

# Interpreting Undesirable Pixels for Image Classification on Black-Box Models

Sin-Han Kang<sup>1</sup>, Hong-Gyu Jung<sup>1</sup> and Seong-Whan Lee<sup>2,1</sup>

<sup>1</sup>Department of Brain and Cognitive Engineering, Korea University, Seoul, Korea

<sup>2</sup>Department of Artificial Intelligence, Korea University, Seoul, Korea

<sup>1</sup>{kangsinhan, hkjung00, sw.lee}@korea.ac.kr

## Abstract

In an effort to interpret black-box models, researches for developing explanation methods have proceeded in recent years. Most studies have tried to identify input pixels that are crucial to the prediction of a classifier. While this approach is meaningful to analyse the characteristic of black-box models, it is also important to investigate pixels that interfere with the prediction. To tackle this issue, in this paper, we propose an explanation method that visualizes undesirable regions to classify an image as a target class. To be specific, we divide the concept of undesirable regions into two terms: (1) factors for a target class, which hinder that black-box models identify intrinsic characteristics of a target class and (2) factors for non-target classes that are important regions for an image to be classified as other classes. We visualize such undesirable regions on heatmaps to qualitatively validate the proposed method. Furthermore, we present an evaluation metric to provide quantitative results on ImageNet.

## 1. Introduction

The tremendous growth of deep networks has brought about the solvability of key problems in computer vision such as object classification [15], [17] and object detection [11], [4]. At the same time, the complexity of models has also increased, making it difficult for humans to understand the decisions of the model. To improve interpretability of black-box models, explanation methods have been proposed in terms of model inspection [3], [16], [10] and outcome explanation [18], [9]. These studies focus on visualizing crucial pixels for a model prediction. In other

*This work was supported by Institute for Information & communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No.2017-0-01779, A machine learning and statistical inference framework for explainable artificial intelligence)*

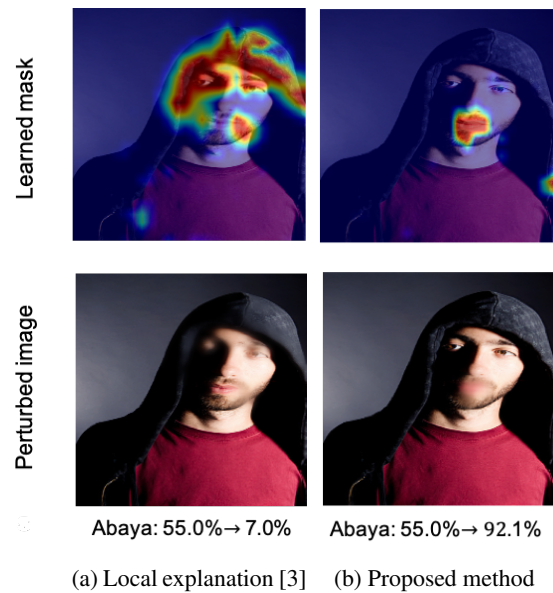


Figure 1: Comparison between a local explanation [3] and the proposed method. (a) [3] produces learned masks that are crucial to classify the abaya class. Thus, perturbing the pixels results in the accuracy significantly decreased. (b) the proposed method generates undesirable pixels for the abaya class. In this case, the perturbation image based on the learned mask improves the accuracy. The change on the accuracy before and after perturbation is presented on the bottom of each image.

words, if we remove those pixels, the prediction accuracy is significantly decreased.

However, to obtain diverse interpretation on black-box models, it is also important to investigate pixels that interfere with the prediction. Thus, in this paper, we aim to find undesirable pixels that can improve the accuracy of a target class by perturbing the pixels. For example, Fig. 1 shows

the difference between [3] and our proposed method. While Fig. 1(a) explains that a hoodie and eyes play a major role to classify the image as an abaya that is a full-length outer garment worn by Muslim women, our method finds regions that help to improve the accuracy. Specifically, Fig. 1(b) interprets a mustache as undesirable pixels, which is generally not seen by women. Thus, perturbing the mustache leads to the accuracy improved for the abaya class.

Figure 1 clearly shows that finding undesirable pixels of a target class can ease off the uncertainty about the decision of black-box models. Thus, we further define undesirable pixels with two different concepts. The first is factors for a target class (F-TC), which hinder that black-box models identify intrinsic characteristics of a target class. The second is factors for non-target classes (F-NTC) that are important regions for an image to be classified as other classes. In the following sections, we will mathematically elaborate on how these two different concepts interpret undesirable pixels for a model prediction. Then, we visually validate our idea on heatmaps and qualitatively evaluate the proposed method on ImageNet.

## 2. Related Works

Class activation map (CAM) [19] and Grad-CAM [14] analyze the decision of neural networks on heatmaps by utilizing activation maps of the last convolution layer in CNNs. Layer-wise relevance propagation (LRP) [1] computes gradients of the prediction score by exploiting a backward operation in neural networks. Model agnostic methods [12], [8] approximate the perimeter decision boundary of a black-box model to a simple model such as logistic regression and decision tree. Local rule-based explanations (LORE) [5] applies a genetic algorithm to build rule-based explanations, offering a set of counterfactual rules. Contrastive explanations methods (CEM) [2] visualize a pertinent positive (PP) and a pertinent negative (PN) by using perturbation. But, PN is useful only when the meaning of classes for different inputs are similar to each other. Lastly, the most similar work to ours is local explanation [3] that learns a mask by perturbing important regions to a prediction. However, these methods do not clearly consider undesirable pixels of an image for a target class.

## 3. Methods

Given an image  $X \in \mathbb{R}^{H \times W \times 3}$ , we generate a blurred image by applying Gaussian blur  $h(X) = g_{\sigma,s}(X)$  where  $\sigma$  and  $s$  are standard deviation and kernel size, respectively. In order to replace specific pixels in  $X$  with blurred pixels, we define a mask  $M \in [0, 1]^d$ , where  $d$  is smaller than  $H \times W$ . Thus, a perturbed image is generated by the masking operator [3] as follows.

$$Q(X; M', h) = X \circ M' + h(X) \circ (1 - M'), \quad (1)$$

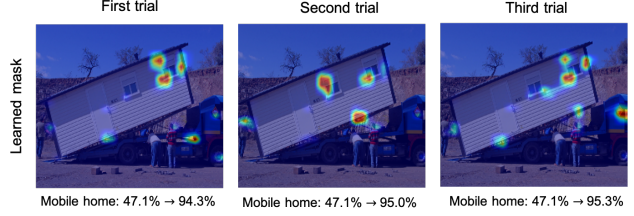


Figure 2: Learned masks generated by Eq. 3 using TV and  $l_1$  norms. When applying Eq. 3 several times, different learned masks are obtained. This leads to inconsistent interpretability.

where  $M' = \text{Inp}(M)$  is an interpolated mask and  $\text{Inp}$  is a bilinear interpolation function.  $\circ$  denotes the element-wise multiplication. Given a black-box model  $f$  and an accuracy  $f_k(X)$  for a target class  $k$ , we expect that the perturbed image makes  $f_k(X) \ll f_k(Q(X; M', h))$ . In other words, the goal is to find an optimal mask  $M^*$  that improves the accuracy for a target class and an objective function can be defined as follows.

$$M^* = \underset{M}{\operatorname{argmax}} f_k(Q(X; M', h)). \quad (2)$$

Since [3] shows that total-variation (TV) norm and  $l_1$  norm can produce reasonable and precise interpretability for the masking operation, we also apply such regularizers to our objective.

$$M^* = \underset{M}{\operatorname{argmax}} f_k(Q(X; M', h)) - \mathcal{R}_M, \quad (3)$$

where  $R_M = \lambda_1 \sum_{i,j} ((M_{i+1,j} - M_{i,j})^\beta + (M_{i,j+1} - M_{i,j})^\beta)^{\frac{1}{\beta}} + \lambda_2 \|1 - M\|_1$ .  $\lambda_1$ ,  $\lambda_2$  and  $\beta$  are hyper-parameters.

However, this objective function generates different masks for each trial as shown in Fig. 2 and do not provide consistency of an explanation. We conjecture that this is due to the softmax operation. Given an output before softmax  $y$ ,  $f_k(Q(X; M', h)) = \frac{\exp(y_k)}{\sum_i \exp(y_i)}$  can be higher when increasing  $\exp(y_k)$  or decreasing  $\sum_i \exp(y_i)$ . That is, improving  $f_k(Q(X; M', h))$  is affected by not only the output for a target class but also those for other classes.

In order to solve this problem, we propose two types of regularizers for obtaining undesirable pixels. We first define factors for a target class (F-TC).

$$\mathcal{R}_{F-TC} = \gamma \left\| \frac{1}{N-1} \sum_{i, i \neq k} \{f'_i(Q(X; M', h)) - f'_i(X)\} \right\|_2, \quad (4)$$

where  $f'_i(\cdot)$  denotes the output before softmax for the  $i$ -th class,  $k$  is the index of the target class,  $\gamma$  is a hyper-

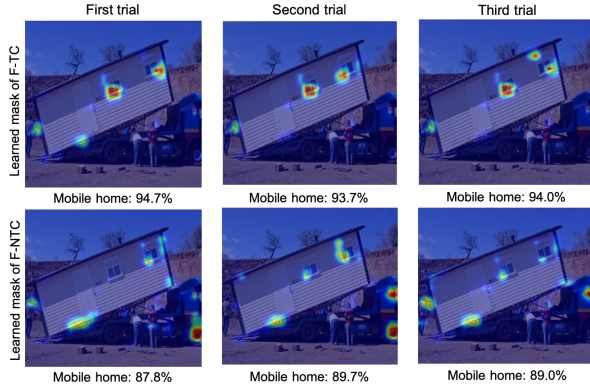


Figure 3: Comparison of F-TC and F-NTC. F-TC explains that undesirable pixels to identify the mobile home are the windows. On the other hands, F-NTC focuses on the parts of the truck. Both methods produce consistent results for each trial. The accuracy of the target class is presented on the bottom of each image.

parameter and  $N$  is the total number of classes. This regularizer forces the objective function into focusing on the target class itself. In other words,  $\mathcal{R}_{F-TC}$  finds the pixels that hinder intrinsic characteristic to be classified as the target class. The final objective function can be expressed as

$$M^* = \operatorname{argmax}_M f_k(Q(X; M', h)) - \mathcal{R}_M - \mathcal{R}_{F-TC}. \quad (5)$$

Secondly, we define factors for non-target classes (F-NTC).

$$\mathcal{R}_{F-NTC} = \gamma \|\{f'_k(Q(X; M', h)) - f'_k(X)\}\|_2, \quad (6)$$

which encourages to find undesirable pixels by focusing on other classes except for the target class. When applying  $\mathcal{R}_{F-NTC}$ , we modify Eq. 3 as follows.

$$M^* = \operatorname{argmin}_M \sum_{i, i \neq k} \{f_i(Q(X; M', h))\} + \mathcal{R}_M + \mathcal{R}_{F-NTC}. \quad (7)$$

In the following section, we show several case studies to understand how these regularizations behave according to their definitions.

## 4. Experiments

### 4.1. Experimental Settings

We use VGG-19 [15] and ResNet-18 [6] pretrained on the ImageNet [13] and solve optimization problems using Adam [7]. We set the learning rate to 0.1 and iterations to

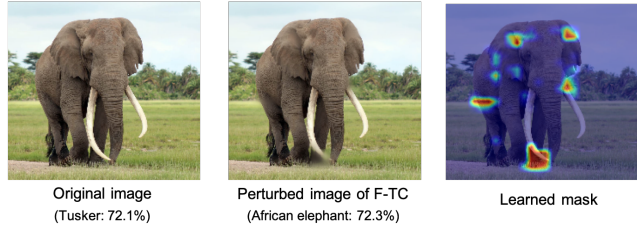


Figure 4: Visualizations of distinctive characteristics between similar classes. The length of the horn plays a major role to distinguish between an African elephant and a tusker. The highest accuracy for each class is presented by brackets.

200. We use the hyper-parameters  $\lambda_1 = 1.7$ ,  $\lambda_2 = 3.0$ ,  $\beta = 2$  and  $\gamma = 0.3$ . A mask  $28 \times 28$  is interpolated by  $224 \times 224$  size by upsampling. The standard deviation  $\sigma$  and kernel size  $s$  for the Gaussian kernel are set to 5 and 11, respectively.

### 4.2. Interpretability

In Sect. 3, we explained a main objective function with (1) TV norm and  $l_1$  norms. Further, additional regularizers such as (2) F-TC and (3) F-NTC were proposed. We now compare interpretability among the three cases. First, as shown in Fig. 2, when merely using TV and  $l_1$  norms, the learned masks are generated irregularly for each trial. This makes us difficult to understand the decision of black-box models. On the other hands, Fig. 3 shows that F-TC and F-NTC provide consistent visual interpretation. Moreover, each regularizer highlights the regions corresponding to their definitions such as Eq. 4 and Eq. 6. Specifically, F-TC explains that the windows are undesirable pixels to identify the intrinsic characteristic of the mobile home. F-NTC explains that the parts of the truck are undesirable pixels since those are more important to classify other classes such as a truck. In this way, our algorithm can be exploited to understand the decision of black-box models.

### 4.3. Qualitative Results

We provide more examples to qualitatively evaluate the proposed method. We used VGG-19 for all experiments of this section.

As illustrated in Fig. 4, the original image is classified as the tusker with the accuracy of 72.1%. When we set the African elephant as a target class, F-TC perturbs the end part of the horn, which results in improving the accuracy for the African elephant class. These results imply that the model generally regards the length of the horn as crucial features to distinguish between the tusker and the African elephant. More importantly, this is consistent with the fact

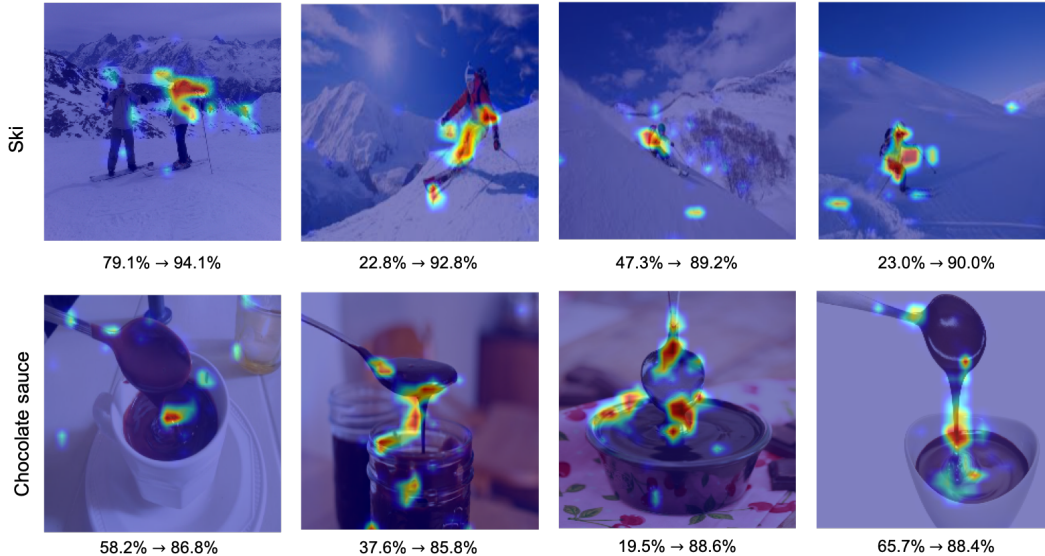


Figure 5: Behaviours of F-TC on VGG-19. The human body connected to ski equipments and the chocolate sauce falling into the cup are found as undesirable pixels to be classified as the target classes. The change on the accuracy before and after perturbation is presented on the bottom of each image.

that the horn of a tusker is longer than an African elephant. Thus, we argue that our method can provide reasonable interpretation about how trained networks distinguish similar classes.

Another example can be shown in Fig. 5. For ski and chocolate sauce classes, F-TC highlights the human body connected to ski equipments and the chocolate sauce that is falling from the spoon. These results suggest that portions connected to a target class have negative effect on a classification for a target class.

#### 4.4. Quantitative Results

We present the following evaluation metric to measure how effectively our method finds undesirable pixels.

$$\phi = \mathbb{E}_X \left[ \left( \frac{f_h(Q(X; M', h)) - f_h(X)}{1 - f_h(X)} \right) * 100 \right], \quad (8)$$

where  $h$  is a class that has the highest accuracy for an image.  $1 - f_h(X)$  is the residual accuracy that can be improved from  $f_h(X)$ . Thus, Eq. 8 measures the relative accuracy improvement. We randomly select 1,000 images from the ImageNet and compare results between F-TC and F-NTC with VGG-19 and ResNet-18. In Table 1, we observe that the accuracy can be effectively improved by perturbing undesirable pixels. We also measure the ratio of the number of undesirable pixels to the image size  $224 \times 224$ . In this case, we use the pixels that have magnitude above a threshold 0.6. Table 2 shows that both F-TC and F-NTC yield a small number of undesirable pixels that are below 4%.

Model	F-TC	F-NTC
VGG-19	48.741	48.979
ResNet-18	44.898	44.688

Table 1: Relative accuracy improvement. The results indicate that perturbing undesirable pixels can effectively improve the classification performance.

Model	F-TC	F-NTC
VGG-19	0.0373	0.0378
ResNet-18	0.0360	0.0398

Table 2: The percentage of undesirable pixels out of total image pixels. A small number of pixels are only used to find undesirable pixels.

## 5. Conclusion

We proposed an explanation method that visualizes undesirable regions for classification. We defined the undesirable regions by two terms. The first is factors for a target class, which hinder that black-box models identify intrinsic characteristics of a target class and the second is factors for non-target classes that are important regions for an image to be classified as other classes. We showed the proposed method successfully found reasonable regions according to their definitions and by perturbing the undesirable pixels, we could improve the accuracy for the target class.

## References

- [1] A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, and W. Samek. Layer-wise relevance propagation for neural networks with local renormalization layers. In *International Conference on Artificial Neural Networks*, pages 63–71. Springer, 2016.
- [2] A. Dhurandhar, P.-Y. Chen, R. Luss, C.-C. Tu, P. Ting, K. Shanmugam, and P. Das. Explanations based on the missing: Towards contrastive explanations with pertinent negatives. In *Advances in Neural Information Processing Systems*, pages 592–603, 2018.
- [3] R. C. Fong and A. Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3429–3437, 2017.
- [4] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [5] R. Guidotti, A. Monreale, S. Ruggieri, D. Pedreschi, F. Turini, and F. Giannotti. Local rule-based explanations of black box decision systems. *arXiv preprint arXiv:1805.10820*, 2018.
- [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [7] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [8] S. M. Lundberg and S.-I. Lee. A unified approach to interpreting model predictions. In *Advances in Neural Information Processing Systems*, pages 4765–4774, 2017.
- [9] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller. Explaining nonlinear classification decisions with deep taylor decomposition. *Pattern Recognition*, 65:211–222, 2017.
- [10] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems*, pages 3387–3395, 2016.
- [11] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [12] M. T. Ribeiro, S. Singh, and C. 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.
- [13] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [14] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017.
- [15] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [16] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3319–3328. JMLR.org, 2017.
- [17] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [18] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, pages 2048–2057, 2015.
- [19] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning deep features for scene recognition using places database. In *Advances in neural information processing systems*, pages 487–495, 2014.