obtain for the difference in the unit normals:

$$n(p) - n(p_0) = \int_{t_0}^t \sum_i \lambda_i(t) u^i(t) e_i(t) dt, \quad (28)$$

where $\lambda_i(t)$ denote the principle curvatures at $\gamma(t)$. By orthonormality of the eigenbasis, it can be easily checked that

$$\begin{aligned} \left\langle \sum_i \lambda_i(t) u^i(t) e_i(t), \sum_j \lambda_j(t) u^j(t) e_j(t) \right\rangle &\leq |\lambda_{max}|^2 \sum_i u^i(t)^2 \\ \Rightarrow \left\| \sum_i \lambda_i(t) u^i(t) e_i(t) \right\| &\leq |\lambda_{max}| \|u(t)\| \end{aligned}$$

Using this relation and the triangle inequality, we then obtain by taking the norm on both sides of (28):

$$\|n(p) - n(p_0)\| \leq |\lambda_{max}| \int_{t_0}^t \|\dot{\gamma}(t)\| dt. \quad (29)$$

This inequality holds for any curve connecting p and p_0 but the tightest bound follows by choosing γ to be the shortest possible path in $\mathcal{U}_\epsilon(p) \cap S$ with length $\int_{t_0}^t \|\dot{\gamma}(t)\| dt$, i.e. the geodesic distance $d_g(p, p_0)$ on $\mathcal{U}_\epsilon(p) \cap S$. The second inequality of the theorem is obtained by the upper bound on the largest principle curvature λ_{max} derived above, i.e. (23).

E.2 Theorem 2

Theorem 4 For one layer neural networks $g(x) = \text{relu}(w^T x)$ and $g_\beta(x) = \text{softplus}_\beta(w^T x)$, it holds that

$$\mathbb{E}_{\epsilon \sim p_\beta} [\nabla g(x - \epsilon)] = \nabla g_{\frac{\beta}{\|w\|}}(x), \quad (30)$$

where $p_\beta(\epsilon) = \frac{\beta}{(e^{\beta\epsilon/2} + e^{-\beta\epsilon/2})^2}$.

Proof: We first show that

$$\text{softplus}_\beta(x) = \mathbb{E}_{\epsilon \sim p_\beta} [\text{relu}(x)] , \quad (31)$$

for a scalar input x . This follows by defining $p(\epsilon)$ implicitly as

$$\text{softplus}_\beta(x) = \int_{-\infty}^{+\infty} p(\epsilon) \text{relu}(x - \epsilon) d\epsilon. \quad (32)$$

Differentiating both sides of this equation with respect to x results in

$$\sigma_\beta(x) = \int_{-\infty}^{+\infty} p(\epsilon) \Theta(x - \epsilon) \, d\epsilon = \int_{-\infty}^x p(\epsilon) \, d\epsilon, \quad (33)$$

where $\Theta(x) = \mathbb{I}(x > 0)$ is the Heaviside step function and $\sigma_\beta(x) = \frac{1}{(1+e^{-\beta x})}$. Differentiating both sides with respect to x again results in

$$p_\beta(x) = p(x). \quad (34)$$

Therefore, (31) holds. For a vector input \vec{x} , we define the distribution of its perturbation $\vec{\epsilon}$ by

$$p_\beta(\vec{\epsilon}) = \prod_i p_\beta(\epsilon_i), \quad (35)$$

where ϵ_i denotes the components of $\vec{\epsilon}$. We will suppress any arrows denoting vector-valued variables in the following in order to ease notation. We choose an orthogonal basis such that

$$\epsilon = \epsilon_p \hat{w} + \sum_i \epsilon_o^{(i)} \hat{w}_o^{(i)} \quad \text{with} \quad \hat{w} \cdot \hat{w}_o^{(i)} = 0 \quad \text{and} \quad w = \|w\| \hat{w}. \quad (36)$$

This allows us to rewrite

$$\begin{aligned} \mathbb{E}_{\epsilon \sim p_\beta} [\text{relu}(w^T(x - \epsilon))] &= \mathbb{E}_{\epsilon \sim p_\beta} [\text{relu}(w^T x - \|w\| \epsilon_p)] \\ &= \int p_\beta(\epsilon_p) (\text{relu}(w^T x - \|w\| \epsilon_p)) \, d\epsilon_p \end{aligned}$$

By changing the integration variable to $\epsilon = \|w\| \epsilon_p$ and using (31), we obtain

$$\text{softplus}_{\frac{\beta}{\|w\|}}(w^T x) = \mathbb{E}_{\epsilon \sim p_\beta} [\text{relu}(w^T(x - \epsilon))] , \quad (37)$$

The theorem then follows by deriving both sides of the equation with respect to x .

F Additional examples for VGG

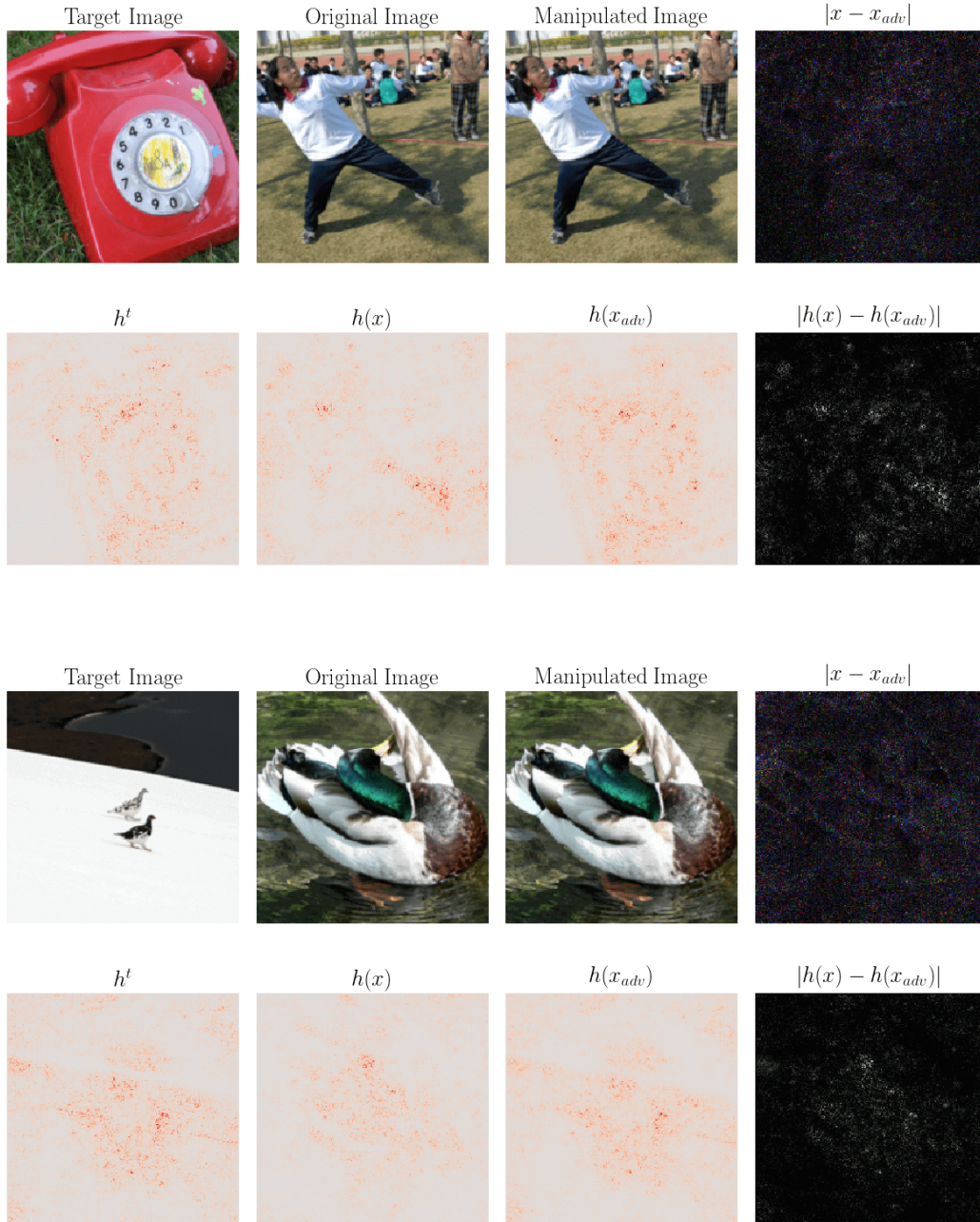


Figure 25: Explanation map produced with the Gradient Explanation Method on VGG.

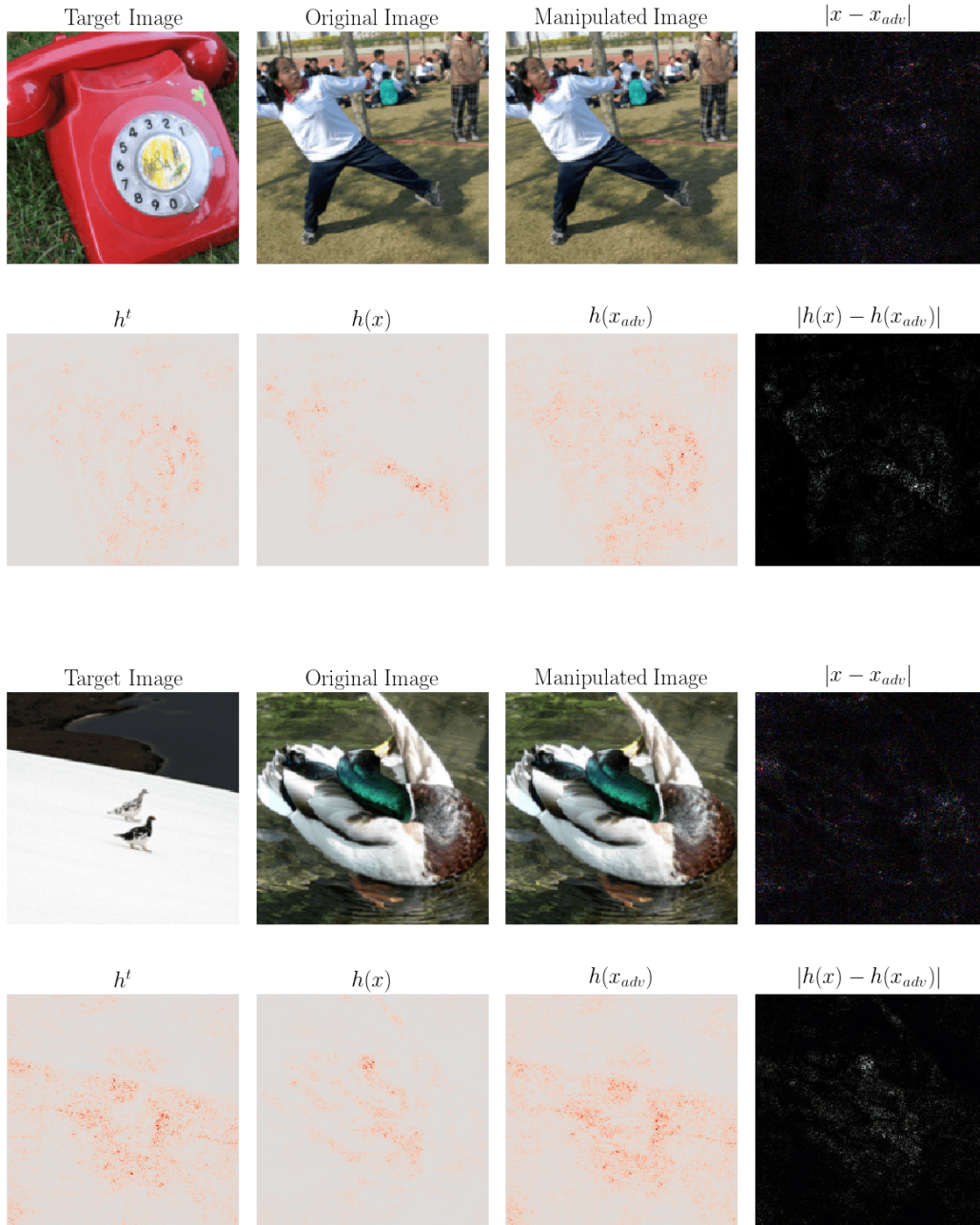


Figure 26: Explanation map produced with the Gradient \times Input Explanation Method on VGG.

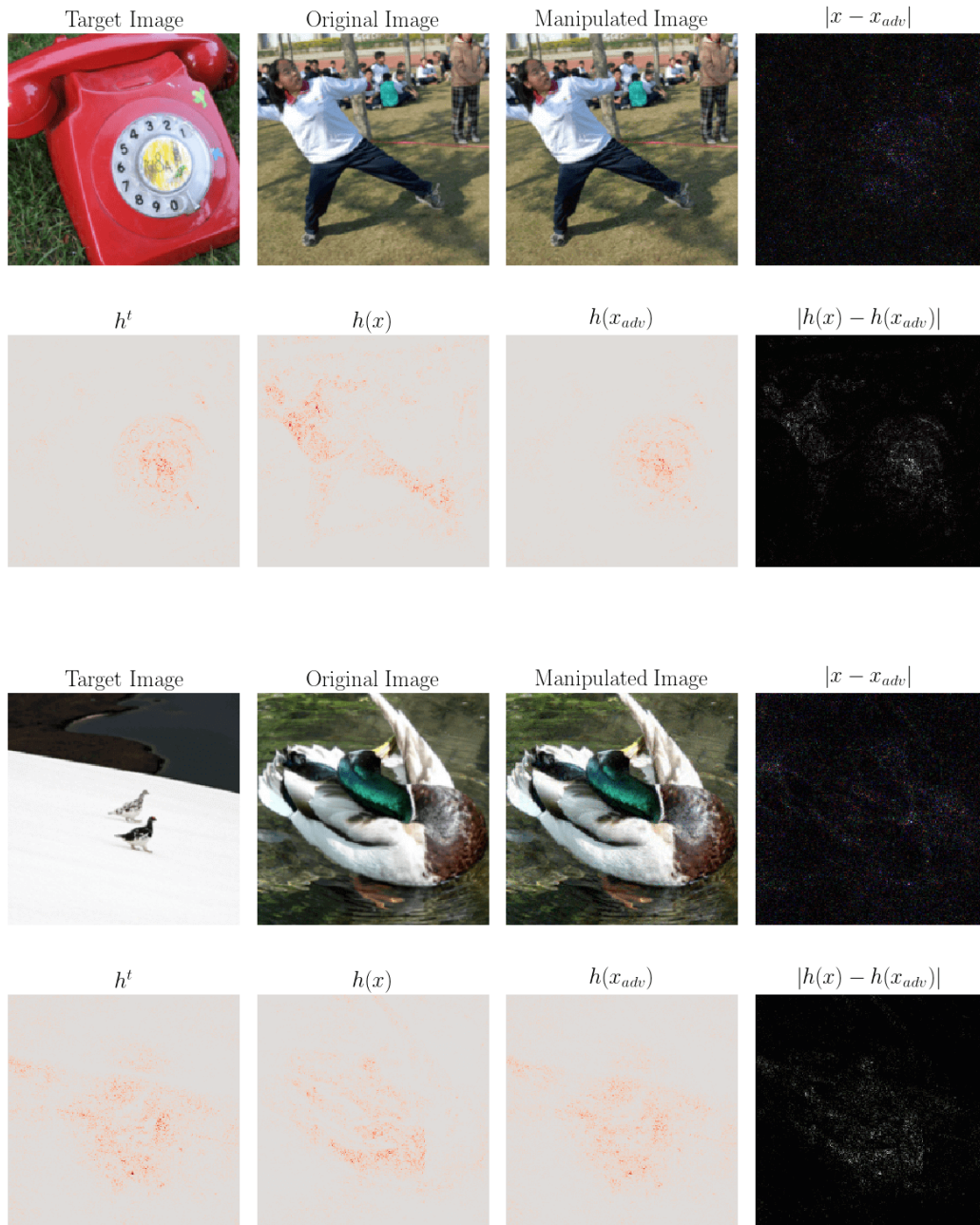


Figure 27: Explanation map produced with the Integrated Gradients Explanation Method on VGG.

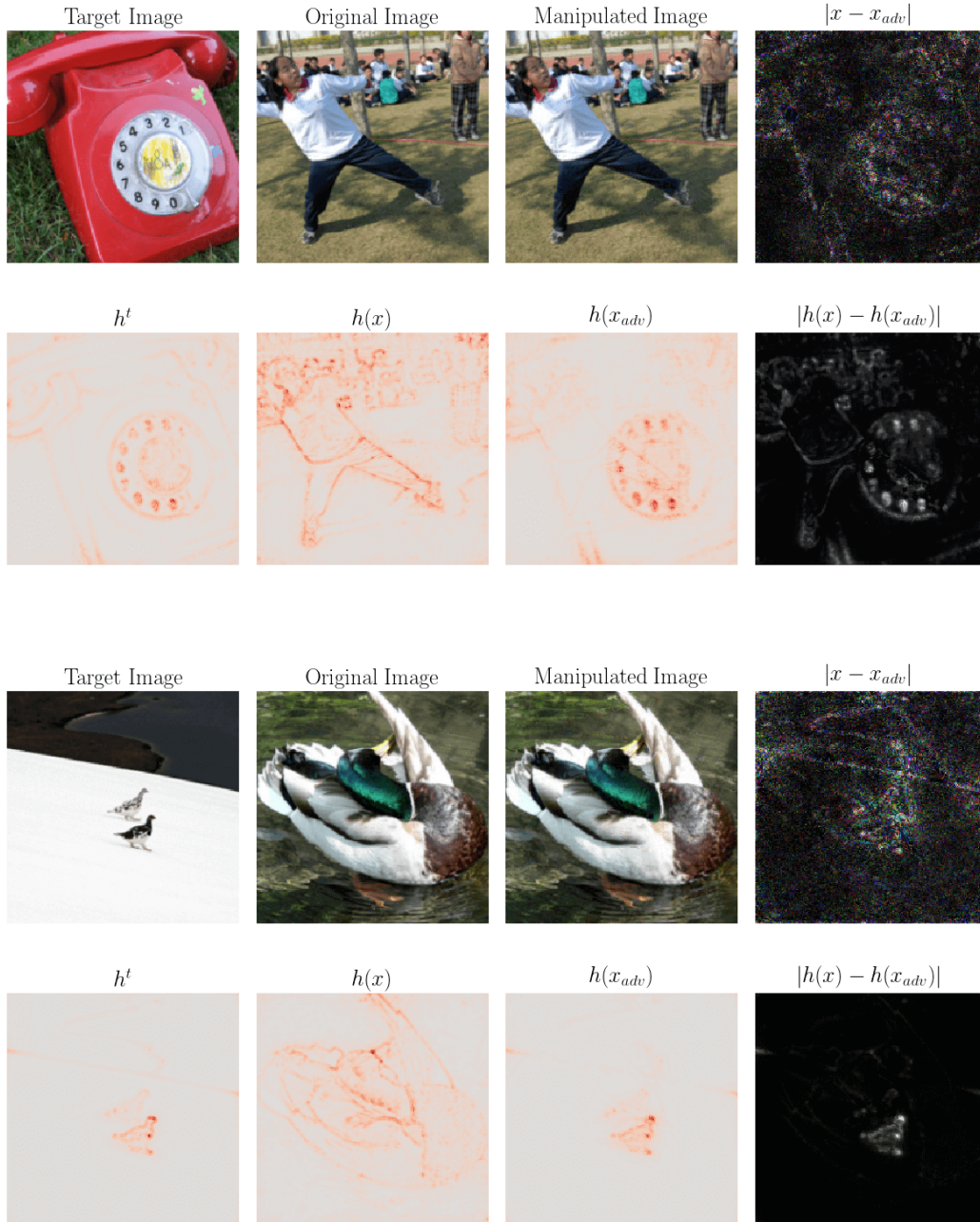


Figure 28: Explanation map produced with the LRP Explanation Method on VGG.

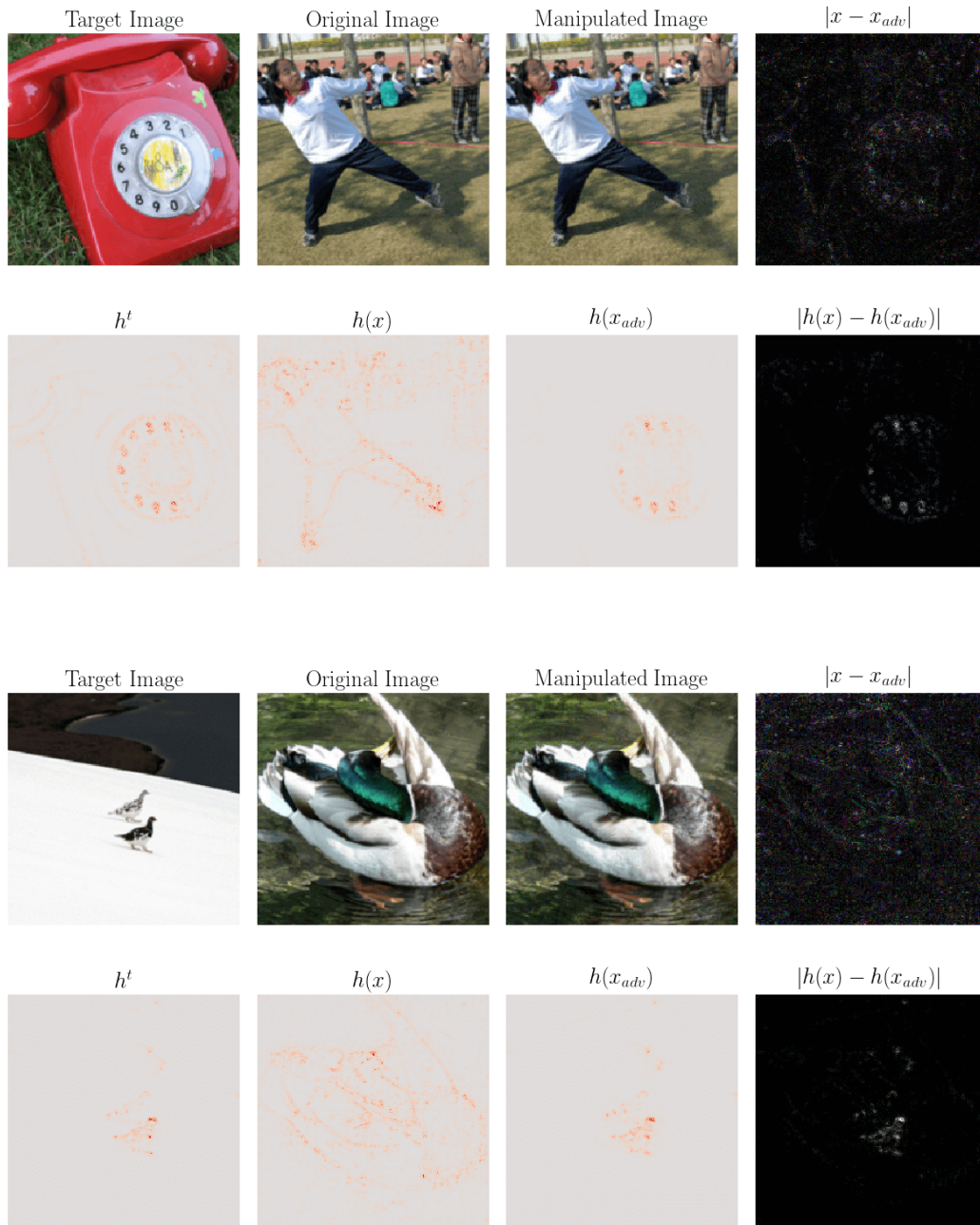


Figure 29: Explanation map produced with the Guided Backpropagation Explanation Method on VGG.

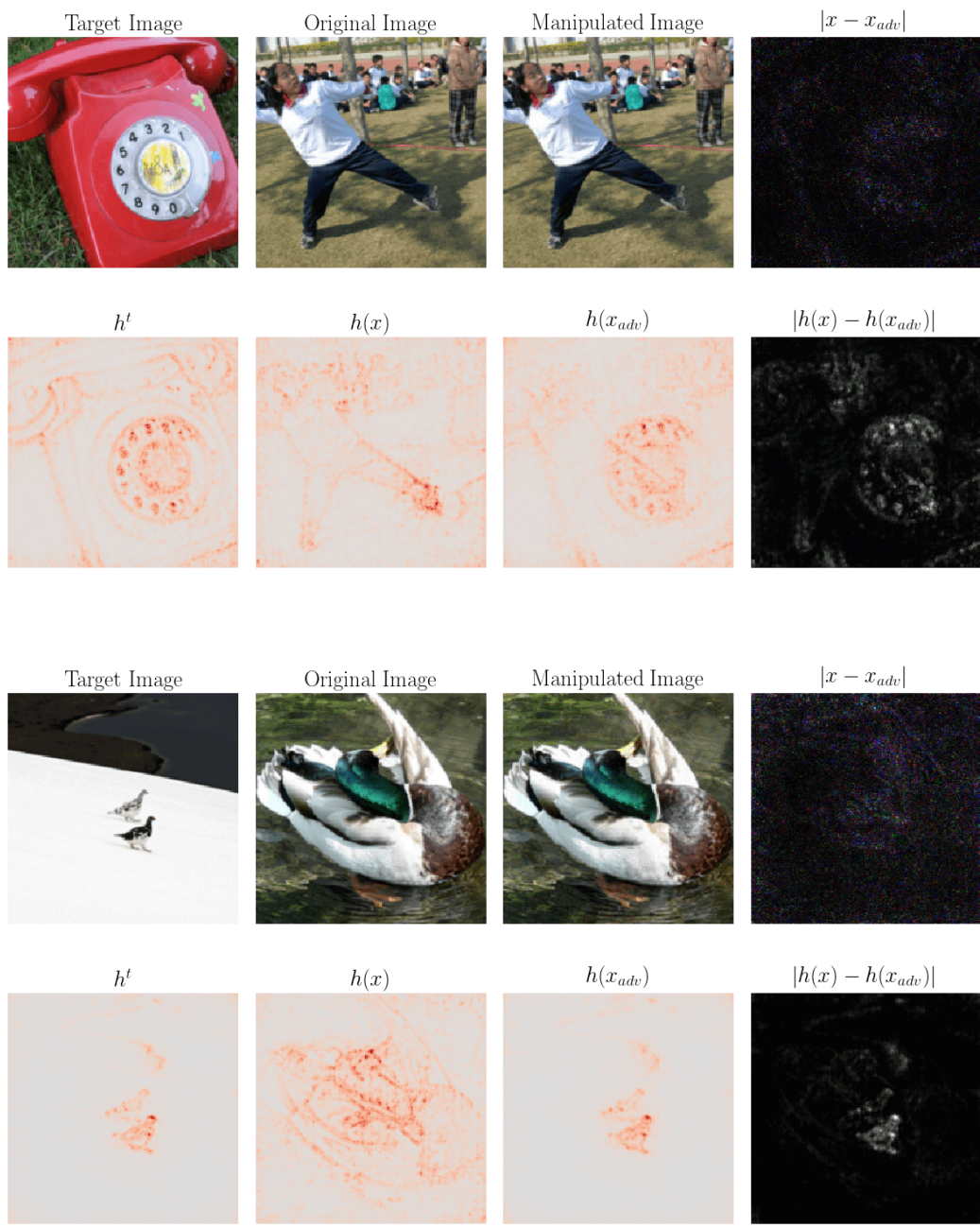


Figure 30: Explanation map produced with the Pattern Attribution Explanation Method on VGG.