

1 We thank the reviewers for their very constructive feedback! We are encouraged that most reviewers (R1,R2,R4)<sup>1</sup> found  
2 our paper to be well-written and that our experiments were thorough/extensive and reproducible. Space prohibits us  
3 addressing all points raised by the reviewers so below we focus on the most pressing points from each reviewer.

4 **@R1,R2,R3,R4:** We will use the final submission’s extra allowed page to clarify all points discussed below.

5 **@R2 Is our absolute-value correction the same as Kiryo et al.’s approach?:** No. Our approach is different in two  
6 ways. First, during optimization, when  $\widehat{R}_u^-(g) - \pi\widehat{R}_p^-(g) < 0$ , Kiryo et al.’s update gradient is:  $\gamma\nabla_{\theta}(-\widehat{R}_u^-(g) + \pi\widehat{R}_p^-(g))$ .  
7 In contrast, our method’s gradient is:  $\nabla_{\theta}(\pi\widehat{R}_p^+(g) - \widehat{R}_u^-(g) + \pi\widehat{R}_p^-(g))$ . Put simply, Kiryo et al. stop optimizing  $\widehat{R}_p^+(g)$   
8 whenever  $\widehat{R}_u^-(g) - \pi\widehat{R}_p^-(g) < 0$  whereas we always optimize  $\widehat{R}_p^+(g)$ . This means that we spend comparatively more  
9 time optimizing the positive-labeled risk. Second, when estimating performance (i.e., risk) on validation data, our  
10 approach penalizes implausible negative risk estimates while Kiryo et al.’s method does not.

11 Empirically, we find that this “soft-constraint” approach to implausible negative risk yields comparable or better models  
12 (see supplemental Sec. E.5). We also show in the supplementals (e.g., Sec. C.(2) and Alg. 3) that all of our methods can  
13 work with the Kiryo et al. approach, albeit with a (much) more complex implementation.

14 **@R4 Assumption of same bias between positives in positive-labeled set  $\mathcal{X}_p$  and training unlabeled set  $\mathcal{X}_{tr-u}$ :** In  
15 our problem setting,  $\mathcal{X}_p$  is an unbiased sample of the positive examples in  $\mathcal{X}_{tr-u}$ , but  $\mathcal{X}_p$  could be arbitrarily different  
16 from the positive examples in test unlabeled set  $\mathcal{X}_{te-u}$ . However, if  $\mathcal{X}_p$  is not representative of the positives in  $\mathcal{X}_{tr-u}$ , then  
17 our two-step methods could still be adapted to handle this case by using any bPU algorithm in step 1 to create surrogate  
18 negative set  $\tilde{\mathcal{X}}_n$ , and our wUU estimator would be used in step 2 as normal. Thank you for this suggestion; we will note  
19 this additional flexibility of our two-step methods in the paper.

20 **@R2 Why do our two-step methods sometimes outperform our joint approach (PURR)?:** Supplemental tables 13  
21 and 15 detail many empirical setups where our PURR estimator outperforms both the baselines, PUC and nnPU\*, and  
22 our two-step approach. Nonetheless, the somewhat counterintuitive finding that a joint aPU learning method is not  
23 consistently the top performer is itself interesting and supports PURR’s inclusion in the paper.

24 As an intuition, note that PURR, with its three risk decompositions/corrections, is strictly harder to optimize than  
25 wUU, aPNU, and nnPU, which each have one correction. This can lead to worse accuracy compared to the two-step  
26 methods, especially on easier problems like MNIST, where each step can be solved accurately on its own. When the  
27 dimensionality is low and dataset sizes are small (see Sec. E.4), PURR is generally the top performer since PURR is  
28 able to extract more information by better combining the limited available data.

29 **@R2 How to empirically compare our methods to the baselines given the implementation differences:** For PUC,  
30 we used the original authors’ implementation, which relies on kernel methods for density-ratio estimation. For baseline  
31 nnPU\* as well as our methods (none of which require density estimation), we followed the standard of most recent  
32 PU learning work (including Kiryo et al. in their nnPU paper), which uses neural networks.

33 However, our experiments show that PUC’s biggest limitation is not its representation: On *unshifted* data (Table 1 row 1  
34 for each dataset), PUC’s performance is close to nnPU\* and slightly outperforms our methods. On *shifted* data (Tab. 1  
35 rows 2 & 3 for each dataset), PUC’s & nnPU\*’s performance degrades while our methods’ performance improves. Thus,  
36 data shift (and methods for handling it) are the biggest factor in performance. See also Sec. E.1 which shows simple  
37 aPU learning tasks that our methods can successfully learn while PUC cannot – specifically with identical models.

38 We will add a “Discussion” subsection to the paper’s “Experimental Results” (Sec. 7) to include the additional context  
39 in our response to this question as well as to the previous question on our joint method PURR.

40 **@R4 Negative-Class Shift in Epidemiology Example:** We expect that negative-class shifts will be small relative to  
41 shifts in the positive class, both in epidemiology and many other domains. Our experiments on TREC spam emails & in  
42 Sec. E.6 show that our methods are still useful even when our assumption is violated and the negative class shifts some.

43 **@R1 Problem Variations:** Sakai & Shimizu discuss additional problem variations including PU learning with arbitrary-  
44 negative shift, which they show is trivially equivalent to standard PU learning using positive set  $\mathcal{X}_p$  and test unlabeled  
45 set  $\mathcal{X}_{te-u}$ . If both classes shift, learning is impossible without additional datasets or assumptions (e.g., consistent  
46 input/output for covariate shift). This expanded discussion of problem variations will be added to the paper.

47 **@R3 Not Reproducible:** Supplemental Sec. D enumerates our complete experimental setup, and our submission  
48 included runnable source code to replicate all experiments. We respectfully believe our experiments are reproducible.

49 **@R3 Notation Used Before Defined?:** We will update the paper to clarify that line 70’s notation is based on the  
50 distribution notation listed at the end of the preceding paragraph (lines 65–66).

---

<sup>1</sup>We denote each reviewer  $RX$ , where  $X$  is the corresponding reviewer ID in CMT. Each reviewer identifier is also color coded.