(a) Synthetic source data: blondness perfectly correlates with the male gender.  (b) Synthetic target data: each gender has a variety of hair colors.  (c) Samples corrected by self-training were mostly mistakenly predicted because of the spurious correlation.

1 We thank the reviewers for the detailed and insightful feedback. The reviewers noted that the paper "target[s] a timely
2 and important problem" and "posits/analyzes a structured model... to study self-training" [R1], "the theoretical claims
3 are well established" [R4], and "the experimental results support the theoretical implications" [R3]. We will address the
4 major points and incorporate others in the next revision.

5 **[R1]:** "missing ... a bound on the target accuracy of the final classifier (in analogy to Theorem 3.2 which studies a
6 simpler setting). Clarification on why this is not provided or difficult to provide ... would be useful."
7 • The strongest result obtained, as noted in Line 211-212, is that the final classifier is a min-entropy solution, or max
8 unsupervised margin solution, **without** using spurious features. Proving this classifier also obtains good accuracy
9 requires analyzing unsupervised margin maximization, which needs stronger data assumptions[1] and is beyond the scope
10 of the paper. Unsupervised margin maximization was proposed 25 years ago (as transductive SVM [5]), but has little
11 theoretical analysis (except [1]). Instead, we focus on removing the spurious features.
12 **[R1]:** "$\gamma$ is not explained", "$x_2$ is assumed to be Gaussian... and $x_1$ is a mixture"
13 • $\gamma/\sigma$ is indeed the SNR. The condition that $w_S$ has $1 - \rho$ accuracy is equivalent to requiring that $w_1\gamma$ is large (Lemma
14 A.1 in Appendix A). $x_2$ should be $x_1$ (a typo) in line 141.
15 **[R2]:** "[theory not surprising because loss] would favor good features due to their correlation with the model", "unclear
16 how... the theory can... provide guidance to empirical discoveries"
17 • We respectfully and strongly disagree. The spurious features are indeed independent with ground-truth target labels,
18 but **no** target label is provided in self-training. The classifier is entirely self-trained on pseudolabels, which can in fact
19 be correlated with the spurious features due to biases from the source-trained initial classifier.
20 • Our theory suggests that self-training on diverse unlabeled datasets can improve model robustness, which guides
21 practitioners to collect larger and more diverse datasets even if labeling is impractical. Recent empirical results [6]
22 are consistent with our theory but do not explain why self-training works. We also show entropy minimization can
23 converge faster than pseudo-labeling (Appendix E.3) which may inform practitioners.
24 **[R3]:** "the [exponential] loss functions [is] not standard. I suggest analyzing log-loss"
25 • We would like to respectfully argue that the distinction between the two losses is unimportant in the context of our
26 paper. In theoretical works, these losses have been regarded as equivalent in various contexts [3, 4], and analyzing
27 exponential loss is convenient for the proof. In Figure 10, the two losses achieve equivalent empirical performance.
28 **[R3]:** "hard to identify which [assumptions] play the essential role in the result and which are just for technical issues",
29 "why the additional source loss on labeled data in self-training is required in the experiment on CelebA"
30 • In the setting of Theorem 3.1, we make two main assumptions: 1. the signal $x_1$ is structured and 2. the source
31 classifier is decent. Both assumptions are essential because they rule out failure cases of self-training in Section 3.1.
32 • The source loss constrains self-training to stay close to the source classifier, ensuring that the classifier correctly uses
33 signal features while entropy minimization removes the spurious ones. Previous empirical works have found the need
34 for explicit constraints to stay close to source classifier [2].
35 **[R4]:** "show that the weights for the spurious features are indeed minimized to zero", "what happens to ... accuracy if
36 the classifier gets stuck ... Will self-training reject low confidence samples?", "any effect of regularization"
37 • Figure 8 shows spurious coordinates go to 0 in the Gaussian mixture experiment. For colored MNIST, Figures 3
38 and 4 show reduced reliance on $x_2$. For our CelebA experiment, visualization of a random sample of examples whose
39 predictions are corrected by self-training shows that self-training removes the spurious feature (see attached figure).
40 • The local minima of the min-entropy objective that self-training converges to can have sub-optimal accuracy. We do
41 not reject any examples in our experiments as our theory does not analyze rejection.
42 • Regularization is necessary; without it, scaling up the source predictions can make the loss arbitrarily small. Our
43 theory analyzes l2 regularization by constraining the classifier norm to 1. We will clarify this point in the next revision.

44 **References:** [1] Derberko et al., Error bounds for transductive learning via compression and clustering, 2004. [2]
45 Shu et al., A dirt-t approach to unsupervised domain adaptation, 2018. [3] Soudry et al., The implicit bias of gradient
46 descent on separable data, 2018. [4] Telgarsky, Margins, shrinkage, and boosting, 2013. [5] Vapnik, The nature of
47 statistical learning theory, 1995. [6] Xie et al., Self-training with noisy student improves imagenet classification, 2020.

---

[1]The authors suspect that a difficult situation is when all the data reside on the hypercube $\{-1, +1\}^d$, where a large unsupervised margin classifier may exist coincidentally due to the discrete nature. We need to make careful assumptions to rule out such cases.