

1 We thank reviewers for the constructive comments. We will include more intuition about the approach, improve the  
2 presentation and fix all minor issues in the final version.

3 Re: Worst-case distribution. Minimizing the loss of the worst-case distribution  $\mathbb{Q}$  implies an optimization over all  
4 distributions within the ball of an appropriate radius  $\epsilon$  (see Eq. (1)), which could also include the unknown real  
5 distribution  $\mathbb{P}_N$ . Though the worst-case  $\mathbb{Q}$  may not be exactly the real  $\mathbb{P}_N$ , *the classifier (i.e.  $\theta$ ) must have fitted  $\mathbb{P}_N$*   
6 *better than (or equivalently with)  $\mathbb{Q}$* , as the classification error over  $\mathbb{Q}$  is the worst. In iterations, the worst-case  $\mathbb{Q}$  will  
7 be dynamically determined by the classifier, and the classifier will fit the real  $\mathbb{P}_N$  increasingly better in an implicit way.

8 To #R1. Re: Reasons of working better. First, generators in existing methods tend to fit the empirical distribution.  
9 Given a bad training set, their generated data could be worse. Second, these generators often produce “easy” samples  
10 by cooperating with the classifier, and a nearly duplicate copy of the given bad training data could be sufficient but will  
11 be useless to estimate the real data distribution. In contrast, our generator plays against the classifier and the capability  
12 of the classifier can be largely enhanced over a distribution ball.

13 Re: G in Eqs. (10) and (14). RHS of Eqs. (10) and (14) indeed depends on G. The last two terms are expectations over  
14 the distribution  $\mathbb{Q}$ , which is approximated by  $G(z)$ , as explained in line 93 of the paper.

15 Re:  $\inf_{\lambda>0}$  in Eqs. (4)-(5). In Eq. (4),  $\inf_{\lambda>0}$  is over both two terms of RHS. A pair of braces will be used to avoid  
16 misleading.  $\lambda$  in Eq. (5) indicates the corresponding optimal for  $\inf_{\lambda>0}$ , and it will change with the bound  $\epsilon$  in Eq. (3).  
17 Since  $\epsilon$  is unknown, it is common to take  $\lambda$  as a hyper parameter to be tuned in experiments (e.g. Theorem 1 in [1]).

18 Re: Good data. This comment is insightful. Though training data are clean, there still probably exists a *subtle* gap  
19 between distributions of training and test data. Moreover, the generator could conduct “data augmentation” for the  
20 classifier. We may thus receive a slightly better result, e.g. 99.42 v.s. 99.35 on MNIST under imbalance= 1.

21 To #R3. Re: Large  $\epsilon$ . A large  $\epsilon$  leads to a set  $\mathbb{B}_\epsilon$  of huge capacity (see Eq. (2)), which could be flooded with distributions  
22 that are far away from both the empirical distribution  $\hat{\mathbb{P}}_N$  and the real data distribution  $\mathbb{P}_N$ . It is therefore reasonable to  
23 set  $\epsilon$  within an appropriate range, as what we usually do with hyper-parameters in machine learning.

24 Re: Difference from adversarial data augmentation. Firstly, adversarial data augmentation usually aims to create adver-  
25 sarial copies of training data by adding perturbations. In contrast, we generate new sample from a distribution. Secondly,  
26 the amount of perturbation (pixels for image) is often constrained in classical methods, while we focus on a distribution  
27 ball with a radius bound. In addition, Triple GAN and Triangle GAN tend to generate samples that can well fit the  
28 classifier (i.e. cooperation between G and C), but we improve the classifier by challenging it with the generator.

29 Re: Eq. (8) and  $\mathbb{Q}_i$ . We more clearly express the second line of Eq. (8) as  $\sup_{\mathbb{Q}_i} \frac{1}{N} \sum_{i=1}^N \int_{\mathcal{X}} \ell_\theta(\mathbf{x}, y) \mathbb{Q}_i(d(\mathbf{x}, y)) - \lambda \cdot$   
30  $\sup_{f \in \mathcal{F}} \left\{ \frac{1}{N} \sum_{i=1}^N [f(\mathbf{x}_i, y_i) - \int_{\mathcal{X}} f(\mathbf{x}, y) \mathbb{Q}_i(d(\mathbf{x}, y))] \right\}$ .  $\mathbb{Q}_i$  is the conditional distribution of  $(\mathbf{x}, y)$  given  $(\mathbf{x}_i, y_i)$ .

31 The joint distribution  $\Pi$  of  $(\mathbf{x}_i, y_i)$  and  $(\mathbf{x}, y)$  with marginals  $\hat{\mathbb{P}}_N$  and  $\mathbb{Q}$  respectively (see Eq. (3)), can be written as  
32  $\Pi = \frac{1}{N} \sum_{i=1}^N \delta_{(\mathbf{x}_i, y_i)} \otimes \mathbb{Q}_i$ . According to the law of total probability, we can factorize  $\mathbb{Q}$  as the first line of Eq. (8).

33 Re: True distribution. Sentences around Line 267 will be rephrased. Minimizing the worst-case expected loss implies  
34 an optimization over all distributions in the  $\epsilon$ -ball where the real data distribution is also expected to be included.

35 Re: Comparison with data augmentation. The results shown in the first line of Table 1 are obtained by a combination  
36 of randomly clip, horizon flip, and rotation. The second line is results of Mixup [2]. Table 1 shows that our method  
37 outperforms both two comparison methods. Mixup also provides reasonable improvement but is not as outstanding as it  
38 on regular datasets. Other GAN-based data augmentation methods have been included in Tables 1 and 2 of the paper.

39 To #R4. Thanks for your support. Source codes will be released.

Table 1: Accuracy (%) on CIFAR-10.

40 To #R5. Re: Explanation of improvement. The algorithm receives per-  
41 formance improvement, as the classifier can be implicitly optimized over  
42 the real data distribution with the help of the worst-case distribution.  
43 Our generator plays against with the classifier within a Wasserstein ball,  
44 so that the capability of the classifier can be enhanced during the iterations.

Method	IF = 10	G(0.2)	Reduced
Combination	85.63	85.25	83.59
Mixup [2]	86.04	85.66	83.91
Ours	<b>86.86</b>	<b>85.87</b>	<b>84.60</b>

45 Re: Eq. (9). Thanks for your advice. We will improve Eq. (9) to make it clearer.

46 Re: JSD and WD. By replacing the critical network D in WGAN with a discriminator, we can easily obtain Eq. (14).  
47 With the help of Theorems 4 and 5, the standard GAN framework can be view as a JSD version of framework defined in  
48 Eq. (10). The difference between them is the distance metric used in Eq. (2).

49 [1] Kloft, M., Brefeld, U., Laskov, P., et al., 2009. Efficient and accurate lp-norm multiple kernel learning. NIPS.

50 [2] Zhang, H., Cisse, M., Dauphin, Y.N. and Lopez-Paz, D., 2017, Mixup: Beyond empirical risk minimization. ICLR.