

1 We thank all the reviewers for the helpful reviews. We respond to each reviewer’s specific questions here.

2 To Reviewer #1:

3 1. *The proposed \mathcal{L}_{DMI} is not limited to DNN. Why does the paper focus on DNN practice?* Yes, \mathcal{L}_{DMI} is not limited
4 to DNN and can be applied to general settings. Our theory is valid in the general settings. Our experimental setting
5 focuses on DNN practice since the noisy-labels problem in data-driven deep learning is especially important as large
6 high-quality data is crucial to data-driven deep learning but may be extremely hard to obtain (see intro), thus a series
7 of works [21, 11, 6, 36] also focus on designing noise-robust functions for DNN. Applying our method to non-DNN
8 settings can be a future direction but we believe our current results not only make a significant theoretical step in the
9 general noisy-label learning problem but also make a significant empirical step in the DNN practice.

10 2. *The sufficiency of the paper’s literature review and experimental comparison are kind of weak.* We have compared
11 our method with some very recent baselines like **GCE** [36] in NIPS 2018, **LCCN** [33] in AAAI 2019 and also cited
12 several works that focus on network structure or learning paradigm (see the second and third paragraphs in related
13 work). We would appreciate if you could provide our missing works and will add them in our final version.

14 To Reviewer #2:

15 1. *Once the latent true variable Y is not identifiable, how to find the optimal classifier for Y ?* Yes, there exists $Y' \neq Y$
16 such that conditioning on Y' , X is also independent of \tilde{Y} . $Y' = \tilde{Y}$ is one example. However, first, training with Y or Y'
17 or \tilde{Y} will obtain the same optimal classifier based on our main theorem. Second, based on the information-monotonicity
18 of \mathcal{L}_{DMI} , if Y is more informative about X than Y' (e.g. when $Y' = \tilde{Y}$), the optimal classifier we obtained will be more
19 similar with Y than Y' . When we assume the ground truth is the most informative among all variables that satisfy
20 the conditional independent assumption, the optimal classifier will be closest to the ground truth. We will add this
21 clarification in our final version.

22 2. *Why it doesn’t work well on small noise?* We have the clean data comparisons since $r = 0$ is the clean case. **DMI**
23 outperforms all other counterparts except **GCE** [36] when $r = \{0.0, 0.1, 0.2\}$ in the Dogs vs. Cats dataset and when
24 $r = 0.0$ in CIFAR-10. **GCE** does a training optimization on **MAE** [6] with some hyperparameters while sacrifices the
25 robustness a little bit theoretically. Our next step may be to employ some training optimizations in our method.

26 3. *A Rate of ... papers.* Thanks for your information! We will cite them and make it clear in the final version.

27 To Reviewer #3:

28 1. *Why converting to two classes and other baselines are missing for Fashion-MNIST?* We use it as an explanation
29 experiment to compare distance-based (we use **CE** as an example here) and information-theoretic (our method) loss
30 functions (see the third paragraph in intro). For clean presentation, we convert to two highly imbalanced classes and
31 only compare with **CE**. We will make it clear in our final version. We attach the comparisons with other methods here
32 (see the following figures) and will include them in Appendix in our final version.

33 2. *Why are CIFAR-10 uniform, Dogs vs. Cats uniform and dog->cat missing?* Due to space limitation, we omit some
34 under-representative experiments although we have done them. For the uniform noise, first we want to clarify that noise
35 patterns are divided into two main cases, diagonally dominant and diagonally non-dominant and uniform noise is a
36 special case of diagonally dominant noise. Thus, we did not emphasize the uniform noise results in the submission. We
37 did not present the results for dog->cat since, unlike the highly imbalanced Fashion-MNIST (90% clothes, 10% bags),
38 Dogs vs. Cats is a balanced dataset (50% dogs, 50% cats) and thus the dog->cat results are very similar to cat->dog
39 results. We attach these results here (see the following figures) and will include them in Appendix in our final version.

