(a) AOPC curve     (b) Weight+Gaussian     (c) AMM for adversarial training     (d) Disguising the bias of a model

1   **AOPC [R1,R2]** We employ Area Over Prediction Curve (AOPC) [Samek *et.al*, `arXiv:1509.06321`], a principled way

2   of quantitatively evaluating the validity of neural network interpretations, to check whether the manipulated model also

3   has been significantly *changed* by fooling the interpretation. From the definition of AOPC, we can conclude that the

4   interpretation is closely tied to the actual prediction procedure of the model if AOPC rapidly increases with the number

5   of perturbation of input pixels. Fig. (a) shows the average AOPC curves on 10K validation images for the original

6   and manipulated DenseNet121 models, $\boldsymbol{w}_o$ and $\boldsymbol{w}_{fool}^*$, with three different perturbation orders; i.e., with respect to

7   the $\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_o)$ scores, $\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_{fool}^*)$ scores, and a random order. We did the Top-$k$ fooling with $\mathcal{I}$ being Grad-CAM. From

8   the figure, we observe that $\boldsymbol{w}_o(\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_o))$ and $\boldsymbol{w}_{fool}^*(\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_o))$ show almost identical AOPC, which suggests that $\boldsymbol{w}_{fool}^*$

9   has *not changed much* from $\boldsymbol{w}_o$ and is making its prediction by focusing on similar parts that $\boldsymbol{w}_o$ bases its prediction.

10   In contrast, the AOPC of both $\boldsymbol{w}_o(\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_{fool}^*))$ and $\boldsymbol{w}_{fool}^*(\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_{fool}^*))$ lie significantly lower, even lower than the

11   case of random perturbation. From this, we can deduce that $\mathbf{h}_c^{\mathcal{I}}(\boldsymbol{w}_{fool}^*)$ is highlighting parts that are less helpful than

12   random pixels for making predictions, hence, is a "wrong" interpretation. We believe the notion of "good/bad/wrong"

13   interpretation (Re:**[R2]**) can be defined by examining the AOPC of $\mathbf{h}_c^{\mathcal{I}}(\cdot)$ and that of random perturbation. Furthermore,

14   the [Basic question] (Re:**[R1]**) says that [Ross *et.al*, 2017] can be contradicting to our result, but depending on the

15   training objective and model capacity, we believe the two phenomenon can both exist.

16   **Robustness of/detecting our fooling [R1,R3]** Detecting or undoing our model manipulation would not be easy as we

17   cannot easily access the original interpretation results or training data. Fig. (b) shows a feasible attempt, inspired by

18   [Roth *et al.*, `arXiv:1902.04818`], fails for detecting our manipulation. Namely, as the adversarial input examples can

19   be detected by adding small Gaussian perturbation to the input, we may also suspect that adding small Gaussian noise

20   to the parameters might reveal our fooling. However, Fig. (b) shows that $\boldsymbol{w}_o$ and $\boldsymbol{w}_{fool}^*$ behave very similarly in terms

21   Top-1 accuracy as we increase the Gaussian noise level, and FSR do not change radically, either. Another possible

22   detection scheme would be to use a black-box interpreter and compare its AOPC with that of the fooled model.

23   **[R1]** ① We will modify the acronym for Simple Gradient and correct the error regarding LRP passing the sanity

24   check. ② We found fooling SmoothGrad / Integrated Gradients while maintaining the accuracy is much harder, as the

25   reviewer has expected. (LIME is a black-box interpreter and is out of scope.) We will discuss this result in our final

26   version. ③ Fig. (c), which shows the accuracy of a ResNet50 model on $\mathcal{D}_{val}$ and PGD($\mathcal{D}_{val}$) (i.e., the PGD-attacked

27   $\mathcal{D}_{val}$), demonstrates that adversarially trained model can be also adversarially manipulated by our method. Namely,

28   starting from a pre-trained $\boldsymbol{w}_o$ (dashed red), we adversarially train the model with PGD attacks ($\epsilon = 1.5$) [Shafahi *et.al*,

29   `arXiv:1904.12843`] to obtain $\boldsymbol{w}_{adv}$ (dashed green), then started our adversarial model manipulation with Location

30   fooling with Grad-CAM. Note the Top-1 accuracy on $\mathcal{D}_{val}$ drops while that on PGD($\mathcal{D}_{val}$) increases during adversarial

31   training phase, and they are maintained during our model manipulation phase (e.g., at blue dashed line). The right panel

32   shows the Grad-CAM interpretations of two images at three distinct phases (see the color-coded boundaries), and we

33   clearly see the success of the Location fooling (blue dashed, third row). ④ We plan to release the code after acceptance.

34   **[R2]** ① We will carry out more thorough investigation on active fooling in the future work. Moreover, $\mathrm{LRP}_T$ is our

35   novel variation of LRP, and the local characteristic of $\mathrm{LRP}_T$ is maintained to LRP since the relevance values of $\mathrm{LRP}_T$

36   are propagated down to the input pixel level. Directly fooling with LRP did not work well. ② We have randomly

37   selected the visualization examples, and have not cherry-picked. ③ We argue the correlation metric also is not enough to

38   accurately measure the success of fooling; e.g., in Location fooling, the mere correlation cannot tell whether the fooling

39   was successful. Instead, we devised FSR to more directly measure the success of intended fooling. For computing FSR,

40   we simply randomly selected 10K images from the ImageNet validation, and $\ell_2$, a common measure of distance in 2-D

41   images, is used in FSR for Center-mass fooling even though $\ell_1$ and $\ell_2$ give almost the same result. ④ For L248/L255

42   comments, please refer to the AOPC results. ⑤ We will emphasize the contrast between the random parametrization.

43   **[R3]** ① Our method does share some similarities with backdooring; e.g., manipulating model parameters to hide some

44   intention while maintaining the accuracy. However, the setting/objective are radically different. Namely, backdooring

45   has nothing to do with the *interpretation* of a model, and it could even be detected by the interpretation methods. But,

46   with our model manipulation, we can even make the backdooring stronger to be not caught by the interpretation methods.

47   ② To support the motivation of our work, we carried out additional experiment on the 'Adult income' classification data

48   in [`UCI ML Repository`]. We trained a classifier with 8 conv layers, $\boldsymbol{w}_o$, and the LRP result in Fig. (d) (blue) shows

49   it assigns high importance on *sensitive* features like 'Race' and 'Sex'. Now, we (or a lazy developer) can manipulate the

50   model with Location fooling that zero-masks the two features and obtain $\boldsymbol{w}_{fool}$ that essentially has the same accuracy

51   of 76.8% as $\boldsymbol{w}_o$ but with a new interpretation that disguises the bias (orange). Obviously, from our discussions on

52   AOPC etc. above, the lazy developer has not resolved the bias in the model but simply fooled the interpretation result.