

1 ===== TO REVIEWER No.2 =====

2 [Q1.1] On similar approaches.

3 [A] Other than the discussions in the related works section, we wish to highlight that, the difference between our
4 work and the usual "data poisoning" works is as follows: In data poisoning setting, the attacker's goal is to add or
5 modify as few samples as possible (i.e., "an attacker usually cannot directly access an existing training database but
6 may provide new training data" [2]), whereas our proposed framework put the constraint on the perturbation levels
7 (as human imperceptible as possible). In addition, to our best knowledge, this is the first scalable solution that can
8 run on big datasets (thanks to the optimization method we designed). Other SOTA methods for poisoning was only
9 tested on datasets of small sizes (e.g., [20] used only 1,000 samples). In other words, a direct comparison is less fair for
10 the standard poisoning algorithms. Nevertheless, we will include some more actual quantitative assessments for such
11 comparison to make the points clearer in the revised version. Thank you for the suggestions. (Please also refer to the
12 similar response to [Q2.1] and [Q3.2].)

13 ===== TO REVIEWER No.4 =====

14 [Q2.1] On comparable baseline methods.

15 [A] Please also see [Q1.1]. We would also like to emphasize that, our work is not about poisoning only but rather on
16 high-dimensional adversarial training examples with bounded perturbation. Therefore, we designed several different
17 experiments with different angles toward this goal. Other existing methods for poisoning seems to be less aligned with
18 our propose, making such comparisons less attractive. Meanwhile, as reviewer 2 also pointed out, we actually did
19 conducted some baseline method such as Random Flip (Figure 5) for compare but only for validating purposes.

20 [Q2.2] "Is modifying 60% a practical setting?"

21 [A] Actually, yes. Unlike the traditional poisoning task, adversarial training data can be used in a good way. For
22 instance, in some applications an agent may agree to release some internal data for academic research, but does not
23 like to enable the data receiver to build a model which performs well on real test data; this can be realized by applying
24 DeepConfuse before the data release. For motivations on privacy and why not just using synthetic data for these tasks,
25 please also see our response to [Q3.2]. We will make the above points clearer in the revised version and to revise the
26 typos/add the omitted related works as suggested. Thank you for the suggestions.

27 ===== TO REVIEWER No.6 =====

28 [Q3.1] How about some visualizations and interpretations for the corresponding models?

29 [A] Thank you for the suggestion. For linear SVMs, we did a quick visualization on the weights trained on clean and
30 adversarial training data as shown below:

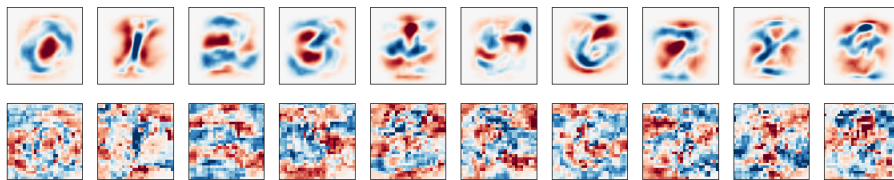


Figure 1: LinearSVM weights visualization for MNIST. Top row: Weights trained on clean training data. Bottom Row: Weights trained on adversarial training data.

31 It can be shown that the weights for the SVM trained on adversarial data indeed went to the opposite direction
32 compared with the corresponding clean model and trend to overfits on image corners. We will definitely add some more
33 experimental results with a more detailed discussion in the revised version as suggested.

34 [Q3.2] On some more motivations.

35 [A] Please also see [Q1.1] and [Q2.1]. For data privacy aspect, our proposed approach is quite different from releasing
36 synthetic data via GANs. Consider a company selling surveillance cameras and the user(usually the police) will store
37 all the photos been taken. These photos cannot be synthetic for obvious reasons. On the other hand, such a company
38 does not want any other third parties to train a classifier using their data. Then DeepConfuse is suitable for this kind of
39 task and the camera company can just send the modified data in real-time. We will elaborate more on the motivation
40 side in the revised version.

41 Thank you all for your thorough and insightful comments again.