# Vulnerability Analysis of Chest X-Ray Image Classification Against Adversarial Attacks

Saeid Asgari Taghanaki[1], Arkadeep Das[1,2], and Ghassan Hamarneh[1]

[1] School of Computing Science, Simon Fraser University, Canada
[2] Department of Mathematics, Indian Institute of Technology, Guwahati, India
{sasgarit,arkadeep das,hamarneh}@sfu.ca

**Abstract.** Recently, there have been several successful deep learning approaches for automatically classifying chest X-ray images into different disease categories. However, there is not yet a comprehensive vulnerability analysis of these models against the so-called adversarial perturbations/attacks, which makes deep models more trustful in clinical practices. In this paper, we extensively analyzed the performance of two state-of-the-art classification deep networks on chest X-ray images. These two networks were attacked by three different categories (ten methods in total) of adversarial methods (both white- and black-box), namely gradient-based, score-based, and decision-based attacks. Furthermore, we modified the pooling operations in the two classification networks to measure their sensitivities against different attacks, on the specific task of chest X-ray classification.

**Keywords:** Adversarial perturbation · chest X-ray classification · deep learning

## 1 Introduction

The chest X-ray is among the top most commonly accessible medical imaging examinations used for affordable screening and diagnosis of numerous lung ailments including pneumothorax, mass, cardiomegaly, effusion, and pneumonia. Owing to huge numbers of patients and increasing burden of lung ailments, the workload of radiologists has significantly multiplied. Hence, with an intention to accelerate/support the predictions of radiologists, many machine (deep) learning classification frameworks have emerged over the past few years.

The availability of a new large scale chest X-ray dataset namely "ChestX-ray14" [20], which comprises 30,805 patients and 112,120 chest X-ray images, makes it feasible to apply deep learning without a need for data augmentation or synthetic data. Recently, different standard classification deep networks (AlexNet [7], VGGNet [14] and ResNet [5]) have been applied to this dataset. Wang et al. [20] applied pre-trained AlexNet, GoogLeNet [17], VGG, and ResNet-50 architectures to classify 8 disease categories. They showed that ResNet-50 achieves superior performance compared to the other applied models. Guendel et al. [4] proposed a local aware dense network for classification of 14 pathology

classes in the ChestX-ray14 dataset. Rajpurkar et al. [12] proposed CheXNet, a 121-layer convolutional neural network trained on ChestX-ray14 for the pneumonia disease detection task, which exceeds average radiologist performance on the F1 metric. Baltruschat et al. [1] proposed a fine-tuned ResNet-50 network which achieved high accuracy on 4 out of the 14 disease classes in the chest X-ray dataset. Yao et al. [21] presented a partial solution to constraints in using LSTMs to leverage inter-dependencies among target labels in predicting 14 pathological classes from chest X-rays.

The *generalizability* of the deep learning methods, i.e. how they perform on unseen chest X-ray test images, have been explored in the above mentioned works to some extent. However, discovery of "adversarial examples" has exposed serious vulnerabilities in even state-of-the-art deep learning systems [8]. As of writing there is no comprehensive study on the *vulnerability* analysis of the state-of-the art classification networks against adversarial perturbations for chest X-rays. Samuel G. et al. [3] considered a single attack, namely projected gradient descent [9,6] on chest X-ray images.

Adversarial images are crafted by adding perturbations, imperceptible to the naked eye, to the clean images to fool machine learning models. Different categories [22] of adversarial attacks on images have been recently developed which have been highly successful in fooling deep neural networks. In the medical image analysis domain, attacks may originate during data-transfer through the Internet or local networks [3]. Even in the case of complete protection from adversarial attacks, training existing deep models with adversarial examples or designing defense mechanisms [23] can improve model generalizability and resilience. In this paper, we present a comprehensive analysis of ten different adversarial attacks on classification of chest X-ray images and investigate how two different standard deep neural networks perform against adversarial perturbations. We perform both white (i.e. producing perturbed images using network A and classifying them by the same network) and black-box (i.e. producing perturbed images using network A and classifying them by network B) attacks.

## 2   Methods

### 2.1   Applied deep networks

We use two state-of-the-art deep models i.e. Inception-ResNet-v2 [16] and Nasnet-Large [24] to evaluate their performance on classification of both clean and perturbed chest X-ray images. Next, we modify the networks by replacing max-pooling operations with average-pooling to analyze whether the modified networks, especially the ones that are based on single/few pixel perturbation, are less sensitive to attacks. We hypothesize that average-pooling may be more resilient to attacks as it captures more global contextual information from the field of view, instead of selecting a single pixel candidate as max-pooling does.

## 2.2   Applied adversarial attacks

We applied three different categories of attacks namely gradient-based, score-based, and decision-based:

- **Gradient-based** attacks linearize the loss (in our case binary cross-entropy) around an input to which the model predictions for a particular class are most sensitive to. These attacks perturb the image with the gradient of the loss w.r.t. the clean image, gradually and efficiently increasing the magnitude until the model predicts a different label for the perturbed image. In our experiments, we have selected five different gradient-based attacks namely, Fast Gradient Sign Method (G1) [8], Projected Gradient Descent (G2) [9], Deep-Fool (G3) [10], Linfinity Basic Iterative Method (G4) [8],Limited-memory Broyden–Fletcher–Goldfarb–Shanno Method (L-BFGS) (G5) [18] and we demonstrate how the models trained on clean images perform against the crafted adversarial examples.
- **Score-based** attacks rely on confidence scores e.g. softmax class probabilities or logits to numerically estimate the gradient. From this group, we apply Local Search (S1) [11] (a black box attack based on the greedy local search algorithm to find pixels for which the model is the most sensitive and perturbing them to misclassify the input) and the Single Pixel (S2) [15] attacks.
- **Decision-based** attacks [2] solely rely on the predicted class or label of the model without requiring gradients or logits. From this group, we applied Gaussian Blur (D1), Contrast Reduction (D2) and Additive Gaussian Noise (D3) in our experiments. In all of the aforementioned attacks, a line-search is performed internally to find minimal perturbations required by the image to turn it into an adversarial example.

We trained both the networks from scratch with a batch size of 32 and 8 for training the Inception-ResNet-v2 and Nasnet-Large, respectively. RMSProp optimizer [19] with a decay of 0.9 and $\epsilon = 1$ and an initial learning rate of 0.045, decayed every 4 epochs using an exponential rate of 0.94 were used for all of our experiments as described in [16,24]. We set all attack parameters as proposed by their authors and utilized Foolbox [13], to craft adversarial examples.

## 3   Dataset

We use ChestX-ray14 dataset [20] which comprises 112,120 gray-scale images with 14 disease labels and 1 no-finding label. We treat all the disease classes as positive and formulate a binary classification task of "disease" vs. "non-disease". We randomly selected 95,128 images for training and 16,792 for validation. We randomly picked 200 unseen images as the test set, with 93 images with chest disease labels and 107 having "no finding" labels. These clean images are used for carrying out different adversarial attacks and the models trained on clean images are evaluated against them.

## 4    Results and discussion

Figure 1 shows the perturbed images produced by the ten different applied attacks. In Figure 2, we visualize a few samples where the perturbations are perceptible by human. We observed that most of the produced images by D1 (i.e Gaussian blur), D2 (i.e. contrast reduction), D3 (i.e additive Gaussian noise), S1 (i.e. local search) attack can be easily detected by the naked eye. We also found that S1 requires relatively more time compared to other methods to find an adversarial image.
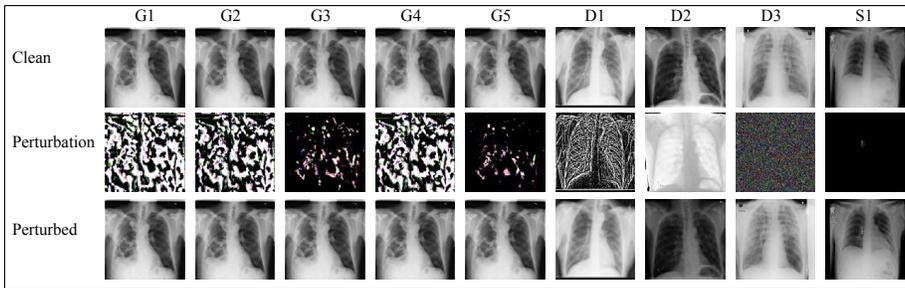


**Fig. 1.** Perturbed images produced by 10 (3 categories) different attacks.

In Tables 1 and 2, we report accuracy and area under ROC for two networks with/without modification for clean and ten different adversarial attacks (white- and black-box). Note that the single pixel attack [15] i.e. S2 (from the score based attacks category) failed to fool the networks for the entire test set which shows the single pixel attack works well on RGB (colored images) but not on gray-scale X-ray images as it is not simple to fool a deep model by changing only a single "gray-scale" pixel.

As reported in Table 1, the gradient based attacks were almost completely successful in fooling both networks (with/without modification) when the victim model for attack was the same reference model, i.e. in a white-box attack scenario. The decision and score based attacks were almost unsuccessful in fooling the models. We observed that Nasnet-Large with average pooling was 18% stronger in comparison to Nasnet-Large with max pooling. Note that the local search attack (S1) completely failed against Nasnet-Large with average-pooling.

In Table 2, we show the performance of both the networks against the black-box attacks i.e. we craft adversarial images with Inception-ResNet-v2, but test them with Nasnet-Large and vice versa. As reported in the table, almost all the methods were partially successful but not as high as white-box attacks. For gradient based black-box attacks, average pooling shows more resiliency against the attacks. We observed that for 23%±9% and 27%±8% of the test samples both the networks failed on the same cases for average and max pooling, respectively.

**Table 1.** Performance of original/modified Inception-ResNet-v2 (IR2) and Nasnet-Large (NL) against ten different *white-box* attacks. In the Table, MP, AP, Acc., AU refer to max-pooling, average-pooling, accuracy, and area under ROC, respectively.

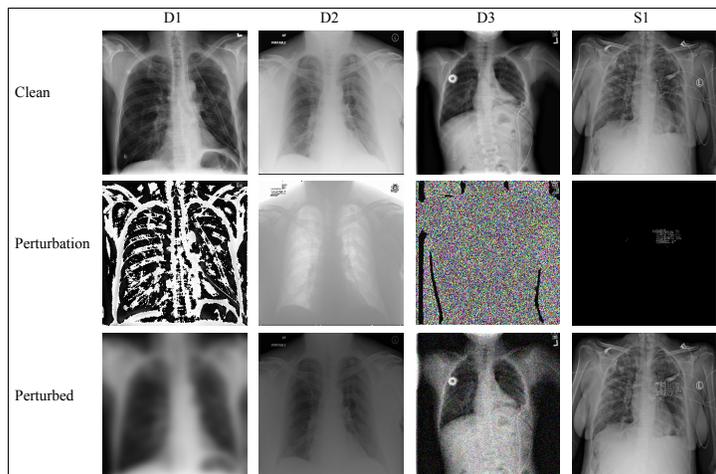| | Model | Metrics | Clean | Gradient | | | | | Decision | | | Score | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | G1 | G2 | G3 | G4 | G5 | D1 | D2 | D3 | S1 | S2 |
| MP | IR2 | Acc. | 0.70 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.04 | 0.10 | 0.32 | 0.65 | 0.70 |
| | | AU | 0.75 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.06 | 0.19 | 0.52 | 0.74 | 0.75 |
| | NL | Acc. | 0.73 | 0.00 | 0.00 | 0.01 | 0.00 | 0.00 | 0.06 | 0.41 | 0.30 | 0.32 | 0.73 |
| | | AU | 0.77 | 0.00 | 0.00 | 0.10 | 0.00 | 0.00 | 0.10 | 0.66 | 0.58 | 0.55 | 0.77 |
| AP | IR2 | Acc. | 0.71 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.04 | 0.24 | 0.14 | 0.62 | 0.71 |
| | | AU | 0.74 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.06 | 0.39 | 0.26 | 0.72 | 0.74 |
| | NL | Acc. | 0.72 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.03 | 0.41 | 0.48 | 0.72 | 0.72 |
| | | AU | 0.74 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 | 0.04 | 0.64 | 0.64 | 0.74 | 0.74 |



**Fig. 2.** Human perceptible adversarial perturbations

**Table 2.** Performance of original/modified Inception-ResNet-v2 (IR2) and Nasnet-Large (NL) against ten different *black-box* attacks. In the Table, MP, AP, Acc., AU refer to max-pooling, average-pooling, accuracy, and area under ROC, respectively.

| | Model | Metrics | Clean | Gradient | | | | | Decision | | | Score | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | G1 | G2 | G3 | G4 | G5 | D1 | D2 | D3 | S1 | S2 |
| MP | IR2 | Acc. | 0.70 | 0.46 | 0.43 | 0.43 | 0.43 | 0.43 | 0.53 | 0.81 | 0.36 | 0.45 | 0.70 |
| | | AU | 0.75 | 0.44 | 0.41 | 0.40 | 0.41 | 0.41 | 0.43 | 0.84 | 0.24 | 0.40 | 0.75 |
| | NL | Acc. | 0.73 | 0.53 | 0.51 | 0.51 | 0.51 | 0.51 | 0.57 | 0.58 | 0.74 | 0.56 | 0.73 |
| | | AU | 0.77 | 0.52 | 0.49 | 0.49 | 0.49 | 0.49 | 0.51 | 0.55 | 0.82 | 0.55 | 0.77 |
| AP | IR2 | Acc. | 0.71 | 0.51 | 0.52 | 0.52 | 0.52 | 0.51 | 0.53 | 0.29 | 0.40 | 0.53 | 0.71 |
| | | AU | 0.74 | 0.49 | 0.49 | 0.49 | 0.47 | 0.50 | 0.47 | 0.24 | 0.40 | 0.52 | 0.74 |
| | NL | Acc. | 0.72 | 0.59 | 0.58 | 0.58 | 0.58 | 0.58 | 0.49 | 0.53 | 0.51 | 0.38 | 0.72 |
| | | AU | 0.74 | 0.59 | 0.58 | 0.58 | 0.58 | 0.58 | 0.46 | 0.52 | 0.46 | 0.39 | 0.74 |

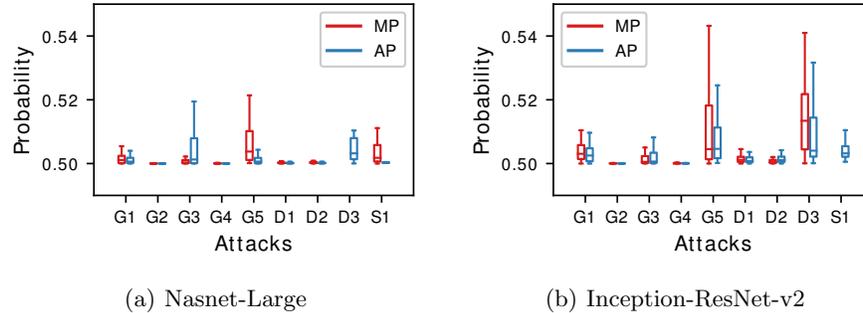(a) Nasnet-Large          (b) Inception-ResNet-v2

**Fig. 3.** Probability values of the original and modified networks after attack

In Figure 3, we show the probability values of the two networks (with/without modification) only after successful attacks on the disease class. Higher ranges/values indicate a stronger attack or a more vulnerable network. As shown in the figure, D3 (i.e. Additive Gaussian Noise), S1 (i.e. Local search) and G5 (i.e. L-BFGS) attacks are highly sensitive to the choice of pooling (max/average) operation. The range of the attack's confidence varies from $\sim 0.50$ to $\sim 0.55$. Note that absence of a box in the figure means there was no successful attack for a disease class in that experiment. In Figure 4, we visualize the accuracy of Inception-
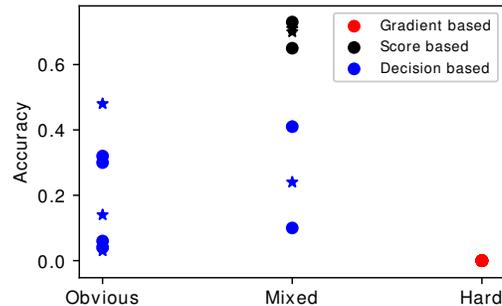


**Fig. 4.** Distribution of adversarial images crafted based on Inception-ResNet-v2 network w.r.t. human perception. The dots and stars in the figure refer to max and average pooling, respectively. The words Obvious, Mixed, and Hard refer to the level of difficulty for a human to perceive the attacks.

ResNet-v2 for different groups of attacks and, in the same plot, we show the perceptibility of each group of the attacks (i.e. the difficulty level for a human to detect a perturbed image). Note that lower accuracy and harder detection

(lower right of the plot) implies a more successful attack. As shown in the figure, gradient based attacks are the most successful ones in terms of fooling both human (i.e. perception) and machine (i.e. accuracy).

## 5   Conclusion

In this paper, we extensively tested the vulnerability of the two state-of-the-art deep classification networks against ten different adversarial attacks on chest X-ray images. We found that the single pixel attack completely failed for gray-level X-ray chest images. We also showed that the pooling operation can make a considerable difference for some attacks, even leading to a complete failure of the attack for a particular class. We also demonstrated that the crafted adversarial images with some of the attacks, e.g. Gaussian blur and contrast reduction methods, can be simply detected with the naked eye. Finally, we showed that the gradient based attacks applied to the chest X-ray images are the most successful in terms of fulling both machine and human. Although both networks, Inception-ResNet-v2 and Nasnet-Large, failed against gradient-based attacks, in general, the latter (with average pooling) was more resilient to decision and score based attacks.

## Acknowledgments

## References

1. Baltruschat, I.M., Nickisch, H., Grass, M., Knopp, T., Saalbach, A.: Comparison of deep learning approaches for multi-label chest x-ray classification. arXiv preprint arXiv:1803.02315 (2018)
2. Brendel, W., Rauber, J., Bethge, M.: Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. arXiv preprint arXiv:1712.04248 (2017)
3. Finlayson, S.G., Kohane, I.S., Beam, A.L.: Adversarial attacks against medical deep learning systems. arXiv preprint arXiv:1804.05296 (2018)
4. Guendel, S., Grbic, S., Georgescu, B., Zhou, K., Ritschl, L., Meier, A., Comaniciu, D.: Learning to recognize abnormalities in chest x-rays with location-aware dense networks. arXiv preprint arXiv:1803.04565 (2018)
5. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 770–778 (2016)
6. Kannan, H., Kurakin, A., Goodfellow, I.: Adversarial logit pairing. arXiv preprint arXiv:1803.06373 (2018)
7. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems. pp. 1097–1105 (2012)

8. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world. arXiv preprint arXiv:1607.02533 (2016)
9. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083 (2017)
10. Moosavi-Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2574–2582 (2016)
11. Narodytska, N., Kasiviswanathan, S.P.: Simple black-box adversarial perturbations for deep networks. arXiv preprint arXiv:1612.06299 (2016)
12. Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C., Shpanskaya, K., et al.: Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. arXiv preprint arXiv:1711.05225 (2017)
13. Rauber, J., Brendel, W., Bethge, M.: Foolbox v0. 8.0: A python toolbox to benchmark the robustness of machine learning models. arXiv preprint arXiv:1707.04131 (2017)
14. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
15. Su, J., Vargas, D.V., Kouichi, S.: One pixel attack for fooling deep neural networks. arXiv preprint arXiv:1710.08864 (2017)
16. Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A.: Inception-v4, inception-resnet and the impact of residual connections on learning. In: AAAI. vol. 4, p. 12 (2017)
17. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1–9 (2015)
18. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
19. Toshev, A., Szegedy, C.: Deeppose: Human pose estimation via deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1653–1660 (2014)
20. Wang, X., Peng, Y., Lu, L., Lu, Z., Bagheri, M., Summers, R.M.: Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In: Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on. pp. 3462–3471. IEEE (2017)
21. Yao, L., Poblenz, E., Dagunts, D., Covington, B., Bernard, D., Lyman, K.: Learning to diagnose from scratch by exploiting dependencies among labels. arXiv preprint arXiv:1710.10501 (2017)
22. Yuan, X., He, P., Zhu, Q., Bhat, R.R., Li, X.: Adversarial examples: Attacks and defenses for deep learning. arXiv preprint arXiv:1712.07107 (2017)
23. Zantedeschi, V., Nicolae, M.I., Rawat, A.: Efficient defenses against adversarial attacks. In: Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security. pp. 39–49. ACM (2017)
24. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. arXiv preprint arXiv:1707.07012 **2**(6) (2017)