

1 We thank all reviewers for the constructive comments! Following the suggestions, here is a list of new experiments.

2 (i) Added a baseline *Truth* for label shift where the importance weight is the validation/train ratio for each class:

3

Acc		ρ 100	Truth: 83.05 (0.58)	DIW3-L: 83.69 (1.21)		ρ 200	Truth: 79.92 (0.46)	DIW3-L: 81.38 (1.24)
-----	--	------------	---------------------	----------------------	--	------------	---------------------	----------------------

4 Truth seems slightly worse since it is just roughly tuned; the difference is *statistically insignificant* by *t*-test.

5 (ii) Computed the ℓ_2 distances between the true weights (by Truth) and the estimated weights (by other methods):

6

Dist		$\rho = 200$	IW: 0.0324 (0.0010)	Reweight: 0.0321 (0.0010)	DIW3-L: 0.0166 (0.0003)
------	--	--------------	---------------------	---------------------------	-------------------------

7 The weights by DIW3-L are *statistically significantly closer* to the true weights. It is a unit test possible under label shift, which is why we wanted the label shift experiments, but we didn't finish them before the deadline.

9 (iii) Tested Adam as the optimizer (we can see that the messages *do not change much* if changing SGD to Adam):

10

(CIFAR-10, 40% symmetric noise) Acc		IW: 44.51 (1.55)	Reweight: 72.96 (0.97)	DIW3-L: 79.80 (0.25)
-------------------------------------	--	------------------	------------------------	----------------------

11 (iv) TODO: add text-data experiments (based on *20News* groups) and real-world experiments (based on *Clothing1M*).
12 For *20News* groups, we need to choose the network architecture, which is not a trivial task. For *Clothing1M*, we
13 are using the common practice ResNet-50, but the issue is that it goes in a speed of 2 days/epoch...

14 (v) TODO: add covariate shift experiments based on MNIST, Fashion-MNIST, or CIFAR-10. It is quite difficult to
15 simulate covariate shifts based on benchmark datasets where we *cannot access/manipulate* $p(x)$ or $p(x|y)$, even
16 though covariate shift should be extremely popular in the wild. Most of covariate shift papers made use of fully
17 controlled toy data where *IW suffices* and there is *no circular dependency*.

18 We have two plans: first, we may try strongly regularized GANs (e.g., very early stopped); second, we may try
19 mixup style data generation (mixing two images taken from two classes). We are thinking of ways to guarantee
20 that the shift is precisely covariate shift and the shift is challenging to learning methods, i.e., $p(x)$ should change
21 significantly, $p(x)$ must have a larger support in training, $p(y|x)$ cannot change at all. Any advice is welcome.

22 In the next version, we will definitely include all the aforementioned experiments!

23 **Use of validation data and motivation of the problem setting** (by R2 & R4 & R5) Our setting is the general IW,
24 with the same data generation as *Learning to reweight* [38] and *MentorNet* [22] (but a different goal from them).

25 First of all, note that there is no *distributional assumption* about the *joint distribution shift*. Given that $p(x, y)$ can
26 shift now, without any further information, the shift is *not identifiable*, or equivalently, the weights are *not estimable*.
27 Hence, we should let the learning methods see some *clean validation data* coming from the test distribution. By careful
28 algorithm design, the learning methods should be able to infer the shift by referring to the clean validation data.

29 On the other hand, what was done in more specific shifts (i.e., covariate shift, label/class-prior shift, and label noise)
30 is that a specific distributional assumption is added instead of the clean validation data to the algorithm design:

- 31 • under covariate shift, only $p(x)$ can shift, and $p(y|x)$ must remain the same;
- 32 • under label shift, only $p(y)$ can shift, and $p(x|y)$ must remain the same;
- 33 • under label noise, only $p(y|x)$ can shift, and $p(x)$ must remain the same.

34 The learning methods should be able to infer *the parameters but not the type* of the shift, because the type is already
35 given by the assumption. In this sense, *IW is more general* than these specific shifts without clean validation data; the
36 *data-driven philosophy/methodology* can tackle more complex shifts where such an assumption is not obvious.

37 Last but not least, the small set of clean validation data can be obtained with the help of domain experts. Since it is
38 small, it would not cost too much of the data-labeling budget. Some real-world label-noise datasets are designed in this
39 manner, for example, *Clothing1M* whose clean training data can serve as our clean validation data.

40 **Baselines are not latest for learning with noisy labels** (by R3) The problem setting of the proposal is IW rather than
41 label noise that is only a special case of IW. The latest label-noise methods *may not be applied* in the more general IW.
42 Among the learning methods under the same problem setting (i.e., how the data look like), [38] *is already the latest*.

43 **Class imbalance to label shift** (by R4) Thanks for pointing the fact out! Yes, class imbalance refers to *cost-sensitive*
44 *learning* under *no label shift*, but the goal is to maximize the AUC, which has been shown as same as to minimize the
45 *balanced error* [Menon et al., ICML'15]. In our experiments, the test data is balanced so that the classification error
46 is also the balanced error. That being said, our terminology following [5] is not as good as label shift. We will call it
47 class-prior or *prior probability shift* following the book [37], since label shift may also be confused with label noise.

48 **On training time and hyperparameter choice** (by R4) Yes, we have used many tricks including regularizations and
49 learning rate decays. In fact, not only weighted ERM with IW but also standard ERM relies on the training time and
50 hyperparameter choice. Our contribution was to find that *complex data form is not necessarily the bottleneck of IW*.

51 **No experiments for the hidden-layer-output transformation version in Algorithm 1** (by R5) In our experiments,
52 *this version corresponds to “-F” methods*, and the other version corresponds to “-L” methods. Indeed, the experimental
53 results you wanted are reported and discussed in Section 5.2 (more specifically, in Table 2 and Figure 5).