

1 We would like to thank all reviewers for your helpful suggestions.

## 2 **Response to Reviewer #1**

3 We are sorry about the confusing descriptions in this submission, in the updated version we will check every symbol to  
4 make sure they are clear. In what follows we answer your questions to help clarify the notations and experiments.

5 1. **“motivation not clear”**: G-SSL and specifically label propagation used to be an important semi-supervised learning  
6 algorithm and it is still not outdated today: for image/text data we can apply the deep feature to construct a graph and  
7 propagate the labels to unlabeled data. Neo4j, an industrial level machine learning library even contains a build-in  
8 module for label propagation and it has already been applied in social network [Speriosu], Pharmacy [Zhang], etc.  
9 Furthermore, G-SSL algorithms are still widely used in social network or recommender systems, where malicious  
10 users can easily create fake account or inject fake information into nodes. Lastly, the clean math of G-SSL provides  
11 provides a good starting point to generalize adversarial machine learning to other SSL methods, including deep SSL.  
12 We will try to include more real-world examples in the revised version.

13 2. **“ $X_l$  and  $X_u$ ”**: following the standard notations, we use  $X \in \mathbb{R}^{n \times d}$  to represent the feature matrix for the whole  
14 data, where each row of  $X$  is a data point. The feature matrix is composed of labeled part and unlabeled part. We  
15 simply set  $X = [X_l; X_u]$ , i.e. the first  $n_l$  rows of  $X$  are labeled data and the following  $n_u$  rows are unlabeled data,  
16 and  $n = n_l + n_u$ . With this in mind, it might be easier to understand  $X_l + \Delta_x$ —this is the feature of labeled data  
17 after perturbation by matrix  $\Delta_x \in \mathbb{R}^{n_l \times d}$ , which means the feature perturbation can be done to all labeled instances.  
18 Moreover, by applying group sparsity constraints on  $\Delta_x$  (consider each row as a group), our algorithm can conduct  
19 perturbation only to a small fraction of instances to greatly change the results.

20 3. **“subscriptions  $i, j$ ”**: each row of  $X$  is a data point.  $x_i$  and  $x_j$  are the  $i$ -th and  $j$ -th row of the matrix  $X$  so  $i$  and  $j$   
21 are actually the index of data ( $i, j \in \{1, 2, \dots, n\}$ ), and they are iterated within the whole dataset (including both  
22 labeled and unlabeled data). However, the constraint in Eq(1) indicates that  $\hat{y}$  are fixed for the labeled part, so the  
23 free variables in (1) are the unlabeled part of  $\hat{y}$ .

24 4. **“Datasets not explained”**: we will add more descriptions of datasets in the next version. Specifically, `mnist` and  
25 `rcv1` are widely used datasets to test non-deep learning models such as SVM/logistic regression/tree based models,  
26 `cadata/E2006` are regression datasets for predicting the house price/stock volatility. All of them are publically  
27 available at `libsvm` repository.

28 5. **“Scalability of our algorithm”**: please note that the computation overhead mainly lies in generating the kernel matrix  
29  $S$  in Eq(3), which is  $\mathcal{O}(n^2d)$ , computing this matrix for 20,000 nodes is certainly doable. For large data, one can  
30 construct a  $k$ -NN graph efficiently by efficient similarity search tools, such as `Faiss` (by Facebook).  $S$  will be sparse  
31 so the computations will be efficient even for very large datasets.

32 6. **“Dimension of MNIST17 should be 784 rather than 780”**: 4 in 784 pixels are zero in every image, so we remove  
33 these four dimensions in the experiment. Removing them will not affect the final results.

34 (Speriosu) Twitter polarity classification with label propagation over lexical links and the follower graph, in *EMNLP* 2011.

35 (Zhang) Label Propagation Prediction of Drug-Drug Interactions Based on Clinical Side Effects, in *Scientific Reports* 2015.

## 36 **Response to Reviewer #2**

37 Thanks for your feedback! We will correct typos and revise our submission according to your suggestions.

## 38 **Response to Reviewer #3**

39 **Regarding the scope of the paper**: The main text of our paper is mainly on poisoning label propagation method. We  
40 further show how to poison another widely used SSL method – manifold regularization in Appendix 4.2.

41 **Regarding the datasets**: the datasets used in the experiment are all real world datasets (downloaded from `libsvm` data  
42 repository) and are widely used for benchmarking classification and regression models. We will provide more details  
43 about these datasets in the revised version of the paper.

## 44 **Response to Reviewer #4**

45 **Regarding the novelty of theory**: Our main theoretical contribution lies in the efficient trust region solver in the  
46 nonconvex case. To our knowledge, prior work is limited to trust region with cubic regularization [Carmon and Duchi],  
47 which is different from our setting (a  $\ell_2$ -norm constraint), and their work cannot be directly applied to our data poisoning  
48 problem. Furthermore, earlier trust region solvers (e.g., Dogleg or Steihaug’s method) can only get an approximate  
49 solution or local minimizer, while our algorithm is guaranteed to return the global minimizer.

50 **Regarding the novelty of algorithm**: Finding good “label flipping” operations via probabilistic method is novel. We agree  
51 that this method has no theoretical guarantee, however we show empirically that this epsilon-greedy solver successfully  
52 avoids some sub-optimal solutions, and beats the greedy method by a significant margin.

53 **Regarding a unified framework**: For completeness, we also include manifold regularization method in the appendix. It  
54 is certainly preferable to design a framework that include all semi-supervised learning approaches. Regarding to this  
55 limitation, we do not intend to solve it completely in this submission and we appeal for follow-up works to give a better  
56 answer.

57 (Carmon and Duchi) Gradient Descent Efficiently Finds the Cubic-Regularized Non-Convex Newton Step, in *ArXiv* 2016.