

**Supplementary Material for:  
Spuriousity Didn't Kill the Classifier:  
Using Invariant Predictions to Harness  
Spurious Features**

**Notes**

- For convenience, we include both the main paper and supplement here.
- The main paper is identical to the one uploaded originally.

---

# Spuriousity Didn't Kill the Classifier: Using Invariant Predictions to Harness Spurious Features

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

2 Machine learning models often fail on out-of-distribution data. To avoid this, many  
3 works have sought to extract features with a *stable* or invariant relationship with  
4 the label across domains, improving robustness by discarding the “spurious” or  
5 *unstable* features whose relationship with the label may change across domains.  
6 However, the discarded unstable features often carry *complementary* information  
7 about the label that could boost performance if used correctly in the test domain.  
8 Our main contribution is to show that it is possible to learn how to use these unstable  
9 features in the test domain *without labels*. In particular, we prove that *pseudo-*  
10 *labels* based on stable features provide sufficient guidance for doing so, provided  
11 that stable and unstable features are conditionally independent given the label.  
12 Along the way, we present a solution to the so-called “marginal problem” from  
13 probability theory, in the special case of conditionally-independent features, which  
14 may be of independent interest. Based on this theoretical insight, we propose Stable  
15 Feature Boosting (SFB), an algorithm for: (i) learning a predictor that separates  
16 stable and conditionally-independent unstable features; and (ii) using the stable-  
17 feature predictions to adapt the unstable-feature predictions in the test domain.  
18 Theoretically, we prove that SFB can learn an asymptotically-optimal predictor  
19 in the test domain without using any test-domain labels, while, empirically, we  
20 demonstrate the effectiveness of SFB on real and synthetic datasets.

## 21 1 Introduction

22 Machine learning systems can be sensitive to distribution shift [25]. Often, this sensitivity is due to a  
23 reliance on “spurious” features whose relationship with the label changes across domains, ultimately  
24 leading to degraded performance in the test domain of interest [20]. To avoid this pitfall, recent  
25 works on out-of-distribution (OOD) generalization have sought predictors which do not rely on these  
26 *spurious* or *unstable* relationships, but instead leverage relationships which are *invariant* or *stable*  
27 across multiple domains [43, 2, 34, 14]. However, despite their instability, spurious features can often  
28 provide additional or *complementary* information about the target label. Thus, if a predictor could be  
29 adjusted to use spurious features optimally in the test domain, it would boost performance substantially.  
30 That is, perhaps we don't need to discard spurious features at all, but rather use them in the right way.

31 As a very simple but illustrative example, consider the CoLoRmNIST dataset [2]. This takes the  
32 original MNIST dataset and first turns it into a binary classification task (digit in 0–4 or 5–9), and  
33 then colorizes it such that digit color (red or green) is a highly-informative but spurious feature. In  
34 particular, as depicted in Fig. 1, the two training domains are constructed such that green digits  
35 generally belong to class 0, while the test domain is constructed such that they generally belong  
36 to class 1. Finally, some label noise is added so that, across all 3 domains, digit shape correctly  
37 determines the label with probability 0.75. In previous works, the goal is to learn an invariant predictor  
38 which uses only shape and avoids using color—a spurious or unstable feature whose relationship  
39 with the label varies across domains. In this work, however, we ask the question: when and how can  
40 these such informative but spurious features be safely harnessed *without labels*? As shown in Fig. 1,  
41 this question is motivated by the fact that the invariant predictor is not Bayes-optimal in many test  
42 domains, since color information can be used to improve predictions in a domain-specific manner.

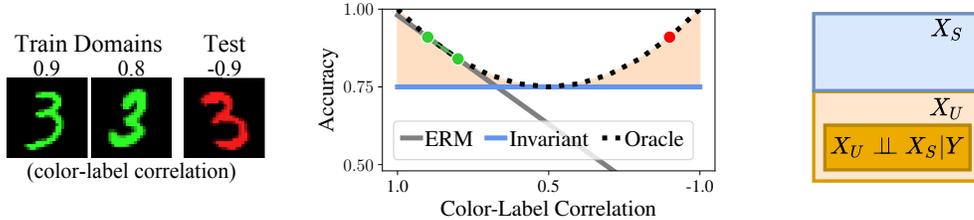


Figure 1: **Invariant (stable) and spurious (unstable) features.** (Left) Illustrative images from the ColorMNIST dataset. (Center) Performance across ColorMNIST test domains of decreasing color-label correlation for: an ERM model; an invariant model; and an oracle model using both the invariant shape *and* spurious color features optimally in the test domain. The shaded region depicts the performance boost from using the spurious feature correctly in the test domain, alongside the invariant feature. Our main contribution shows how this can be done *without labels*. (Right) Generally, invariant models use only the *stable* component  $X_S$  of  $X$ , discarding the spurious or *unstable* component  $X_U$ . We prove that predictions based on  $X_S$  can be used to harness a sub-component of  $X_U$ , highlighted in darkened orange, to reliably improve test-domain performance.

43 **Structure and contributions.** To answer this question, the remainder of this paper is organised as  
 44 follows. We first discuss related work in § 2, providing context and high-level motivation for our  
 45 proposed approach. In § 3, we then formalise the notion of stable and unstable features, showing how  
 46 unstable features can be harnessed *with* test labels, and end with a number of challenges in doing so  
 47 *without labels*. Next, in the theory of § 4, we provide concrete answers to these questions, before  
 48 using our theoretical insights to propose a Stable Feature Boosting (SFB) algorithm with guarantees  
 49 in § 5. Finally, § 6 presents our experimental results. Our main contributions can be summarised as:

- 50 • **Algorithmic:** We propose the Stable Feature Boosting (SFB) algorithm for using stable/invariant  
 51 predictions to reliably harness unstable/spurious features *without test-domain labels*. To the best of  
 52 our knowledge, SFB is the first method to do so.
- 53 • **Theoretical:** SFB is grounded in a novel theoretical result (Thm 4.4) giving sufficient conditions  
 54 under which test-domain adaptation is provably possible without labels. Under these conditions,  
 55 Thm 4.5 shows that, given enough unlabeled data, SFB learns the optimal adapted classifier.
- 56 • **Experimental:** Our experiments on synthetic and real-world data demonstrate the effectiveness  
 57 of SFB, even in practical scenarios where it is unclear if its assumptions are satisfied.

## 58 2 Related Work

59 **Domain generalization, robustness and invariant prediction.** A fundamental starting point for  
 60 work in domain generalization and robustness is the observation that certain “stable” features, often  
 61 direct causes of the label, may have an invariant relationship with the label across domains [43, 2, 58,  
 62 49, 39, 65]. However, such stable or causal predictors often discard highly-informative but unstable  
 63 information about the label. Rothenhäusler et al. [47] show that we may need to trade-off stability  
 64 and predictiveness, with the causal predictor often too conservative. Eastwood et al. [14] seek such a  
 65 trade-off via an interpretable probability-of-generalization parameter. The current work is motivated  
 66 by the idea that one might avoid such a trade-off by changing how spurious features are used at test  
 67 time, rather than discarding them at training time.

68 **Test-domain adaptation with labels.** Fine-tuning part of a model using a small number of labelled  
 69 test-domain examples is a common way to deal with distribution shift [16, 17, 13]. More recently,  
 70 it has been shown that simply retraining the last layer of an ERM-trained model outperforms more  
 71 robust feature-learning methods on spurious correlation benchmarks [46, 31, 64]. In particular, Jiang  
 72 and Veitch [30] do so when using a conditional-independence assumption not too dissimilar to ours.  
 73 However, all of these works require labels in the test domain, while we seek to adapt *without labels*.

74 **Learning with noisy labels.** An intermediate goal in our work, namely learning a model to  
 75 predict  $Y$  from  $X_U$  using pseudo-labels based on  $X_S$ , is an instance of *learning with noisy labels*,  
 76 a widely studied problem [50, 42, 8, 51, 37, 55]. Specifically, under the complementarity assumption  
 77 ( $X_S \perp\!\!\!\perp X_U | Y$ ), the accuracy of the pseudo-labels on each class is independent of  $X_U$ , placing  
 78 us in the so-called *class-conditional random noise model* [50, 42, 8]. As we discuss in Section 4,  
 79 our theoretical insights about the special structure of pseudo-labels complement existing results on  
 80 learning under this model. Our bias-correction (Eq. (4.1)) for  $P_{Y|X_U}$  is also closely related to the  
 81 “method of unbiased estimators” [42]. However, rather than correcting the loss used in ERM, our  
 82 post-hoc bias correction applies to any calibrated classifier. Moreover, our ultimate goal, learning  
 83 a predictor of  $Y$  *jointly* using  $X_S$  and  $X_U$ , is not captured by learning with noisy labels.

Table 1: **Related work.** \*QRM [14] includes a continuous hyperparameter  $\alpha \in [0, 1]$  trading off between robustness and using more information from  $X$ .

Method	Components of $X$ Used			Robust	No test-domain labels
	Stable	Complementary	All		
ERM [56]	✓	✓	✓	✗	✓
DARE [46]	✓	✓	✓	✓	✗
IRM [2]	✓	✗	✗	✓	✓
ACTIR [30]	✓	✓	✗	✓	✗
QRM [14]	✓	✓*	✓*	✓*	✓
SFB (Ours)	✓	✓	✗	✓	✓

84 **Co-training.** Our use of stable-feature pseudo-labels to train a classifier based on a disjoint subset  
 85 of (unstable) features is reminiscent of co-training [10]. Both methods benefit from conditional  
 86 independence of the two feature subsets given the label to ensure that they provide complementary  
 87 information.<sup>1</sup> The key difference is that while co-training requires (a small number of) labeled  
 88 samples from the *same distribution as the test data*, our method instead uses labeled data from  
 89 *a different distribution* (training domains), along with the assumption of a stable feature. Further  
 90 related work is discussed in Appendix H.

### 91 3 Stable and Unstable Features

92 **Setup.** We consider the problem of domain generalization (DG) [9, 40, 23] where predictors are  
 93 trained on data from multiple training domains and with the goal of performing well on data from  
 94 unseen test domains. For example, in the Camelyon17 dataset, the task is to predict if a given image of  
 95 cells contains tumor tissue, and domains correspond to the different hospitals in which the images were  
 96 captured ([5], see Fig. 4 of Appendix E). More formally, we consider datasets  $D^e = \{(X_i^e, Y_i^e)\}_{i=1}^{n_e}$   
 97 collected from  $m$  different training domains or *environments*  $\mathcal{E}_{\text{tr}} := \{E_1, \dots, E_m\}$ , with each dataset  
 98  $D^e$  containing data pairs  $(X_i^e, Y_i^e)$  sampled i.i.d. from  $\mathbb{P}(X^e, Y^e)$ .<sup>2</sup> The goal is then to learn a predictor  
 99  $f(X)$  that performs well on data from a larger set of all possible domains  $\mathcal{E}_{\text{all}} \supset \mathcal{E}_{\text{tr}}$ .

100 **Average performance: use all features.** The first approaches to DG sought predictors that perform  
 101 well *on average* over domains [9, 40] using empirical risk minimization (ERM, Vapnik 57). However,  
 102 predictors that perform well on average provably lack robustness [41], potentially performing quite  
 103 poorly on large subsets of  $\mathcal{E}_{\text{all}}$ . In particular, minimizing the average error leads predictors to make  
 104 use of any features which are informative about the label (on average), including “spurious” or  
 105 “shortcut” [20] features whose relationship with the label is subject to change across domains. In test  
 106 domains where these feature-label relationships change in new or more severe ways than observed  
 107 during training, this usually leads to significant performance drops or even complete failure [63, 6].

108 **Worst-case or robust performance: use only stable features.** To mitigate this lack of robustness,  
 109 subsequent works have sought predictors that only use *stable or invariant* features, i.e., those which  
 110 have a stable or invariant relationship with the label across domains [43, 2]. In particular, Arjovsky  
 111 et al. [2] learn features which have an invariant functional relationship with the label by enforcing that  
 112 the classifier on top of these features is optimal for all domains simultaneously. We henceforth use  
 113 *stable features* and  $X_S$  to refer to these features, and stable predictors to refer to predictors which use  
 114 only these features. Analogously, we use *unstable features*  $X_U$  to refer to features with an unstable or  
 115 changing relationship with the label across domains. Finally, note that  $X_S$  and  $X_U$  form a partition of  
 116 the components of  $X$  which are *informative about  $Y$* , as depicted in Fig. 1.

#### 117 3.1 Harnessing unstable features with labels

118 A stable predictor  $f_S(X)$  is unlikely to be the best predictor in any given domain. As depicted by the  
 119 orange regions of Fig. 1, this is because it excludes unstable features  $X_U$  which are informative about  
 120  $Y$  and can boost performance *if used correctly*. The main question we will address in the present work  
 121 is how we can harness  $X_U$  to reliably boost the performance of  $f_S(X)$  in a new domain  $e$ . To explore  
 122 this question, we assume that we are indeed able to learn a stable predictor using prior methods, e.g.,  
 123 IRM [2], and, for now, that we have access to labelled examples in this new domain which can be  
 124 used to update or re-learn the domain-specific relation between  $X_U$  and  $Y$ .

125 **Boosting the stable predictor.** To begin, note that we need only update the  $X_U$ - $Y$  relation since,  
 126 by definition, the  $X_S$ - $Y$  relation is stable across domains. We will thus seek a feature space which

<sup>1</sup>See Krogel and Scheffer [33] and Theorem 1 of Blum and Mitchell [10] for discussion of this assumption.

<sup>2</sup>We drop the domain superscript  $e$  when referring to random variables from any environment.

127 separates  $X_S$  and  $X_U$ , allowing only the unstable  $X_U$ - $Y$  relation to be updated. To do so, let us first  
 128 decompose a predictor  $f$  into a feature representation  $\Phi$  and classifier  $h$ , with  $f = h \circ \Phi$ , and then  
 129 describe the boosted joint predictor  $f^e(X)$  in domain  $e$  as:

$$f^e(X) = f_S(X) + f^e(X) = h_S(\Phi_S(X)) + h^e(\Phi_U(X)) \quad (3.1)$$

$$= h_S(X_S) + h^e(X_U). \quad (3.2)$$

130 Here, both  $f_S(X)$  and  $f^e(X)$  produce logits, meaning that the unstable predictor  $f^e(X)$  essentially  
 131 adds a domain-specific adjustment to the stable predictor  $f^s(X)$  in logit space. As illustrated by  
 132 Eqs. (3.1) and (3.2), the role of  $\Phi_S$  and  $\Phi_U$  is to extract  $X_S$  and  $X_U$ , respectively, from the observed  
 133 features  $X$ . Note that the stable predictor  $f_S$  and classifier  $h_S$ , as well as the feature extractors  $\Phi_S$   
 134 and  $\Phi_U$  are shared across domains  $e$ , whereas the unstable classifier  $h_U^e$  is not. In principle,  $h_U^e$  could  
 135 take any form, so long as we have enough labelled examples to learn it. In practice, however, we  
 136 generally take  $h_U^e$  to be a linear classifier for sample efficiency.

137 **Adapting  $h_U^e$  with labels.** Given a new domain  $e$  with labelled examples, we can boost the perfor-  
 138 mance of our stable predictor by adapting  $h_U^e$  to minimize the joint-predictor loss. Specifically, letting  
 139  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  be a loss function (e.g., cross-entropy) and  $R_e(f) = \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) | E = e]$   
 140 the statistical risk of a predictor  $f : \mathcal{X} \rightarrow \mathcal{Y}$  in domain  $e$ , we can adapt  $h_U^e$  to solve:

$$\min_{h_U} \sum_{e \in \mathcal{E}_U} R_e(\sigma \circ ((h_S \circ \Phi_S) + (h_U \circ \Phi_U))) \quad (3.3)$$

141 Note that Jiang and Veitch [30, Eq. 2.1] proposed a similar joint predictor for using *labelled* test-  
 142 domain examples to update a domain-specific component. However, they do not explicitly separate  
 143 stable and unstable features  $X_S$  and  $X_U$ , which will later prove crucial for our approach *without labels*.

### 144 3.2 Harnessing unstable features *without labels*

145 The previous section made clear how we can safely harness  $X_U$  when we have test-domain labels. We  
 146 now consider the main question of this work—*can we safely harness  $X_U$  without test-domain labels?*  
 147 More specifically, how can we update the unstable classifier  $h_U^e$  to capture the new  $X_U$ - $Y$  relation  
 148 given only unlabelled test-domain examples  $\{X_i^e\}_{i=1}^{n_e}$ ? We could, of course, simply select a *fixed*  
 149 unstable classifier  $h_U^e$  by relying solely on the training domains (e.g., by minimizing average error),  
 150 and hope that this works for the test-domain  $X_U$ - $Y$  relation. However, by definition of  $X_U$  being  
 151 unstable, this is clearly not a robust or reliable approach—the focus of our efforts, as illustrated in  
 152 Table 1. As in § 3.1, we assume that we are able to learn a stable predictor  $f_S$  using prior methods [2].

153 **From stable predictions to robust pseudo-labels.** While we do not have labels in the test domain,  
 154 we *do* have stable predictions. By definition, such predictions are imperfect (i.e., *noisy*) but robust,  
 155 and can be used to form *pseudo-labels*  $\hat{Y}_i = \arg \max_j f_S(X_i)_j$ , with  $f_S(X_i)_j$  denoting the  $j^{\text{th}}$  logit  
 156 of the stable prediction for  $X_i$ . Can we somehow use these noisy but robust pseudo-labels to guide  
 157 our updating of  $h_U^e$ , and, ultimately, our use of  $X_U$  in the test domain?

158 **From joint to unstable-only risk.** Unfortunately, if we try to use our robust pseudo-labels as if they  
 159 were true labels—updating  $h_U^e$  to minimize the joint risk as in Eq. (3.3)—we get a trivial solution  
 160 of  $h_U^e(\cdot) = 0$ . If our loss  $\ell$  is accuracy, this trivial solution is clear since  $h_U^e(\cdot) = 0$  achieves 100%  
 161 accuracy. For cross-entropy, the same trivial solution exists, as we show in Prop. D.1 of Appendix D.  
 162 Thus, we cannot minimize a joint loss involving  $f_S$ 's predictions when using  $f_S$ 's pseudo-labels.  
 163 Instead, we must consider updating  $h_U^e$  to minimize the unstable-only risk  $R_e(\sigma \circ h_S \circ \Phi_S)$ .

164 **More questions than answers.** While this new procedure *could* work, it raises many questions about  
 165 *when* it will work, or, more precisely, the conditions under which it can be used to safely harness  $X_U$ .  
 166 We now summarise these questions before addressing them in the next section (§ 4):

- 167 1. **Does it make sense to minimize the unstable-only risk?** In particular, when can we minimize  
 168 the unstable-predictor risk *alone* or separately, and then arrive at the optimal joint predictor? This  
 169 cannot always work; e.g., for independent  $X_S, X_U \sim \text{Bernoulli}(1/2)$  and  $Y = X_S \text{ XOR } X_U$ ,  $Y$   
 170 is independent of each of  $X_S$  and  $X_U$  and hence cannot be predicted from either alone.
- 171 2. **Can we just add the logits as before?** Building on question 1, if we separately optimize the  
 172 predictions of the unstable classifier  $h_U^e$  using the pseudo-labels  $\hat{Y}$ , does it make sense to simply  
 173 add the logits afterwards as in Eq. (3.2)? Intuitively, simply adding the stable and unstable logits as  
 174 before would require them both to be “of the same scale”, or, more precisely, properly calibrated.  
 175 Do we have any reason to believe that, after training on  $h_S$ 's pseudo-labels,  $h_U^e$  will properly  
 176 calibrated and thus can be integrated with  $h_S$  as in Eq. (3.2)?

177 3. **Can the student outperform the teacher?** Stable predictions likely make mistakes—indeed, this  
 178 is the motivation for trying to improve them. Is it possible to correct these mistakes with unstable  
 179 features, thus improving performance? In particular, is it possible to learn an unstable “student”  
 180 predictor that outperforms its own supervision signal or teacher? Perhaps surprisingly, we show  
 181 that, for certain types of features, the answer is yes; in fact, even a very weak stable predictor,  
 182 with performance just above chance, can be used to learn an *optimal* unstable classifier in the test  
 183 domain given enough unlabeled data.

#### 184 4 Theory: When can we safely harness unstable features without labels?

185 Suppose we have already identified a stable feature  $X_S$  and a potentially unstable feature  $X_U$  (we  
 186 will return to the question of how to learn  $X_S$  and  $X_U$  themselves in Section 5, after identifying the  
 187 additional conditions we would like  $X_S$  and  $X_U$  to satisfy). In this section, we analyze the problem of  
 188 using  $X_S$  to leverage  $X_U$  without labels in the test domain. We first reduce this to a special case of the  
 189 so-called “marginal problem” in probability theory, i.e., the problem of identifying a joint distribution  
 190 based on information about its marginals. In the special case where two variables are conditionally  
 191 independent given a third, we show this problem can be solved exactly; this solution, which may be  
 192 of interest beyond the context of domain adaptation, motivates our test-domain adaptation algorithm  
 193 (Algorithm 1), presented in Section 5, and forms the basis of Theorem 4.5 showing that Algorithm 1  
 194 converges to the best possible classifier given enough unlabeled data.

195 To formalize our assumptions, we first pose a population-level model of our domain generalization  
 196 setup. Let  $E$  be a random variable denoting the *environment*. Given an environment  $E$ , the stable  
 197 feature  $X_S$ , the unstable feature  $X_U$ , and the label  $Y$  are distributed according to  $P_{X_S, X_U, Y|E}$ . Given,  
 198 this we can formalize the three key assumptions underlying our approach.

199 We first formalize the notion of a stable feature, motivated in the previous section:

200 **Definition 4.1** (Stable and Unstable Predictors).  $X_S$  is a stable predictor of  $Y$  if  $P_{Y|X_S}$  does not depend  
 201 on  $E$ ; equivalently, if  $Y$  and  $E$  are conditionally independent given  $X_S$  ( $Y \perp\!\!\!\perp E|X_S$ ). Conversely,  
 202  $X_U$  is an unstable predictor of  $Y$  if  $P_{Y|X_U}$  depends on  $E$ ; equivalently, if  $Y$  and  $E$  are conditionally  
 203 dependent given  $X_U$  ( $Y \not\perp\!\!\!\perp E|X_U$ ).

204 Next, we state our complementarity assumption, which we will show is key to justifying the approach  
 205 of separately learning the relationships  $X_S$ - $Y$  and  $X_U$ - $Y$  and then combining them:

206 **Definition 4.2** (Complementary Features).  $X_S$  and  $X_U$  are complementary predictors of  $Y$  if  $X_S \perp\!\!\!\perp$   
 207  $X_U|(Y, E)$ ; i.e., if  $X_S$  and  $X_U$  contain no redundant information beyond that contained in  $Y$  and  $E$ .

208 Finally, it is fairly intuitive that, to provide a useful signal for test-domain adaptation, the stable  
 209 feature needs to be predictive of the label in the test domain. Formally, we assume

210 **Definition 4.3** (Informative Stable Predictor).  $X_S$  is said to be informative of  $Y$  in environment  $E$  if  
 211  $X \not\perp\!\!\!\perp Y|E$  (i.e.,  $X_S$  is predictive of  $Y$  within the environment  $E$ ).

212 We will discuss the roles of these assumptions, and how they relate to the motivating questions  
 213 at the end of Section 3.2, in greater detail after stating the main result (Theorem 4.4) that utilizes  
 214 them. Note that, to keep our results as general as possible, we avoid assuming a particular causal  
 215 generative model underlying data. However, the conditional (in)dependence assumptions above can  
 216 be interpreted as constraints on such a causal model, and, in Appendix D.1, we formally characterize  
 217 the set of causal generative models that are consistent with our assumptions. Notably, we show that  
 218 our setting generalizes those of several existing works assuming specific causal generative models or  
 219 constraints on possible distribution shifts [45, 59, 30].

220 **Reduction to the marginal problem with complementary features.** Since we assumed the feature  
 221  $X_S$  is stable,  $P_{Y|X_S, E} = P_{Y|X_S}$  in the test domain is the same as in the training domains. Hence, let us  
 222 suppose we have used the training data to learn this relationship, and hence know  $P_{Y|X_S}$ . Suppose  
 223 also that we have enough unlabeled test-domain data to learn  $P_{X_S, X_U|E}$  in the test environment  $E$ .

224 Recall that our goal in test-domain adaptation is to predict  $Y$  from  $(X_S, X_U)$  in the test domain  $E$ .  
 225 The remainder of our discussion will take place entirely conditioned on  $E$  being the test domain, and  
 226 hence we will omit this dependence from the notation. If we could express  $P_{Y|X_S, X_U}$  in terms of  
 227  $P_{Y|X_S}$  and  $P_{X_S, X_U}$ , we could then use  $P_{Y|X_S, X_U}$  to optimally predict  $Y$  from  $(X_S, X_U)$ . Thus, our task

228 thus becomes to reconstruct  $P_{Y|X_S, X_U}$  from  $P_{Y|X_S}$  and  $P_{X_S, X_U}$ . This is an instance of the classical  
 229 “marginal problem” from probability theory [27, 28, 18], which asks under which conditions we can  
 230 recover the joint distribution of a set of random variables given information about its marginals. In  
 231 general, although one can place bounds on the conditional distribution  $P_{Y|X_U}$ , it cannot be completely  
 232 inferred from  $P_{Y|X_S}$  and  $P_{X_S, X_U}$  [18]. However, the following section demonstrates that, *under the*  
 233 *additional assumptions that  $X_S$  and  $X_U$  are complementary and  $X_S$  is informative*, we can exactly  
 234 recover  $P_{Y|X_S, X_U}$  from  $P_{Y|X_S}$  and  $P_{X_S, X_U}$ .

#### 235 4.1 Solving the marginal problem with complementary features

236 To simplify notation, suppose the label  $Y$  is binary, taking values in  $\{0, 1\}$ ; the multiclass extension  
 237 is detailed in Appendix C. The following result then shows how to reconstruct  $P_{Y|X_S, X_U}$  from  $P_{Y|X_S}$   
 238 and  $P_{X_S, X_U}$  when  $X_S$  and  $X_U$  are complementary and  $X_S$  is informative.

239 **Theorem 4.4** (Solution to the marginal problem with binary labels and complementary features).  
 240 Consider three random variables  $X_1$ ,  $X_2$ , and  $Y$ , where (i)  $Y$  is binary ( $\{0, 1\}$ -valued), (ii)  $X_1$  and  
 241  $X_2$  are complementary features for  $Y$  (i.e.,  $X_1 \perp\!\!\!\perp X_2|Y$ ), and (iii)  $X_1$  is informative of  $Y$  ( $X_1 \not\perp\!\!\!\perp Y$ ).  
 242 Then, the joint distribution of  $(X_1, X_2, Y)$  can be written in terms of the joint distributions of  $(X_1, Y)$   
 243 and  $(X_1, X_2)$ . Specifically, suppose  $\hat{Y}|X_S \sim \text{Bernoulli}(\Pr[Y = 1|X_S])$  is a pseudo-label<sup>3</sup>, and  
 244  $\epsilon_0 := \Pr[\hat{Y} = 0|Y = 0]$  and are the conditional probabilities that  $\hat{Y}$  and  $Y$  agree, given  $Y = 0$  and  
 245  $Y = 1$ , respectively. Then, we have  $\epsilon_0 + \epsilon_1 > 1$ ,

$$\Pr[Y = 1|X_2] = \frac{\Pr[\hat{Y} = 1|X_2] + \epsilon_0 - 1}{\epsilon_0 + \epsilon_1 - 1}, \quad \text{and} \quad (4.1)$$

$$246 \Pr[Y = 1|X_1, X_2] = \sigma(\text{logit}(\Pr[Y = 1|X_1]) + \text{logit}(\Pr[Y = 1|X_2]) - \text{logit}(\Pr[Y = 1])). \quad (4.2)$$

247 Intuitively, suppose we train a model to predict a pseudo-label  $\hat{Y}$ , generated based on feature  $X_1$ ,  
 248 from a feature  $X_2$ . Assuming  $X_1$  and  $X_2$  are complementary, Eq (4.1) shows how to transform this  
 249 into a prediction of the true label  $Y$ , correcting for biases caused by possible disagreement between  $\hat{Y}$   
 250 and  $Y$ . Meanwhile, Eq. (4.2) shows how to integrate predictors based on  $X_1$  and  $X_2$  while accounting  
 251 for redundancy in the two predictors.  
 252

253 **The role of complementarity.** The assumption that  $X_1$  and  $X_2$  are complementary plays two separate  
 254 but equally crucial roles in Theorem 4.4. First, if  $X_1$  and  $X_2$  only share information about  $Y$ , then,  
 255 when we train a model to predict  $\hat{Y}$  (which depends only on  $X_1$ ) from  $X_2$ , the model will only learn  
 256 to predict information about  $Y$  (rather than other relationships between  $X_1$  and  $X_2$ ). This insight is  
 257 key to justifying the bias-correction formula (Eq. (4.1)). Second, by ensuring that the only interaction  
 258 between  $X_1$  and  $X_2$  is due to  $Y$  itself, complementarity implies that  $P_{Y|X_1, X_2}$  is decomposable into  
 259  $P_{Y|X_1}$  and  $P_{Y|X_2}$ . Specifically, one can simply add estimates of  $P_{Y|X_1}$  and  $P_{Y|X_2}$  in logit-space while  
 260 subtracting a correction term based on the marginal distribution of  $Y$  (see Eq. (4.2)).

261 **The role of informativeness.** It is intuitive that informativeness ( $X_1 \not\perp\!\!\!\perp Y$ ) is necessary; for the  
 262 pseudo-labels to be useful,  $X_1$  must help predict  $Y$ . More surprisingly, informativeness is *sufficient*  
 263 for Theorem 4.4, i.e., *any* dependence between  $X_1$  and  $Y$  allows us to fully learn the relationship  
 264 between  $X_2$  and  $Y$ . This gives an affirmative answer to our question, *Can the student outperform*  
 265 *the teacher?*, from Section 3.2. This is not to say that a strong relationship between  $X_1$  and  $Y$  is not  
 266 helpful; while informativeness is equivalent to  $\epsilon_0 + \epsilon_1 > 1$  (see Lemma A.2 in Appendix A.1), a  
 267 weak relationship corresponds to  $\epsilon_0 + \epsilon_1 \approx 1$ , making the bias-correction 4.1 unstable. Notably, this  
 268 only affects the (unlabeled) *sample complexity* of learning  $P_{Y|X_2}$ , not *consistency* (Theorem 4.5).

269 Appendix A.2 provides further discussion of Theorem 4.4, including its relationship with existing  
 270 work on learning from noisy labels and possible applications beyond domain adaptation.

#### 271 4.2 A provably consistent algorithm for unsupervised test-domain adaptation

272 To see why Theorem 4.4 is useful for test-domain adaptation, observe that stability of  $X_1$  implies that  
 273 the conditional distribution  $P_{Y|X_1}$  is the same in the training and test domains. Hence,  $P_{Y|X_1}$  can be  
 274 learned using labeled data. Meanwhile, the joint distribution  $P_{X_1, X_2}$  in the test domain can be learned  
 275 using only *unlabeled* test-domain data. Theorem 4.4 thus implies that we can learn  $P_{Y|X_1, X_2}$  in the test  
 276 domain using only labeled data from the training domains and unlabeled data from the test domain.

<sup>3</sup>Though discrete, our *stochastic* pseudo-labels differ from hard ( $\hat{Y} = 1\{\Pr[Y = 1|X_S] > 1/2\}$ ) or soft  
 pseudo-labels often used in practice [19, 35, 48]. By capturing irreducible error in  $Y$ , stochastic pseudo-labels  
 ensure  $\Pr[Y|X_2]$  is well-calibrated, allowing us to combine  $\Pr[Y|X_1]$  and  $\Pr[Y|X_2]$  in Eq. (4.2).

---

**Algorithm 1:** Bias-corrected unsupervised domain adaptation procedure.

---

**Input:** Regression function  $\eta_S(x_S) = \Pr[Y = 1|X_S = x_S]$ , subroutine regressor,  $n$  unlabeled samples  $\{(X_{S,i}, X_{U,i})\}_{i=1}^n$  from the test domain  
**Output:** Estimate  $\hat{\eta}_n : \mathcal{X}_S \times \mathcal{X}_U \rightarrow [0, 1]$  of  $\Pr[Y = 1|X_S = x_S, X_U = x_U]$

```
1 for  $i \in [n]$  do // generate pseudolabels
2 | Sample  $\hat{Y}_i \sim \text{Bernoulli}(\eta_S(X_{S,i}))$ 
3  $\hat{\eta}_{U,n} \leftarrow \text{regressor}(\{(X_{U,i}, \hat{Y}_i)\}_{i=1}^n)$ 
4  $n_1 \leftarrow \sum_{i=1}^n \hat{Y}_i$ ;  $\hat{\beta}_{1,n} \leftarrow \text{logit}(\frac{n_1}{n})$ 
5  $(\hat{\epsilon}_{0,n}, \hat{\epsilon}_{1,n}) \leftarrow (\frac{1}{n-n_1} \sum_{i=1}^n (1 - \hat{Y}_i)(1 - \eta_S(X_{S,i})), \frac{1}{n_1} \sum_{i=1}^n \hat{Y}_i \eta_S(X_{S,i}))$ 
6 return  $(\hat{\eta}_n(x_S, x_U) \mapsto \sigma(\text{logit}(\eta_S(x_S)) + \text{logit}(\frac{\min\{\hat{\epsilon}_{1,n}, \max\{1 - \hat{\epsilon}_{0,n}, \hat{\eta}_{U,n}(x_U)\}\} + \hat{\epsilon}_{0,n} - 1}{\hat{\epsilon}_{0,n} + \hat{\epsilon}_{1,n} - 1}) - \hat{\beta}_{1,n}))$ 
```

---

277 Based on this reasoning, Alg. 1 presents our proposed unsupervised test-domain adaptation method.  
278 Intuitively, given a stable soft-classifier  $\eta_S$ , Algorithm 1 simply implements a finite-sample version  
279 of the bias-correction and combination equations (Eqs. (4.1) and 4.2) in Theorem 4.4. Algorithm 1  
280 also comes with the following guarantee:

281 **Theorem 4.5** (Consistency Guarantee, Informal). *Assume (i)  $X_S$  is stable, (ii)  $X_S$  and  $X_U$  are*  
282 *complementary, and (iii)  $X_S$  is informative of  $Y$  in the test domain. If  $\hat{\eta}_{U,n} \rightarrow \Pr[\hat{Y} = 1|X_U]$  as*  
283  *$n \rightarrow \infty$ , then  $\hat{\eta}_n \rightarrow \Pr[Y = 1|X_S, X_U]$ .*

284 In words, as the amount of unlabeled data from the test domain increases, if the regressor on Line 3  
285 of Algorithm 1 is able to learn to predict the pseudo-label  $\hat{Y}$ , then the test-domain classifier output by  
286 Algorithm 1 will learn to predict the true label  $Y$  in the test domain. Convergence in Theorem 4.5  
287 occurs  $P_{X_S, X_U}$ -almost everywhere, both weakly (in prob.) and strongly (a.s.), depending on the mode  
288 of convergence of  $\hat{\eta}_{U,n}$ . Due to space constraints, formal statements and proofs are in Appendix B.

## 289 5 Algorithm: Stable Feature Boosting

290 We now use our theoretical insights from § 4 to pick up where we left off in § 3.2, ultimately arriving  
291 at a practical algorithm for harnessing unstable features without labels. We start by describing the  
292 training-domain algorithm, where our goal is to learn stable and complementary features, and then  
293 describe the test-domain adaptation algorithm, where our goal is to correctly adapt the unstable  
294 classifier  $h_U^e$  using the stable predictions (or pseudo-labels).

295 **Recap and learning goals.** In Eq. (3.1) of § 3.1 we described a joint predictor  $f^e(X) = f_S(X) +$   
296  $f_U^e(X)$  which can reliably boost the performance of the unstable predictor  $f_S$ —so long as we have  
297 labels in the test domain to update the unstable or domain-specific classifier  $h_U^e$ . In § 3.2, we ran into  
298 some problems when trying to update  $h_U^e$  without labels, and ended the section with a number of  
299 questions about when it’s possible to use the stable predictions of  $f_S$  to update  $h_U^e$ . In § 4, we provided  
300 concrete answers to these questions, proving that informativeness ( $f_S$  carries some information about  
301  $Y$ ) and complementarity (the stable and unstable features are conditionally independent given  $Y$ )  
302 suffice for learning the optimal  $h_U^e$  from  $f_S$ ’s predictions (asymptotically). Moreover, § 4 showed  
303 that, if we can indeed learn informative stable features  $X_S$  and complementary features  $X_C$ , then we  
304 can employ the bias-corrected adaptation algorithm of Alg. 1 (or Alg. 2 for the multi-class case) to  
305 update  $h_U^e$ . Thus, our training-time goal is now to extract  $X_S$  and  $X_C$  from the observed  $X$ , such that  
306 we can harness  $X_C$  in the test domain. More specifically, we have the following learning goals:

- 307 1.  $f_S(X)$  is a stable, well-calibrated predictor with good performance.<sup>4</sup>
- 308 2. In a given domain  $e$ ,  $f_U^e(X)$  boosts the performance of  $f_S(X)$  using complementary features.

---

<sup>4</sup>While Theorem 4.4 only assumes the stable feature is informative, as discussed in Section 4.1, a more accurate stable predictor improves sample efficiency of SFB.

Table 2: OOD accuracies. Mean and standard errors are over 100, 5, 5 seeds (Synthetic, Camelyon17, PACS).

Algorithm	Synthetic	Camelyon17	PACS			
	-	-	P	A	C	S
ERM	9.9 ± 0.1	90.2 ± 1.1	93.0 ± 0.7	79.3 ± 0.5	74.3 ± 0.7	65.4 ± 1.5
IRM	74.9 ± 0.1	90.2 ± 1.1	93.3 ± 0.3	78.7 ± 0.7	75.4 ± 1.5	65.6 ± 2.5
ACTIR	74.8 ± 0.4	77.7 ± 1.7 <sup>†</sup>	94.8 ± 0.1	<b>82.5 ± 0.4</b>	<b>76.6 ± 0.6</b>	62.1 ± 1.3
SFB w/o adapt	74.7 ± 1.2	89.8 ± 1.2	93.7 ± 0.6	78.1 ± 1.1	73.7 ± 0.6	69.7 ± 2.3
SFB w. adapt	<b>89.2 ± 2.9</b>	<b>90.3 ± 0.7</b>	<b>95.8 ± 0.6</b>	80.4 ± 1.3	<b>76.6 ± 0.6</b>	<b>71.8 ± 2.0</b>

309 **Objective function.** To achieve the above learning goals, we propose the following objective:

$$\min_{\Phi, h_S, h_U^e} \sum_{e \in \mathcal{E}_r} R_e(\sigma \circ h_S \circ \Phi_S) + R_e(\sigma \circ ((h_S \circ \Phi_S) + (h_U^e \circ \Phi_U))) \quad (5.1)$$

$$+ \lambda_S \cdot P_{\text{Stab}}(\Phi_S, h_S) + \lambda_C \cdot P_{\text{Comp}}(\Phi_S(X), \Phi_U(X)) \quad (5.2)$$

310 Here,  $P_{\text{Stability}}(\Phi_S, h_S)$  is a penalty encouraging stability of  $\Phi_S(X)$  (i.e.,  $Y \perp\!\!\!\perp E | \Phi_S(X)$ ), while  
 311  $P_{\text{Comp}}(\Phi_S(x_i), \Phi_U(x_i))$  is a penalty encouraging complementarity of  $\Phi_S(X)$  and  $\Phi_U(X)$  (i.e.,  
 312  $\Phi_S(X) \perp\!\!\!\perp \Phi_U(X) | Y$ ). Several approaches have been proposed for enforcing stability [43, 2, 15,  
 313 47, 58, 39, 65], e.g., IRM [2], while complementarity can be enforced by a generic conditional-  
 314 dependence penalty, e.g., the conditional Hilbert-Schmidt Independence Criterion [21, HSIC] or  
 315 cheaper proxy methods like that of Jiang and Veitch [30, §3.1].  $\lambda_S \in [0, \infty)$  and  $\lambda_C \in [0, \infty)$  are  
 316 regularization hyperparameters. In principle, an additional hyperparameter  $\gamma \in [0, 1]$  could control  
 317 the relative weighting of stable and joint risks, i.e.,  $\gamma R_e(h_S \circ \Phi_S)$  and  $(1 - \gamma) R_e((h_S \circ \Phi_S) + (h_U \circ$   
 318  $\Phi_U))$ . However, in practice, we found this to be unnecessary.

319 **Post-hoc calibration.** Finally, as discussed in Section 4.2, correctly combining the stable and unstable  
 320 predictions at adaptation time requires them to be properly calibrated. Thus, after optimizing the  
 321 objective (5.2), we also suggest applying a standard post-processing step that improves the calibration  
 322 of the stable classifier  $h_S \circ \Phi_S$ , e.g., simple temperature scaling [24].

323 **Adapting without labels.** Armed with a stable predictor  $f_S = h_S \circ \Phi_S$  and complementary features  
 324  $\Phi_U(X)$ , our goal is now to adapt the unstable classifier  $h_U^e$  in the test domain to safely harness (or  
 325 make optimal use of)  $\Phi_U(X)$ . To do so, we’ll make use of the bias-corrected adaptation algorithm  
 326 of Alg. 1 (or Alg. 2 for the multi-class case) which takes as input the stable classifier  $h_S$  and  
 327 unlabelled test-domain dataset  $\{\Phi_S(x_i), \Phi_U(x_i)\}_{i=1}^{n_e^T}$ . This adaptation procedure returns the adapted  
 328 joint classifier  $\hat{f}^{eT}$  (the logit of  $\hat{\eta}_n$  in Line 6 of Alg 1) finally used for prediction in the test domain.

## 329 6 Experiments

330 We now evaluate the performance of our algorithm on synthetic and real-world datasets requiring out-  
 331 of-distribution generalization. Fig. 4 depicts samples from the datasets considered, while Appendix G  
 332 gives further experimental details. Code will be made available upon acceptance.

333 **Synthetic dataset.** We first consider an anti-causal synthetic dataset based on that of [30, §6.1]  
 334 where data is generated according to the following structural equations:  $Y \leftarrow \text{Rad}(0.5)$ ,  $X_S \leftarrow$   
 335  $Y \cdot \text{Rad}(0.75)$ , and  $X_U \leftarrow Y \cdot \text{Rad}(\beta^e)$ , where the input  $X = (X_S, X_U)$  and  $\text{Rad}(\beta)$  means that a  
 336 random variable is  $-1$  with probability  $1 - \beta$  and  $+1$  with probability  $\beta$ . Following [30, §6.1], we  
 337 create two training domains with  $\beta_e \in \{0.95, 0.7\}$ , one validation domain with  $\beta_e = 0.6$  and one  
 338 test domain with  $\beta_e = 0.1$ . The idea here is that, during training, prediction based on the stable  $X_S$   
 339 results in lower accuracy (75%) than prediction based on the unstable  $X_U$  (82.5%). Thus, models  
 340 optimizing for prediction accuracy only—and not stability—will use  $X_U$  and ultimately end up with  
 341 only 10% in the test domain. Importantly, while the stable predictor achieves 75% accuracy in the  
 342 test domain, performance can be improved to 90% if  $X_U$  can be used correctly.

343 Following [30], we use a simple 3-layer network and choose hyperparameters using the validation-  
 344 domain performance: see Appendix G for further details. As shown in Table 2, ERM performs poorly  
 345 as it uses the unstable feature  $X_S$ , while IRM [2], ACTIR [30] and our SFB algorithm all do well by us-  
 346 ing only the stable feature  $X_S$ . Critically, only our SFB is capable of harnessing  $X_U$  in the test domain  
 347 *without labels*, leading to a near-optimal boost in performance. In Appendix F.1, we also consider  
 348 a synthetic dataset where our conditional independence assumption  $X_S \perp\!\!\!\perp X_U | Y$  does not hold.

349 **ColorMNIST.** We now consider the ColorMNIST dataset of Arjovsky et al. [2], described in § 1 and  
 350 depicted in Fig. 1 (left). Experimentally, we follow the setup of Eastwood et al. [14, §6.1], including  
 351 a simple 3-layer network: see Appendix G for further implementation details.

Table 3: CMNIST test accuracies.

Algorithm	Test Acc.
ERM	27.9 ± 1.5
IRM	69.7 ± 0.9
V-REx	71.6 ± 0.5
EQRm	71.4 ± 0.4
SFB (Ours) w/o adapt.	70.6 ± 1.8
SFB (Ours) w. adapt.	88.1 ± 1.8
Oracle w/o adapt.	72.1 ± 0.7
Oracle w. adapt.	89.9 ± 0.1

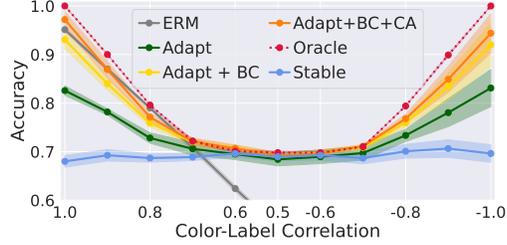


Figure 2: CMNIST results. Oracle: ERM trained on labelled test-domain data. All other curves (but ERM) refer to our algorithm. ‘Stable’: unadapted, ‘BC’: bias-corrected, and ‘CA’: calibrated.

352 Table 3 shows that: (i) SFB learns a stable predictor with performance comparable to other invariant-  
 353 prediction methods like IRM [2], V-REx [34] and EQRm [14]; and (ii) only SFB is capable of  
 354 harnessing the spurious color features in the test domain *without labels*, leading to a near-optimal  
 355 boost in performance. Note that “Oracle w/o adapt.” refers to an ERM model trained on grayscale  
 356 images, while “Oracle w. adapt” refers to an ERM model trained on labelled test-domain data. In  
 357 addition, Fig. 2 shows that: (i) both bias-correction (BC) and post-hoc calibration (CA) improve  
 358 adaptation performance; and (ii) without labels, our SFB algorithm can harness the spurious color  
 359 feature near-optimally in test domains of varying color-label correlation—the original goal we set out  
 360 to achieve, depicted in Fig. 1. Further results and ablations are provided in Appendix F.2.

361 **PACS.** We now consider PACS [36]—a 7-class image-classification dataset consisting of 4 domains:  
 362 photos (P), art (A), cartoons (C) and sketches (S), with examples shown in Fig. 4 of Appendix E.  
 363 For each domain, we test model performances after training on the other three domains. Following  
 364 [23, 30], we choose hyperparameters using leave-one-domain-out cross-validation.

365 Table 2 shows that our SFB algorithm’s stable (i.e., without-adaptation) performance is comparable  
 366 to that of the other invariance-seeking methods: IRM and ACTIR. One exception is the sketch  
 367 domain (S), the most severe shift based on performance drop, where our stable predictor performs  
 368 best. Another exception is that ACTIR’s stable predictor performs better on domains A and C.  
 369 Most notable, however, is: (i) the consistent boost in performance that SFB gets from unsupervised  
 370 adaptation; and (ii) the fact that SFB performs best or joint-best on 3 of the 4 domains. Together,  
 371 these results indicate that SFB can be useful on real-world datasets where it is unclear whether or not  
 372 our conditional-independence assumption holds.

373 **Camelyon17.** Finally, we consider the Camelyon17 [5] dataset from the WILDS benchmark [32], a  
 374 medical dataset with histopathology images from 5 hospitals which use different staining and imaging  
 375 techniques (see Fig. 4 of Appendix E). The goal is to determine whether or not a given image contains  
 376 tumour tissue, making it a binary classification task. We follow the train-validation-test split of  
 377 WILDS, using 3 domains for training and 1 each for validation and testing. Following Jiang and Veitch  
 378 [30], we use an ImageNet-pretrained ResNet18. See Appendix G.3 for further implementation details.

379 Table 2 shows mixed results. On the one hand, adapting gives SFB a small performance boost and  
 380 reduces the variance across random seeds. On the other hand, the adapted performance is on par  
 381 with both IRM and the simpler ERM method. In line with [23], we found that a properly-tuned ERM  
 382 model can be difficult to beat on real-world datasets, particularly when they don’t contain severe  
 383 distribution shift. While we conducted this proper tuning for ERM, IRM and SFB (see Appendix G.3),  
 384 doing so for ACTIR was non-trivial. We thus report the result from their paper [30, Tab. 1], which is  
 385 likely lower due to hyperparameter selection (they report  $\approx 70\%$  accuracy for ERM and IRM).

## 386 7 Discussion

387 This work demonstrated, both theoretically and practically, how to adapt spurious but informative  
 388 features to new test domains using only a stable, complementary training signal. Our proposed Stable  
 389 Feature Boosting algorithm can provide significant performance gains compared to only using stable  
 390 features or using unadapted spurious features, without requiring any true labels in the test domain. In  
 391 theory, the most significant limitation of SFB is its assumption of complementarity (i.e., conditional  
 392 independence of spurious features and stable features, given the label). Importantly, our experimental  
 393 results suggest that SFB may be robust to violations of complementarity in practice; on real-world  
 394 datasets such as PACS or Camelyon17, where there is no reason to believe complementarity holds,  
 395 SFB performs at least as well or better than unadapted methods such as ERM and IRM.

396 **References**

- 397 [1] Abney, S. (2002). Bootstrapping. In *Proceedings of the 40th annual meeting of the Association*  
398 *for Computational Linguistics*, pages 360–367. [Cited on page 28.]
- 399 [2] Arjovsky, M., Bottou, L., Gulrajani, I., and Lopez-Paz, D. (2020). Invariant risk minimization.  
400 arXiv:1907.02893. [Cited on pages 1, 2, 3, 4, 8, 9, 23, 24, and 27.]
- 401 [3] Ba, J. and Caruana, R. (2014). Do deep nets really need to be deep? *Advances in neural*  
402 *information processing systems*, 27. [Cited on page 17.]
- 403 [4] Balcan, M.-F., Blum, A., and Yang, K. (2004). Co-training and expansion: Towards bridging  
404 theory and practice. *Advances in neural information processing systems*, 17. [Cited on page 28.]
- 405 [5] Bandi, P., Geessink, O., Manson, Q., Van Dijk, M., Balkenhol, M., Hermsen, M., Bejnordi,  
406 B. E., Lee, B., Paeng, K., Zhong, A., et al. (2018). From detection of individual metastases  
407 to classification of lymph node status at the patient level: the camelyon17 challenge. *IEEE*  
408 *Transactions on Medical Imaging*, 38(2):550–560. [Cited on pages 3, 9, and 24.]
- 409 [6] Beery, S., Van Horn, G., and Perona, P. (2018). Recognition in terra incognita. In *Proceedings of*  
410 *the European Conference on Computer Vision*, pages 456–473. [Cited on page 3.]
- 411 [7] Bickel, S., Brückner, M., and Scheffer, T. (2009). Discriminative learning under covariate shift.  
412 *Journal of Machine Learning Research*, 10(9). [Cited on page 23.]
- 413 [8] Blanchard, G., Flaska, M., Handy, G., Pozzi, S., and Scott, C. (2016). Classification with  
414 asymmetric label noise: Consistency and maximal denoising. *Electronic Journal of Statistics*,  
415 10:2780–2824. [Cited on pages 2 and 17.]
- 416 [9] Blanchard, G., Lee, G., and Scott, C. (2011). Generalizing from several related classification tasks  
417 to a new unlabeled sample. In *Advances in Neural Information Processing Systems*, volume 24.  
418 [Cited on page 3.]
- 419 [10] Blum, A. and Mitchell, T. (1998). Combining labeled and unlabeled data with co-training. In  
420 *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100.  
421 [Cited on pages 3 and 28.]
- 422 [11] Buciluă, C., Caruana, R., and Niculescu-Mizil, A. (2006). Model compression. In *Proceedings*  
423 *of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*,  
424 pages 535–541. [Cited on page 17.]
- 425 [12] Bui, M.-H., Tran, T., Tran, A., and Phung, D. (2021). Exploiting domain-specific features to  
426 enhance domain generalization. In *Advances in Neural Information Processing Systems*, volume 34.  
427 [Cited on page 28.]
- 428 [13] Eastwood, C., Mason, I., and Williams, C. (2021). Unit-level surprise in neural networks. In *I*  
429 *(Still) Can't Believe It's Not Better! NeurIPS 2021 Workshop*. [Cited on page 2.]
- 430 [14] Eastwood, C., Robey, A., Singh, S., von Kügelgen, J., Hassani, H., Pappas, G. J., and Schölkopf,  
431 B. (2022). Probable domain generalization via quantile risk minimization. In *Advances in Neural*  
432 *Information Processing Systems*. [Cited on pages 1, 2, 3, 8, 9, and 27.]
- 433 [15] Eastwood, C. and Williams, C. K. (2018). A framework for the quantitative evaluation of  
434 disentangled representations. In *International Conference on Learning Representations*. [Cited on  
435 page 8.]
- 436 [16] Fei-Fei, L., Fergus, R., and Perona, P. (2006). One-shot learning of object categories. *IEEE*  
437 *Transactions on Pattern Analysis and Machine Intelligence*, 28(4):594–611. [Cited on page 2.]
- 438 [17] Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of  
439 deep networks. In *International Conference on Machine Learning*, pages 1126–1135. [Cited on  
440 page 2.]
- 441 [18] Fréchet, M. (1951). Sur les tableaux de corrélation dont les marges sont données. *Ann. Univ.*  
442 *Lyon, 3<sup>e</sup> e serie, Sciences, Sect. A*, 14:53–77. [Cited on page 6.]

- 443 [19] Galstyan, A. and Cohen, P. R. (2008). Empirical comparison of “hard” and “soft” label  
444 propagation for relational classification. In *Inductive Logic Programming: 17th International*  
445 *Conference, ILP 2007, Corvallis, OR, USA, June 19-21, 2007, Revised Selected Papers 17*, pages  
446 98–111. Springer. [Cited on page 6.]
- 447 [20] Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., and Wichmann,  
448 F. A. (2020). Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2:665–673.  
449 [Cited on pages 1 and 3.]
- 450 [21] Gretton, A., Bousquet, O., Smola, A., and Schölkopf, B. (2005). Measuring statistical depen-  
451 dence with hilbert-schmidt norms. In *Algorithmic Learning Theory: 16th International Conference,*  
452 *ALT 2005, Singapore, October 8-11, 2005. Proceedings 16*, pages 63–77. Springer. [Cited on  
453 page 8.]
- 454 [22] Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., and Schölkopf, B. (2009).  
455 Covariate shift by kernel mean matching. *Dataset shift in machine learning*, 3(4):5. [Cited on  
456 page 23.]
- 457 [23] Gulrajani, I. and Lopez-Paz, D. (2020). In search of lost domain generalization. *arXiv preprint*  
458 *arXiv:2007.01434*. [Cited on pages 3, 9, 24, and 27.]
- 459 [24] Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. (2017). On calibration of modern neural  
460 networks. In *International Conference on Machine Learning*, pages 1321–1330. [Cited on pages 8,  
461 26, 27, and 28.]
- 462 [25] Hendrycks, D. and Dietterich, T. (2019). Benchmarking neural network robustness to common  
463 corruptions and perturbations. In *International Conference on Learning Representations*. [Cited  
464 on page 1.]
- 465 [26] Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the knowledge in a neural network.  
466 *arXiv preprint arXiv:1503.02531*. [Cited on page 17.]
- 467 [27] Hoeffding, W. (1940). Masstabinvariante korrelationstheorie. *Schriften des Mathematischen*  
468 *Instituts und Instituts für Angewandte Mathematik der Universität Berlin*, 5:181–233. [Cited on  
469 page 6.]
- 470 [28] Hoeffding, W. (1941). Masstabinvariante korrelationsmasse für diskontinuierliche verteilungen.  
471 *Archiv für mathematische Wirtschafts-und Sozialforschung*, 7:49–70. [Cited on page 6.]
- 472 [29] Iwasawa, Y. and Matsuo, Y. (2021). Test-time classifier adjustment module for model-agnostic  
473 domain generalization. In *Advances in Neural Information Processing Systems*. [Cited on page 28.]
- 474 [30] Jiang, Y. and Veitch, V. (2022). Invariant and transportable representations for anti-causal  
475 domain shifts. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K., editors, *Advances in Neural*  
476 *Information Processing Systems*. [Cited on pages 2, 3, 4, 5, 8, 9, 23, 25, 27, and 28.]
- 477 [31] Kirichenko, P., Izmailov, P., and Wilson, A. G. (2022). Last layer re-training is sufficient for  
478 robustness to spurious correlations. In *Advances in Neural Information Processing Systems*. [Cited  
479 on page 2.]
- 480 [32] Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W.,  
481 Yasunaga, M., Phillips, R. L., Gao, I., Lee, T., David, E., Stavness, I., Guo, W., Earnshaw, B. A.,  
482 Haque, I. S., Beery, S., Leskovec, J., Kundaje, A., Pierson, E., Levine, S., Finn, C., and Liang, P.  
483 (2021). WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on*  
484 *Machine Learning*. [Cited on pages 9, 24, and 27.]
- 485 [33] Krogel, M.-A. and Scheffer, T. (2004). Multi-relational learning, text mining, and semi-  
486 supervised learning for functional genomics. *Machine Learning*, 57:61–81. [Cited on page 3.]
- 487 [34] Krueger, D., Caballero, E., Jacobsen, J.-H., Zhang, A., Binas, J., Zhang, D., Priol, R. L.,  
488 and Courville, A. (2021). Out-of-distribution generalization via risk extrapolation (rex). In  
489 *International Conference on Machine Learning*, volume 139, pages 5815–5826. [Cited on pages 1,  
490 9, and 27.]

- 491 [35] Lee, D.-H. et al. (2013). Pseudo-label: The simple and efficient semi-supervised learning  
492 method for deep neural networks. In *Workshop on Challenges in Representation Learning, ICML*,  
493 volume 3. [Cited on pages 6 and 28.]
- 494 [36] Li, D., Yang, Y., Song, Y.-Z., and Hospedales, T. M. (2017a). Deeper, broader and artier domain  
495 generalization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.  
496 [Cited on pages 9 and 24.]
- 497 [37] Li, Y., Yang, J., Song, Y., Cao, L., Luo, J., and Li, L.-J. (2017b). Learning from noisy labels  
498 with distillation. In *Proceedings of the IEEE international conference on computer vision*, pages  
499 1910–1918. [Cited on page 2.]
- 500 [38] Liang, J., Hu, D., and Feng, J. (2020). Do we really need to access the source data? Source  
501 hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine  
502 Learning (ICML)*, pages 6028–6039. [Cited on page 28.]
- 503 [39] Makar, M., Packer, B., Moldovan, D., Blalock, D., Halpern, Y., and D’Amour, A. (2022).  
504 Causally motivated shortcut removal using auxiliary labels. In *International Conference on  
505 Artificial Intelligence and Statistics*, pages 739–766. PMLR. [Cited on pages 2 and 8.]
- 506 [40] Muandet, K., Balduzzi, D., and Schölkopf, B. (2013). Domain generalization via invariant  
507 feature representation. In *International Conference on Machine Learning*, pages 10–18. [Cited on  
508 page 3.]
- 509 [41] Nagarajan, V., Andreassen, A., and Neyshabur, B. (2021). Understanding the failure modes  
510 of out-of-distribution generalization. In *International Conference on Learning Representations*.  
511 [Cited on page 3.]
- 512 [42] Natarajan, N., Dhillon, I. S., Ravikumar, P. K., and Tewari, A. (2013). Learning with noisy  
513 labels. *Advances in neural information processing systems*, 26. [Cited on pages 2 and 17.]
- 514 [43] Peters, J., Bühlmann, P., and Meinshausen, N. (2016). Causal inference by using invariant  
515 prediction: identification and confidence intervals. *Journal of the Royal Statistical Society. Series  
516 B (Statistical Methodology)*, pages 947–1012. [Cited on pages 1, 2, 3, and 8.]
- 517 [44] Rish, I. et al. (2001). An empirical study of the naive bayes classifier. In *IJCAI 2001 workshop  
518 on empirical methods in artificial intelligence*, volume 3, pages 41–46. [Cited on page 28.]
- 519 [45] Rojas-Carulla, M., Schölkopf, B., Turner, R., and Peters, J. (2018). Invariant models for causal  
520 transfer learning. *The Journal of Machine Learning Research*, 19(1):1309–1342. [Cited on pages 5,  
521 23, and 28.]
- 522 [46] Rosenfeld, E., Ravikumar, P., and Risteski, A. (2022). Domain-adjusted regression or: ERM  
523 may already learn features sufficient for out-of-distribution generalization. *arXiv preprint  
524 arXiv:2202.06856*. [Cited on pages 2 and 3.]
- 525 [47] Rothenhäusler, D., Meinshausen, N., Bühlmann, P., and Peters, J. (2021). Anchor regression:  
526 Heterogeneous data meet causality. *Journal of the Royal Statistical Society: Series B (Statistical  
527 Methodology)*, 83(2):215–246. [Cited on pages 2 and 8.]
- 528 [48] Rusak, E., Schneider, S., Pachitariu, G., Eck, L., Gehler, P. V., Bringmann, O., Brendel, W., and  
529 Bethge, M. (2022). If your data distribution shifts, use self-learning. *Transactions on Machine  
530 Learning Research*. [Cited on pages 6 and 28.]
- 531 [49] Schölkopf, B. (2022). Causality for machine learning. In *Probabilistic and Causal Inference:  
532 The Works of Judea Pearl*, pages 765–804. Association for Computing Machinery. [Cited on  
533 pages 2 and 23.]
- 534 [50] Scott, C., Blanchard, G., and Handy, G. (2013). Classification with asymmetric label noise:  
535 Consistency and maximal denoising. In *Conference on learning theory*, pages 489–511. PMLR.  
536 [Cited on page 2.]
- 537 [51] Song, H., Kim, M., Park, D., Shin, Y., and Lee, J.-G. (2022). Learning from noisy labels with  
538 deep neural networks: A survey. *IEEE Transactions on Neural Networks and Learning Systems*.  
539 [Cited on page 2.]

- 540 [52] Sugiyama, M. and Kawanabe, M. (2012). *Machine learning in non-stationary environments: Introduction to covariate shift adaptation*. MIT press. [Cited on page 23.]  
541
- 542 [53] Sugiyama, M., Krauledat, M., and Müller, K.-R. (2007). Covariate shift adaptation by im-  
543 portance weighted cross validation. *Journal of Machine Learning Research*, 8(5). [Cited on  
544 page 23.]
- 545 [54] Sun, Q., Murphy, K., Ebrahimi, S., and D’Amour, A. (2022). Beyond invariance: Test-time  
546 label-shift adaptation for distributions with "spurious" correlations. [Cited on page 28.]
- 547 [55] Tanaka, D., Ikami, D., Yamasaki, T., and Aizawa, K. (2018). Joint optimization framework for  
548 learning with noisy labels. In *Proceedings of the IEEE conference on computer vision and pattern  
549 recognition*, pages 5552–5560. [Cited on page 2.]
- 550 [56] Vapnik, V. (1991). Principles of risk minimization for learning theory. *Advances in neural  
551 information processing systems*, 4. [Cited on pages 3 and 23.]
- 552 [57] Vapnik, V. N. (1998). *Statistical Learning Theory*. Wiley, New York, NY. [Cited on page 3.]
- 553 [58] Veitch, V., D’Amour, A., Yadlowsky, S., and Eisenstein, J. (2021). Counterfactual invariance  
554 to spurious correlations: Why and how to pass stress tests. In *Advances in Neural Information  
555 Processing Systems*. [Cited on pages 2 and 8.]
- 556 [59] von Kügelgen, J., Mey, A., and Loog, M. (2019). Semi-generative modelling: Covariate-shift  
557 adaptation with cause and effect features. In *The 22nd International Conference on Artificial  
558 Intelligence and Statistics*, pages 1361–1369. PMLR. [Cited on pages 5 and 23.]
- 559 [60] von Kügelgen, J., Sharma, Y., Gresele, L., Brendel, W., Schölkopf, B., Besserve, M., and  
560 Locatello, F. (2021). Self-Supervised Learning with Data Augmentations Provably Isolates  
561 Content from Style. In *Advances in Neural Information Processing Systems*. [Cited on page 28.]
- 562 [61] Wang, D., Shelhamer, E., Liu, S., Olshausen, B., and Darrell, T. (2021). Tent: Fully test-time  
563 adaptation by entropy minimization. In *International Conference on Learning Representations*.  
564 [Cited on page 28.]
- 565 [62] Wang, W. and Zhou, Z.-H. (2010). A new analysis of co-training. In *ICML*, volume 2, page 3.  
566 [Cited on page 28.]
- 567 [63] Zech, J. R., Badgeley, M. A., Liu, M., Costa, A. B., Titano, J. J., and Oermann, E. K. (2018).  
568 Variable generalization performance of a deep learning model to detect pneumonia in chest  
569 radiographs: a cross-sectional study. *PLoS Medicine*, 15(11). [Cited on page 3.]
- 570 [64] Zhang, J., Lopez-Paz, D., and Bottou, L. (2022). Rich feature construction for the optimization-  
571 generalization dilemma. In *International Conference on Machine Learning*. [Cited on pages 2  
572 and 27.]
- 573 [65] Zheng, J. and Makar, M. (2022). Causally motivated multi-shortcut identification & removal.  
574 *Advances in neural information processing systems*. [Cited on pages 2 and 8.]

575

# 576 Appendices

577

## 578 Table of Contents

---

579	<b>A Proof and further discussion of Theorem 4.4</b>	<b>15</b>
580	A.1 Proof of Theorem 4.4 . . . . .	15
581	A.2 Further discussion of Theorem 4.4 . . . . .	17
582	<b>B Proof of Theorem 4.5</b>	<b>18</b>
583	<b>C Multiclass Case</b>	<b>20</b>
584	<b>D Supplementary Results</b>	<b>22</b>
585	D.1 Causal Perspectives . . . . .	23
586	<b>E Datasets</b>	<b>24</b>
587	<b>F Further Experiments</b>	<b>24</b>
588	F.1 Synthetic dataset . . . . .	24
589	F.2 ColorMNIST . . . . .	26
590	<b>G Implementation Details</b>	<b>27</b>
591	G.1 ColorMNIST . . . . .	27
592	G.2 PACS . . . . .	27
593	G.3 Camelyon17 . . . . .	27
594	G.4 Synthetic . . . . .	28
595	<b>H Further Related Work</b>	<b>28</b>
596	<b>I Limitations</b>	<b>28</b>

---

597  
598  
599

600 **A Proof and further discussion of Theorem 4.4**

601 **A.1 Proof of Theorem 4.4**

602 In this section, we prove our main results regarding the marginal generalization problem presented in  
603 Section 4, namely Theorem 4.4. For the reader's convenience, we restate Theorem 4.4 here:

604 **Theorem 4.4** (Marginal generalization with for binary labels and complementary features). *Consider*  
605 *three random variables  $X_1$ ,  $X_2$ , and  $Y$ , where*

- 606 1.  $Y$  is binary ( $\{0, 1\}$ -valued),  
607 2.  $X_1$  and  $X_2$  are complementary features for  $Y$  (i.e.,  $X_1 \perp\!\!\!\perp X_2 | Y$ ), and  
608 3.  $X_1$  is informative of  $Y$  ( $X_1 \not\perp\!\!\!\perp Y$ ).

609 *Then, the joint distribution of  $(X_1, X_2, Y)$  can be written in terms of the joint distributions of  $(X_1, Y)$*   
610 *and  $(X_1, X_2)$ . Specifically, if  $\hat{Y}|X_S \sim \text{Bernoulli}(\Pr[Y = 1|X_S])$  is pseudo-label and*

$$\epsilon_0 := \Pr[\hat{Y} = 0|Y = 0] \quad \text{and} \quad \epsilon_1 := \Pr[\hat{Y} = 1|Y = 1] \quad (\text{A.1})$$

611 *are the conditional probabilities that  $\hat{Y}$  and  $Y$  agree, given  $Y = 0$  and  $Y = 1$ , respectively, then,*

- 612 1.  $\epsilon_0 + \epsilon_1 > 1$ ,  
613 2.  $\Pr[Y = 1|X_2] = \frac{\Pr[\hat{Y} = 1|X_2] + \epsilon_0 - 1}{\epsilon_0 + \epsilon_1 - 1}$ , and  
614 3.  $\Pr[Y = 1|X_1, X_2] = \sigma(\text{logit}(\Pr[Y = 1|X_1]) + \text{logit}(\Pr[Y = 1|X_2]) - \text{logit}(\Pr[Y = 1]))$ .

615 Before proving Theorem 4.4, we provide some examples demonstrating that the complementarity  
616 and informativeness assumptions in Theorem 4.4 cannot be dropped.

617 **Example A.1.** Suppose  $X_1$  and  $X_2$  have independent  $\text{Bernoulli}(1/2)$  distributions. Then,  $X_1$  is  
618 informative of both of the binary variables  $Y_1 = X_1 X_2$  and  $Y_2 = X_1(1 - X_2)$  and both have identical  
619 conditional distributions given  $X_1$ , but  $Y_1$  and  $Y_2$  have different conditional distributions given  $X_2$ :

$$\Pr[Y_1 = 1|X_2 = 0] = 0 \neq 1/2 = \Pr[Y_2 = 1|X_2 = 0].$$

620 Thus, the complementarity condition cannot be omitted.

621 On the other hand,  $X_1$  and  $X_2$  are complementary for both  $Y_3 = X_2$  and an independent  $Y_4 \sim$   
622  $\text{Bernoulli}(1/2)$  and both  $Y_3$  and  $Y_4$  both have identical conditional distributions given  $X_1$ , but  $Y_1$   
623 and  $Y_2$  have different conditional distributions given  $X_2$ :

$$\Pr[Y_3 = 1|X_2 = 1] = 1/2 \neq 1 = \Pr[Y_4 = 1|X_2 = 1].$$

624 Thus, the informativeness condition cannot be omitted.

625 Before proving Theorem 4.4, we prove Lemma A.2, which allows us to safely divide by the quantity  
626  $\epsilon_0 + \epsilon_1 - 1$  in the formula for  $\Pr[Y = 1|X_2]$ , under the condition that  $X_1$  is informative of  $Y$ .

627 **Lemma A.2.** *In the setting of Theorem 4.4, let  $\epsilon_0$  and  $\epsilon_1$  be the class-wise pseudo-label accuracies*  
628 *defined in as in Eq. (A.1). Then,  $\epsilon_0 + \epsilon_1 = 1$  if and only if  $X_1$  and  $Y$  are independent.*

629 Note that the entire result also holds, with almost identical proof, in the multi-environment setting of  
630 Sections 3 and 5, conditioned on a particular environment  $E$ .

631 *Proof.* We first prove the forwards implication. Suppose  $\epsilon_0 + \epsilon_1 = 1$ . If  $\Pr[Y = 1] \in \{0, 1\}$ , then  
632  $X_1$  and  $Y$  are trivially independent, so we may assume  $\Pr[Y = 1] \in (0, 1)$ . Then,

$$\begin{aligned} \mathbb{E}[\hat{Y}] &= \epsilon_1 \Pr[Y = 1] + (1 - \epsilon_0)(1 - \Pr[Y = 1]) && \text{(Law of Total Expectation)} \\ &= (\epsilon_0 + \epsilon_1 - 1) \Pr[Y = 1] + 1 - \epsilon_0 \\ &= 1 - \epsilon_0 && (\epsilon_0 + \epsilon_1 = 1) \\ &= \mathbb{E}[\hat{Y}|Y = 0]. && \text{(Definition of } \epsilon_0) \end{aligned}$$

633 Since  $Y$  is binary and  $\Pr[Y = 1] \in (0, 1)$ , it follows that  $\mathbb{E}[\hat{Y}] = \mathbb{E}[\hat{Y}|Y = 0] = \mathbb{E}[\hat{Y}|Y = 1]$ ; i.e.,  
634  $\mathbb{E}[\hat{Y}|Y] \perp\!\!\!\perp Y$ . Since  $\hat{Y}$  is binary, its distribution is specified entirely by its mean, and so  $\hat{Y} \perp\!\!\!\perp Y$ . It  
635 follows that the covariance between  $\hat{Y}$  and  $Y$  is 0:

$$\begin{aligned} 0 &= \mathbb{E}[(Y - \mathbb{E}[Y])(\hat{Y} - \mathbb{E}[\hat{Y}])] \\ &= \mathbb{E}[\mathbb{E}[(Y - \mathbb{E}[Y])(\hat{Y} - \mathbb{E}[\hat{Y}])|X_1]] && \text{(Law of Total Expectation)} \\ &= \mathbb{E}[\mathbb{E}[Y - \mathbb{E}[Y]|X_1] \mathbb{E}[\hat{Y} - \mathbb{E}[\hat{Y}]|X_1]] && (Y \perp\!\!\!\perp \hat{Y}|X_1) \\ &= \mathbb{E}[(\mathbb{E}[Y - \mathbb{E}[Y]|X_1])^2], \end{aligned}$$

636 where the final equality holds because  $\hat{Y}$  and  $Y$  have identical conditional distributions given  $X_1$ .  
637 Since the  $\mathcal{L}_2$  norm of a random variable is 0 if and only if the variable is 0 almost surely, it follows  
638 that,  $P_{X_1}$ -almost surely,

$$0 = \mathbb{E}[Y - \mathbb{E}[Y]|X_1] = \mathbb{E}[Y|X_1] - \mathbb{E}[Y],$$

639 so that  $\mathbb{E}[Y|X_1] \perp\!\!\!\perp X_1$ . Since  $Y$  is binary, its distribution is specified entirely by its mean, and so  
640  $Y \perp\!\!\!\perp X_1$ , proving the forwards implication.

641 To prove the reverse implication, suppose  $X_1$  and  $Y$  are independent. Then  $\hat{Y}$  and  $Y$  are also  
642 independent. Hence,

$$\epsilon_1 = \mathbb{E}[\hat{Y}|Y = 1] = \mathbb{E}[\hat{Y}|Y = 0] = 1 - \epsilon_0,$$

643 so that  $\epsilon_0 + \epsilon_1 = 1$ . □

644 We now use Lemma A.2 to prove Theorem 4.4:

645 *Proof.* To begin, note that  $\hat{Y}$  has the same conditional distribution given  $X_1$  as  $Y$  (i.e.,  $P_{\hat{Y}|X_1} = P_{Y|X_1}$ )  
646 and that  $\hat{Y}$  is conditionally independent of  $Y$  given  $X_1$  ( $\hat{Y} \perp\!\!\!\perp Y|X_1$ ). Then, since

$$\Pr[\hat{Y} = 1] = \mathbb{E}[\Pr[Y = 1|X_1]] = \Pr[Y = 1], \quad (\text{A.2})$$

647 we have

$$\begin{aligned} \epsilon_1 = \Pr[\hat{Y} = 1|Y = 1] &= \frac{\Pr[Y = 1, \hat{Y} = 1]}{\Pr[Y = 1]} && \text{(Definition of } \epsilon_1) \\ &= \frac{\Pr[Y = 1, \hat{Y} = 1]}{\Pr[\hat{Y} = 1]} && \text{(Eq. (A.2))} \\ &= \frac{\mathbb{E}_{X_1}[\Pr[Y = 1, \hat{Y} = 1|X_1]]}{\mathbb{E}_{X_1}[\Pr[\hat{Y} = 1|X_1]]} && \text{(Law of Total Expectation)} \\ &= \frac{\mathbb{E}_{X_1}[\Pr[Y = 1|X_1] \Pr[\hat{Y} = 1|X_1]]}{\mathbb{E}_{X_1}[\Pr[\hat{Y} = 1|X_1]]} && (\hat{Y} \perp\!\!\!\perp Y|X_1) \\ &= \frac{\mathbb{E}_{X_1}[(\Pr[Y = 1|X_1])^2]}{\mathbb{E}_{X_1}[\Pr[Y = 1|X_1]]} && (P_{\hat{Y}|X_1} = P_{Y|X_1}) \end{aligned}$$

648 entirely in terms of the conditional distribution  $P_{Y|X_1}$  and the marginal distribution  $P_{X_1}$ . Similarly,

649  $\epsilon_0$  can be written as  $\epsilon_0 = \frac{\mathbb{E}_{X_1}[(\Pr[Y = 0|X_1])^2]}{\mathbb{E}_{X_1}[\Pr[Y = 0|X_1]]}$ . Meanwhile, by the law of total expectation, and

650 the assumption that  $X_1$  (and hence  $\hat{Y}$ ) is conditionally independent of  $X_2$  given  $Y$ , the conditional  
651 distribution  $P_{\hat{Y}|X_2}$  of  $\hat{Y}$  given  $X_2$  can be written as

$$\begin{aligned} &\Pr[\hat{Y} = 1|X_2] \\ &= \Pr[\hat{Y} = 1|Y = 0, X_2] \Pr[Y = 0|X_2] + \Pr[\hat{Y} = 1|Y = 1, X_2] \Pr[Y = 1|X_2] \\ &= \Pr[\hat{Y} = 1|Y = 0] \Pr[Y = 0|X_2] + \Pr[\hat{Y} = 1|Y = 1] \Pr[Y = 1|X_2] \\ &= (1 - \epsilon_0)(1 - \Pr[Y = 1|X_2]) + \epsilon_1 \Pr[Y = 1|X_2] \\ &= (\epsilon_0 + \epsilon_1 - 1) \Pr[Y = 1|X_2] + 1 - \epsilon_0. \end{aligned}$$

652 By Lemma A.2, the assumption  $X_1 \not\perp\!\!\!\perp Y$  implies  $\epsilon_0 + \epsilon_1 \neq 1$ . Hence, re-arranging the above  
653 equality gives us the conditional distribution  $P_{Y|X_2}$  of  $Y$  given  $X_2$  purely in terms of the conditional  
654  $P_{Y|X_1}$  and  $P_{X_1, X_2}$ :

$$\Pr[Y = 1|X_2 = X_2] = \frac{\Pr[\hat{Y} = 1|X_2 = X_2] + \epsilon_0 - 1}{\epsilon_0 + \epsilon_1 - 1}.$$

655 It remains now to write the conditional distribution  $P_{Y|X_1, X_2}$  in terms of the conditional distributions  
656  $P_{Y|X_1}$  and  $P_{Y|X_2}$  and the marginal  $P_Y$ . Note that

$$\begin{aligned} \frac{\Pr[Y = 1|X_1, X_2]}{\Pr[Y = 0|X_1, X_2]} &= \frac{\Pr[X_1, X_2|Y = 1] \Pr[Y = 1]}{\Pr[X_1, X_2|Y = 0] \Pr[Y = 0]} && \text{(Bayes' Rule)} \\ &= \frac{\Pr[X_1|Y = 1] \Pr[X_2|Y = 1] \Pr[Y = 1]}{\Pr[X_1|Y = 0] \Pr[X_2|Y = 0] \Pr[Y = 0]} && \text{(Complementarity)} \\ &= \frac{\Pr[Y = 1|X_1] \Pr[Y = 1|X_2] \Pr[Y = 0]}{\Pr[Y = 0|X_1] \Pr[Y = 0|X_2] \Pr[Y = 1]}. && \text{(Bayes' Rule)} \end{aligned}$$

657 It follows that the logit of  $\Pr[Y = 1|X_1, X_2]$  can be written as the sum of a term depending only on  
658  $X_1$ , a term depending only on  $X_2$ , and a constant term:

$$\begin{aligned} \text{logit}(\Pr[Y = 1|X_1, X_2]) &= \log \frac{\Pr[Y = 1|X_1, X_2]}{1 - \Pr[Y = 1|X_1, X_2]} \\ &= \log \frac{\Pr[Y = 1|X_1, X_2]}{\Pr[Y = 0|X_1, X_2]} \\ &= \log \frac{\Pr[Y = 1|X_1]}{\Pr[Y = 0|X_1]} + \log \frac{\Pr[Y = 1|X_2]}{\Pr[Y = 0|X_2]} - \log \frac{\Pr[Y = 1]}{\Pr[Y = 0]} \\ &= \text{logit}(\Pr[Y = 1|X_1]) + \text{logit}(\Pr[Y = 1|X_2]) - \text{logit}(\Pr[Y = 1]). \end{aligned}$$

659 Since the sigmoid  $\sigma$  is the inverse of logit,

$$\Pr[Y = 1|X_1, X_2] = \sigma(\text{logit}(\Pr[Y = 1|X_1]) + \text{logit}(\Pr[Y = 1|X_2]) - \text{logit}(\Pr[Y = 1])),$$

660 which, by Eq. (4.1), can be written in terms of the conditional distribution  $P_{Y|X_1}$  and the joint  
661 distribution  $P_{X_1, X_2}$ .  $\square$

## 662 A.2 Further discussion of Theorem 4.4

663 **Connections to learning from noisy labels.** Theorem 4.4 leverages two theoretical insights about  
664 the special structure of pseudo-labels that complement results in the literature on learning from noisy  
665 labels. First, Blanchard et al. [8] showed that learning from noisy labels is possible if and only if the  
666 total noise level is below the critical threshold  $\epsilon_0 + \epsilon_1 > 1$ ; in the case of learning from pseudo-labels,  
667 we show (see Lemma A.2 in Appendix A.1) that this is satisfied if and only if  $X_S$  is informative of  $Y$   
668 (i.e.,  $Y \not\perp\!\!\!\perp X_S$ ). Second, methods for learning under label noise commonly assume knowledge of  $\epsilon_0$   
669 and  $\epsilon_1$  [42], which is unrealistic in many applications; however, for pseudo-labels sampled from a  
670 known conditional probability distribution  $P_{Y|X_S}$ , one can express these noise levels we show (as part  
671 of Theorem 4.4) that the class-conditional noise levels can be easily estimated.

672 **Possible applications of Theorem 4.4 beyond domain adaptation** The reason we wrote Theorem  
673 4.4 in the more general setting of the marginal problem rather than in the specific context of  
674 domain adaptation is that we envision possible applications to a number of problems besides domain  
675 adaptation. For example, suppose that, after learning a calibrated machine learning model  $M_1$  using  
676 a feature  $X_1$ , we observe an additional feature  $X_2$ . In the case that  $X_1$  and  $X_2$  are complementary,  
677 Theorem 4.4 justifies using the student-teacher paradigm [11, 3, 26] to train a model for predicting  $Y$   
678 from  $X_2$  (or from  $(X_1, X_2)$  jointly) based on predictions from  $M_1$ . This could be useful if we don't  
679 have access to labeled pairs  $(X_2, Y)$ , or if retraining a model using  $X_1$  would require substantial  
680 computational resources or access to sensitive or private data. Exploring such approaches could be a  
681 fruitful direction for future work

682 **B Proof of Theorem 4.5**

683 This appendix provides a proof of Theorem 4.5, which provides conditions under which our proposed  
684 domain adaptation procedure (Alg. 1) is consistent.

685 We state provide a formal version of Theorem 4.5:

686 **Theorem 4.5** (Consistency of the bias-corrected classifier). *Assume*

- 687 1.  $X_S$  is stable,  
688 2.  $X_S$  and  $X_U$  are complementary, and  
689 3.  $X_S$  is informative of  $Y$  (i.e.,  $X_S \not\perp Y$ ).

690 Let  $\hat{\eta}_n : \mathcal{X}_S \times X_U \rightarrow [0, 1]$  given by

$$\hat{\eta}_n(x_S, x_U) = \sigma \left( f_S(x_S) + \text{logit} \left( \frac{\hat{\eta}_{U,n}(x_U) + \hat{\epsilon}_{0,n} - 1}{\hat{\epsilon}_{0,n} + \hat{\epsilon}_{1,n} - 1} \right) - \beta_1 \right), \quad \text{for all } (x_S, x_U) \in \mathcal{X}_S \times \mathcal{X}_U,$$

691 denote the bias-corrected regression function estimate proposed in Alg. 1, and let  $\hat{h}_n : \mathcal{X}_S \times \mathcal{X}_U \rightarrow$   
692  $\{0, 1\}$  given by

$$\hat{h}_n(x_S, x_U) = 1\{\hat{\eta}(x_S, x_U) > 1/2\}, \quad \text{for all } (x_S, x_U) \in \mathcal{X}_S \times \mathcal{X}_U,$$

693 denote the corresponding hard classifier. Let  $\eta_U : \mathcal{X}_U \rightarrow [0, 1]$ , given by  $\eta_U(x_U) = \Pr[Y =$   
694  $1 | X_U = x_U, E = 1]$  for all  $x_U \in \mathcal{X}_U$ , denote the true regression function over  $X_U$ , and let  $\hat{\eta}_{U,n}$   
695 denote its estimate as assumed in Line 3 of Alg. 1. Then, as  $n \rightarrow \infty$ ,

- 696 (a) if, for  $P_{X_U}$ -almost all  $x_U \in \mathcal{X}_U$ ,  $\hat{\eta}_{U,n}(x_U) \rightarrow \eta_U(x_U)$  in probability, then  $\hat{\eta}_n$  and  $\hat{h}_n$  are  
697 weakly consistent (i.e.,  $\hat{\eta}_n(x_S, x_U) \rightarrow \eta(x_S, x_U)$   $P_{X_S, X_U}$ -almost surely and  $R(\hat{h}_n) \rightarrow R(h^*)$  in  
698 probability).  
699 (b) if, for  $P_{X_U}$ -almost all  $x_U \in \mathcal{X}_U$ ,  $\hat{\eta}_{U,n}(x_U) \rightarrow \eta_U(x_U)$  almost surely, then  $\hat{\eta}_n$  and  $\hat{h}_n$  are  
700 strongly consistent (i.e.,  $\hat{\eta}_n(x_S, x_U) \rightarrow \eta(x_S, x_U)$   $P_{X_S, X_U}$ -almost surely and  $R(\hat{h}_n) \rightarrow R(h^*)$   
701 a.s.).

702 Before proving Theorem 4.5, we provide a few technical lemmas. The first shows that almost-  
703 everywhere convergence of regression functions implies convergence of the corresponding classifiers  
704 in classification risk:

705 **Lemma B.1.** Consider a sequence of regression functions  $\eta, \eta_1, \eta_2, \dots : \mathcal{X} \rightarrow [0, 1]$ . Let  $h, h_1, h_2, \dots :$   
706  $\mathcal{X} \rightarrow \{0, 1\}$  denote the corresponding classifiers

$$h(x) = 1\{\eta(x) > 1/2\} \quad \text{and} \quad h_i(x) = 1\{\eta_i(x) > 1/2\}, \quad \text{for all } i \in \mathbb{N}, x \in \mathcal{X}.$$

- 707 (a) If  $\eta_n(x) \rightarrow \eta(x)$  for  $P_X$ -almost all  $x \in \mathcal{X}$  in probability, then  $R(h_n) \rightarrow R(h^*)$  in probability.  
708 (b) If  $\eta_n(x) \rightarrow \eta(x)$  for  $P_X$ -almost all  $x \in \mathcal{X}$  almost surely as  $n \rightarrow \infty$ , then  $R(h_n) \rightarrow R(h)$   
709 almost surely.

710 *Proof.* Note that, since  $h_n(x) \neq h(x)$  implies  $|\eta_n(x) - \eta(x)| \geq |\eta(x) - 1/2|$ ,

$$1\{h_n(x) \neq h(x)\} \leq 1\{|\eta_n(x) - \eta(x)| \geq |\eta(x) - 1/2|\}. \quad (\text{B.1})$$

711 We utilize this observation to prove both (a) and (b).

712 **Proof of (a)** Let  $\delta > 0$ . By Inequality (B.1) and partitioning  $\mathcal{X}$  based on whether  $|2\eta(X) - 1| \leq$   
713  $\delta/2$ ,

$$\begin{aligned} & \mathbb{E}_X [ |2\eta(X) - 1| 1\{h_n(X) \neq h(X)\} ] \\ & \leq \mathbb{E}_X [ |2\eta(X) - 1| 1\{|\eta_n(X) - \eta(X)| \geq |\eta(X) - 1/2|\} ] \\ & = \mathbb{E}_X [ |2\eta(X) - 1| 1\{|\eta_n(X) - \eta(X)| \geq |\eta(X) - 1/2|\} 1\{|2\eta(X) - 1| > \delta/2\} ] \\ & \quad + \mathbb{E}_X [ |2\eta(X) - 1| 1\{|\eta_n(X) - \eta(X)| \geq |\eta(X) - 1/2|\} 1\{|2\eta(X) - 1| \leq \delta/2\} ] \\ & \leq \mathbb{E}_X [ 1\{|\eta_n(X) - \eta(X)| > \delta/2\} ] + \delta/2. \end{aligned}$$

714 Hence,

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \Pr_{\eta_n} [\mathbb{E}_X [ |2\eta(X) - 1| 1\{h_n(X) \neq h(X)\} ] > \delta] \\
& \leq \lim_{n \rightarrow \infty} \Pr_{\eta_n} [\mathbb{E}_X [ 1\{|\eta_n(X) - \eta(X)| > \delta/2\} ] > \delta/2] \\
& \leq \lim_{n \rightarrow \infty} \frac{2}{\delta} \mathbb{E}_{\eta_n} [\mathbb{E}_X [ 1\{|\eta_n(X) - \eta(X)| > \delta/2\} ]] \quad (\text{Markov's Inequality}) \\
& = \lim_{n \rightarrow \infty} \frac{2}{\delta} \mathbb{E}_X [\mathbb{E}_{\eta_n} [ 1\{|\eta_n(X) - \eta(X)| > \delta/2\} ]] \quad (\text{Fubini's Theorem}) \\
& = \frac{2}{\delta} \mathbb{E}_X \left[ \lim_{n \rightarrow \infty} \Pr_{\eta_n} [ |\eta_n(X) - \eta(X)| > \delta/2 ] \right] \quad (\text{Dominated Convergence Theorem}) \\
& = 0. \quad (\eta_n(X) \rightarrow \eta(X), P_X\text{-a.s., in probability})
\end{aligned}$$

715 **Proof of (b)** For any  $x \in \mathcal{X}$  with  $\eta(x) \neq 1/2$ , if  $\eta_n(x) \rightarrow \eta(x)$  then  $1\{|\eta_n(x) - \eta(x)| \geq$   
716  $|\eta(x) - 1/2|\} \rightarrow 0$ . Hence, by Inequality (B.1), the dominated convergence theorem (with  $|2\eta(x) -$   
717  $1| 1\{|\eta_n(x) - \eta(x)| \geq |\eta(x) - 1/2|\} \leq 1$ ), and the assumption that  $\eta_n(x) \rightarrow \eta(x)$  for  $P_X$ -almost  
718 all  $x \in \mathcal{X}$  almost surely,

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \mathbb{E}_X [ |2\eta(X) - 1| 1\{h_n(X) \neq h(X)\} ] \\
& \leq \lim_{n \rightarrow \infty} \mathbb{E}_X [ |2\eta(X) - 1| 1\{|\eta_n(X) - \eta(X)| \geq |\eta(X) - 1/2|\} ] \\
& = \mathbb{E}_X \left[ \lim_{n \rightarrow \infty} |2\eta(X) - 1| 1\{|\eta_n(x) - \eta(x)| \geq |\eta(x) - 1/2|\} \right] \\
& = 0, \quad \text{almost surely.}
\end{aligned}$$

719 □

720 Our next lemma concerns an edge case in which the features  $X_S$  and  $X_U$  provide perfect but  
721 contradictory information about  $Y$ , leading to Equation (4.2) being ill defined. We show that this can  
722 happen only with probability 0 over  $(X_S, X_U) \sim P_{X_S, X_U}$  can thus be safely ignored:

723 **Lemma B.2.** Consider two predictors  $X_S$  and  $X_U$  of a binary label  $Y$ . Then,

$$\Pr_{X_S, X_U} [\mathbb{E}[Y|X_S] = 1 \text{ and } \mathbb{E}[Y|X_U] = 0] = \Pr_{X_S, X_U} [\mathbb{E}[Y|X_S] = 0 \text{ and } \mathbb{E}[Y|X_U] = 1] = 0.$$

724

725 *Proof.* Suppose, for sake of contradiction, that the event

$$A := \{(x_S, x_U) : \mathbb{E}[Y|X_S = x_S] = 1 \text{ and } \mathbb{E}[Y|X_U = x_U] = 0\}$$

726 has positive probability. Then, the conditional expectation  $\mathbb{E}[Y|A]$  is well-defined, giving the  
727 contradiction

$$1 = \mathbb{E}_{X_S} [\mathbb{E}[Y|E, X_S]] = \mathbb{E}[Y|A] = \mathbb{E}_{X_U} [\mathbb{E}[Y|E, X_U]] = 0.$$

728 The case  $\mathbb{E}[Y|X_S] = 0$  and  $\mathbb{E}[Y|X_U] = 1$  is similar. □

729 We now utilize Lemmas B.1 and B.2 to prove Theorem 4.5.

730 *Proof.* By Lemma B.1, it suffices to prove that  $\hat{\eta}(x_S, x_U) \rightarrow \eta(x_S, x_U)$ , for  $P_{X_S, X_U}$ -almost all  
731  $(x_S, x_U) \in \mathcal{X}_S \times \mathcal{X}_U$ , in probability (to prove (a)) and almost surely (to prove (b)).

732 **Finite case** We first consider the case when both  $\Pr[Y|X_S = x_S], \Pr[Y|X_U = x_U] \in (0, 1)$ , so  
733 that  $f_S(x_S)$  and  $\text{logit} \left( \frac{\hat{\eta}(x_U) + \epsilon_0 - 1}{\epsilon_0 + \epsilon_1 - 1} \right)$  are both finite. Since

$$\begin{aligned}
& \hat{\eta}_{S,U}(x_S, x_U) - \eta_{S,U}(x_S, x_U) \\
& = \sigma \left( f_S(x_S) + \text{logit} \left( \frac{\hat{\eta}_{U,1}(x_U) + \hat{\epsilon}_0 - 1}{\hat{\epsilon}_0 + \hat{\epsilon}_1 - 1} \right) - \hat{\beta}_{1,n} \right) - \sigma \left( f_S(x_S) + \text{logit} \left( \frac{\tilde{\eta}(x_U) + \epsilon_0 - 1}{\epsilon_0 + \epsilon_1 - 1} \right) - \beta_1 \right),
\end{aligned}$$

734 where the sigmoid  $\sigma : \mathbb{R} \rightarrow [0, 1]$  is continuous, by the continuous mapping theorem and the  
735 assumption that  $\hat{\eta}_{U,1}(x_U) \rightarrow \tilde{\eta}(x_U)$ , to prove both of these, it suffices to show:

736 (i)  $\hat{\epsilon}_0 \rightarrow \epsilon_0$  and  $\hat{\epsilon}_1 \rightarrow \epsilon_1$  almost surely as  $n \rightarrow \infty$ .

737 (ii)  $\hat{\beta}_{1,n} \rightarrow \beta_1 \in (-\infty, \infty)$  almost surely as  $n \rightarrow \infty$ .

738 (iii) The mapping  $(a, b, c) \mapsto \text{logit}\left(\frac{a+b-1}{b+c-1}\right)$  is continuous at  $(\hat{\eta}(x_U), \epsilon_0, \epsilon_1)$ .

739 We now prove each of these in turn.

740 **Proof of (i)** Since  $\hat{Y}_i \perp\!\!\!\perp Y_i | X_S$  and  $0 < \Pr[\hat{Y} = 1]$ , by the strong law of large numbers and the  
741 continuous mapping theorem,

$$\hat{\epsilon}_1 = \frac{1}{n_1} \sum_{i=1}^n \hat{Y}_i \sigma(f_S(X_i)) = \frac{\frac{1}{n} \sum_{i=1}^n \hat{Y}_i \sigma(f_S(X_i))}{\frac{1}{n} \sum_{i=1}^n \hat{Y}_i} \rightarrow \frac{\mathbb{E}[\sigma(f_S(X)) 1\{\hat{Y} = 1\}]}{\Pr[\hat{Y} = 1]} = \mathbb{E}[\sigma(f_S(X)) | \hat{Y} = 1] = \epsilon_1,$$

742 almost surely as  $n \rightarrow \infty$ . Similarly, since  $\Pr[\hat{Y} = 0] = 1 - \Pr[\hat{Y} = 1] > 0$ ,  $\hat{\epsilon}_0 \rightarrow \epsilon_0$  almost surely.

743 **Proof of (ii)** Recall that

$$\hat{\beta}_{1,n} = \text{logit}\left(\frac{1}{n} \sum_{i=1}^n \hat{Y}_i\right).$$

744 By the strong law of large numbers,  $\frac{1}{n} \sum_{i=1}^n \hat{Y}_i \rightarrow \Pr[\hat{Y} = 1 | E = 1] = \Pr[Y = 1 | E = 1]$ .  
745 Since we assumed  $\Pr[Y = 1 | E = 1] \in (0, 1)$ , it follows that the mapping  $a \mapsto \text{logit}(a)$  is  
746 continuous at  $a = \Pr[Y = 1 | E = 1]$ . Hence, by the continuous mapping theorem,  $\hat{\beta}_{1,n} \rightarrow$   
747  $\text{logit}(\Pr[Y = 1 | E = 1]) = \beta_1$  almost surely.

748 **Proof of (iii)** Since the logit function is continuous on the open interval  $(0, 1)$  and we assumed  
749  $\epsilon_0 + \epsilon_1 > 1$ , it suffices to show that  $0 < \hat{\eta}(x_U) + \epsilon_0 - 1 < \epsilon_0 + \epsilon_1 - 1$ . Since, according to  
750 Theorem 4.4,

$$\hat{\eta}(x_U) = (\epsilon_0 + \epsilon_1 - 1)\eta^*(x_U) + 1 - \epsilon_0,$$

751 this holds as long as  $0 < \eta^*(x_U) < 1$ , as we assumed for  $P_{X_U}$ -almost all  $x_U \in \mathcal{X}_U$ .

752 **Infinite case** We now address the case where either  $\Pr[Y | X_S = x_S] \in \{0, 1\}$  or  $\Pr[Y | X_U =$   
753  $x_U] \in \{0, 1\}$ . By Lemma B.2, only one of these can happen at once,  $P_{X_S, X_U}$ -almost surely. Hence,  
754 since  $\lim_{n \rightarrow \infty} \hat{\beta}_{1,n}$  is also finite almost surely, if  $\Pr[Y | X_S = x_S] \in \{0, 1\}$ , then  $\hat{\eta}(x_S, x_U) =$   
755  $\sigma(\text{logit}(\Pr[Y | X_S = x_S])) = \eta(x_S, x_U)$ , while, if  $\Pr[Y | X_U = x_U] \in \{0, 1\}$ , then  $\hat{\eta}(x_S, x_U) \rightarrow$   
756  $\sigma(\text{logit}(\Pr[Y | X_U = x_U])) = \eta(x_S, x_U)$ , in probability or almost surely, as appropriate.  $\square$

## 757 C Multiclass Case

758 In the main paper, to simplify notation, we presented our unsupervised test-domain adaptation method  
759 in the case of binary labels  $Y$ . However, in many cases, including several of our experiments in  
760 Section 6, the label  $Y$  can take more than 2 distinct values. Hence, in this section, we show how to  
761 generalize our method to the multiclass setting and then present the exact procedure (Alg. 2) used in  
762 our multiclass experiments in Section 6.

763 Suppose we have  $K \geq 2$  classes. We “one-hot encode” these classes, so that  $Y$  takes values in the set

$$\mathcal{Y} = \{(1, 0, \dots, 0), (0, 1, 0, \dots, 0), \dots, (0, \dots, 0, 1)\} \subseteq \{0, 1\}^K.$$

764 Let  $\epsilon \in [0, 1]^{\mathcal{Y} \times \mathcal{Y}}$  with

$$\epsilon_{y, y'} = \Pr[\hat{Y} = y | Y = y']$$

765 denote the class-conditional confusion matrix of the pseudo-labels. Then, we have

$$\mathbb{E}[\hat{Y} | X_2] = \sum_{y \in \mathcal{Y}} \mathbb{E}[\hat{Y} | Y = y, X_2] \Pr[Y = y | X_2] \quad (\text{Law of Total Expectation})$$

$$= \sum_{y \in \mathcal{Y}} \mathbb{E}[\hat{Y} | Y = y] \Pr[Y = y | X_2] \quad (\text{Complementary})$$

$$= \epsilon \mathbb{E}[Y | X_2]; \quad (\text{Definition of } \epsilon)$$

766 in particular, when  $\epsilon$  is invertible,

$$\mathbb{E}[Y|X_2] = \epsilon^{-1} \mathbb{E}[\hat{Y}|X_2],$$

767 giving a multiclass equivalent of Eq. (4.1) in Theorem 4.4. We also have

$$\begin{aligned} \epsilon_{y,y'} &= \Pr[\hat{Y} = y|Y = y'] = \frac{\Pr[\hat{Y} = y, Y = y']}{\Pr[Y = y']} = \frac{\mathbb{E}[\Pr[\hat{Y} = y, Y = y'|X_1]]}{\mathbb{E}[\Pr[Y = y'|X_1]]} \\ &= \frac{\mathbb{E}[\Pr[\hat{Y} = y|X_1] \Pr[Y = y'|X_1]]}{\mathbb{E}[\Pr[Y = y'|X_1]]} \\ &= \frac{\mathbb{E}[\eta_{1,y}(X_1)\eta_{1,y'}(X_1)]}{\mathbb{E}[\eta_{1,y'}(X_1)]}, \end{aligned}$$

768 suggesting the estimate

$$\hat{\epsilon}_{y,y'} = \frac{\sum_{i=1}^n \hat{\eta}_{S,y}(X_{S,i})\hat{\eta}_{S,y'}(X_{S,i})}{\sum_{i=1}^n \hat{\eta}_{S,y'}(X_{S,i})} = \sum_{i=1}^n \hat{\eta}_{S,y}(X_{S,i}) \frac{\hat{\eta}_{S,y'}(X_{S,i})}{\sum_{i=1}^n \hat{\eta}_{S,y'}(X_{S,i})}$$

769 of each  $\epsilon_{y,y'}$ , or, in matrix notation,

$$\hat{\epsilon} = \eta_S^T(X_S) \text{Normalize}(\eta_S(X_S)),$$

770 where  $\text{Normalize}(X)$  scales each column of  $X$  to sum to 1. This gives us an multiclass equivalent  
771 of Line 4 in Alg. 1.

772 The multiclass versions of Eq. (4.2) and Line 6 of Alg. 1 are slightly less straightforward. Specifically,  
773 whereas, in the binary case, we used the fact that  $\Pr[X_S, X_U|Y \neq 1] = \Pr[X_S, X_U|Y = 0] =$   
774  $\Pr[X_S|Y = 0] \Pr[X_U|Y = 0] = \Pr[X_S|Y \neq 1] \Pr[X_U|Y \neq 1]$  (by complementarity), in the  
775 multiclass case, we do not have  $\Pr[X_S, X_U|Y \neq 1] = \Pr[X_S|Y \neq 1] \Pr[X_U|Y \neq 1]$ . However,  
776 following similar reasoning as in the proof of Theorem 4.4, we have

$$\begin{aligned} \frac{\Pr[Y = y|X_S, X_U, E]}{\Pr[Y \neq y|X_S, X_U, E]} &= \frac{\Pr[Y = y|X_S, X_U, E]}{\sum_{y' \neq y} \Pr[Y = y'|X_S, X_U, E]} \\ &= \frac{\Pr[X_S, X_U|Y = y, E] \Pr[Y = y|E]}{\sum_{y' \neq y} \Pr[X_S, X_U|Y = y', E] \Pr[Y = y'|E]} \quad (\text{Bayes' Rule}) \\ &= \frac{\Pr[X_S|Y = y, E] \Pr[X_U|Y = y, E] \Pr[Y = y|E]}{\sum_{y' \neq y} \Pr[X_S|Y = y', E] \Pr[X_U|Y = y', E] \Pr[Y = y'|E]} \quad (X_S \perp\!\!\!\perp X_U|Y) \\ &= \frac{\Pr[Y = y|X_S, E] \Pr[Y = y|X_U, E]}{\sum_{y' \neq y} \Pr[Y = y'|X_S, E] \Pr[Y = y'|X_U, E]} \cdot \frac{\Pr[Y = y|E]}{\Pr[Y = y'|E]}. \quad (\text{Bayes' Rule}) \end{aligned}$$

777 Hence,

$$\begin{aligned} \text{logit}(\Pr[Y = y|X_S, X_U, E]) &= \log \left( \frac{\Pr[Y = y|X_S, E] \Pr[Y = y|X_U, E]}{\sum_{y' \neq y} \Pr[Y = y'|X_S, E] \Pr[Y = y'|X_U, E]} \cdot \frac{\Pr[Y = y|E]}{\Pr[Y = y'|E]} \right) \\ &= \log \left( \frac{C_y}{\sum_{y' \neq y} C_{y'}} \right) = \log \left( \frac{\frac{C_y}{\|C\|_1}}{\sum_{y' \neq y} \frac{C_{y'}}{\|C\|_1}} \right) = \text{logit} \left( \frac{C_y}{\|C\|_1} \right), \end{aligned}$$

778 for  $C \in \mathbb{R}^{\mathcal{Y}}$  defined by

$$C_y = \frac{\eta_{S,y}(X_S)\eta_{U,y}(X_U)}{\Pr[Y = y]} \quad \text{for each } y \in \mathcal{Y}.$$

779 In particular, applying the sigmoid function to each side, we have

$$\Pr[Y|X_S, X_U] = \frac{C}{\|C\|_1}.$$

780 We can estimate  $C_y$  by

$$\hat{C}_y = \frac{\eta_{S,y}(X_S)\eta_{U,y}(X_U)}{\frac{1}{n}\sum_{i=1}^n \eta_{S,y}(X_{S,i})}.$$

781 In matrix notation, this is

$$\hat{C} = \frac{\eta_S(X_S) \circ \eta_U(X_U)}{\frac{1}{n}\sum_{i=1}^n \eta_S(X_{S,i})},$$

782 where  $\circ$  denotes element-wise multiplication. Putting these derivations together gives us our multi-  
 783 class version of Alg. 1, presented in Alg. 2, where  $\Delta^{\mathcal{Y}} = \{z \in [0, 1]^K : \sum_{y \in \mathcal{Y}} z_y = 1\}$  denotes the  
 784 standard probability simplex over  $\mathcal{Y}$ .

---

**Algorithm 2:** Multiclass bias-corrected unsupervised domain adaptation procedure.

---

**Input:** Regression function  $\eta_S : \mathcal{X} \rightarrow \Delta^{\mathcal{Y}}$ , subroutine `regressor`,  $n$  unlabeled samples  $\{(X_{S,i}, X_{U,i})\}_{i=1}^n$  from the test domain

**Output:** Estimate  $\hat{\eta}_n : \mathcal{X}_S \times \mathcal{X}_U \rightarrow \Delta^{\mathcal{Y}}$  of regression function  $\eta_y(x_S, x_U) = \Pr[Y = y | X_S = x_S, X_U = x_U]$

```

1 for  $i \in [n]$  do // generate pseudolabels
2 |   Sample  $\hat{Y}_i \sim \text{Categorical}(\eta_S(X_{S,i}))$  //  $\hat{Y} \in \{0, 1\}^{n \times K}$  is one-hot encoded
3  $\tilde{\eta}_{U,n} \leftarrow \text{regressor}(\{(X_{U,i}, \hat{Y}_i)\}_{i=1}^n)$  // regress pseudolabels over  $X_U$ 
4  $\hat{\epsilon} \leftarrow \eta_S^T(X_S) \text{Normalize}(\eta_S^T(X_S))$  // Estimate  $\epsilon_{y,y'} = \Pr[\hat{Y} = y | Y = y]$ 
5  $\hat{\eta}_{U,n} \leftarrow (x_U \mapsto \max\{0, \min\{1, \epsilon^{-1} \tilde{\eta}_{U,n}(x_U)\}, \})$  // Unstable predictor
6 for  $y \in [K]$  do
7 |    $C_y \leftarrow \left( (x_S, x_U) \mapsto \frac{\eta_{S,y}(x_S) \circ \hat{\eta}_{U,n,y}(x_U)}{\frac{1}{n} \sum_{i=1}^n \eta_{S,y}(X_{S,i})} \right)$ 
8  $\hat{\eta}_{S,U,n} \leftarrow \left( (x_S, x_U) \mapsto \frac{C(x_S, x_U)}{\|C(x_S, x_U)\|_1} \right)$  // Joint predictor
9 return  $(\hat{\eta}_{U,n}, \hat{\eta}_{S,U,n})$ 

```

---

## 785 D Supplementary Results

786 **Proposition D.1.** Suppose  $\hat{Y}|f_S(X) \sim \text{Bernoulli}(\sigma(f_S(X)))$ , such that  $\hat{Y} \perp\!\!\!\perp f_U(X)|f_S(X)$ . Then,

$$0 \in \arg \min_{f_U: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}[\ell(\hat{Y}, \sigma(f_S(X) + f_U(X)))],$$

787 where  $\ell(x, y) = -x \log y - (1 - x) \log(1 - y)$  denotes the cross-entropy loss.

788 Suppose  $\hat{Y}|f_S(X) \sim \text{Bernoulli}(\sigma(f_S(X)))$ , such that  $\hat{Y} \perp\!\!\!\perp f_U(X)|f_S(X)$ . Then,

$$\begin{aligned}
& - \mathbb{E}[\ell(\hat{Y}, \sigma(f_S(X) + f_U(X)))] \\
&= \mathbb{E}[\mathbb{E}[\ell(\hat{Y}, \sigma(f_S(X) + f_U(X)))] \quad \text{(Law of Total Expectation)} \\
&= \mathbb{E}[\mathbb{E}[\hat{Y} \log \sigma(f_S(X) + f_U(X)) \\
&\quad + (1 - \hat{Y}) \log(1 - \sigma(f_S(X) + f_U(X)) | f_S(X))] \\
&= \mathbb{E}[\mathbb{E}[\hat{Y} | f_S(X_S)] \mathbb{E}[\log \sigma(f_S(X) + f_U(X)) | f_S(X_S)] \\
&\quad + \mathbb{E}[(1 - \hat{Y}) | f_S(X_S)] \mathbb{E}[\log(1 - \sigma(f_S(X) + f_U(X)) | f_S(X))] \quad (\hat{Y} \perp\!\!\!\perp f_U(X) | f_S(X)) \\
&= \mathbb{E}[\sigma(f_S(X)) \log \sigma(f_S(X) + f_U(X)) \\
&\quad + (1 - \sigma(f_S(X))) \log(1 - \sigma(f_S(X) + f_U(X)))] \quad (\hat{Y} | f_S(X) \sim \text{Bernoulli}(\sigma(f_S(X)))).
\end{aligned}$$

789 Since the cross-entropy loss is differentiable and convex, any  $f_U(X)$  satisfying  $0 =$   
 790  $\frac{d}{df_U(X)} \mathbb{E}[\ell(\hat{Y}, \sigma(f_S(X) + f_U(X)))]$  is a minimizer. Indeed, under the mild assumption that the ex-

791 pection and derivative commute, for  $f_U(X) = 0$ ,

$$\begin{aligned} \frac{d}{df_U(X)} \mathbb{E}[\ell(\hat{Y}, \sigma(f_S(X) + f_U(X)))] &= -\mathbb{E} \left[ \frac{\sigma(f_S(X))}{\sigma(f_S(X) + f_U(X))} + \frac{1 - \sigma(f_S(X))}{1 - \sigma(f_S(X) + f_U(X))} \right] \\ &= -\mathbb{E} \left[ \frac{\sigma(f_S(X))}{\sigma(f_S(X))} + \frac{1 - \sigma(f_S(X))}{1 - \sigma(f_S(X))} \right] = 0. \end{aligned}$$

## 792 D.1 Causal Perspectives

793 The stability, complementarity, and informativeness assumptions in Theorem 4.4 can be interpreted  
794 as constraints on the causal relationships between the variables  $X_S$ ,  $X_U$ ,  $Y$ , and  $E$ . We conclude this  
795 section with a result with a characterization of causal directed acyclic graphs (DAGs) that are consis-  
796 tent with these assumptions. In particular, this result shows that our assumptions are satisfied in the  
797 “anti-causal” and “cause-effect” settings assumed in prior work [45, 59, 30], as well as work assuming  
798 only covariate shift (i.e., changes in the distribution of  $X$  without changes in the conditional  $P_{Y|X}$ ).  
799

800 **Proposition D.2** (Possible Causal DAGs). *Consider an environment*  
801 *variable  $E$ , two covariates  $X_U$  and  $X_S$ , and a label  $Y$ . Assume there*  
802 *are no other hidden confounders (i.e., causal sufficiency). First,*  
803 *assume:*

- 804 1)  $E$  is a root (i.e., none of  $X_U$ ,  $X_S$ , and  $Y$  is an ancestor of  $E$ ).
- 805 2)  $X_S$  is informative of  $Y$  (i.e.,  $X_S \not\perp\!\!\!\perp Y|E$ ).
- 806 3)  $X_S$  and  $X_U$  are complementary predictors of  $Y$ ; i.e.,  $X_S \perp\!\!\!\perp$   
807  $X_U|Y, E$ .
- 808 4)  $X_S$  is stable (i.e.,  $E \perp\!\!\!\perp Y|X_S$ ).

809 *These are the four structural assumptions under which Theorems 4.4*  
810 *and 4.5 show that the SFB algorithm learns the conditional distri-*  
811 *bution  $P_{Y|X_1, X_2}$  in the test domain. Additionally, suppose*

- 812 5)  $X_U$  is unstable (i.e.,  $E \not\perp\!\!\!\perp Y|X_U$ ), *This is the case in which*  
813 *empirical risk minimization [ERM 56] may suffer bias due to*  
814 *distribution shift, and hence when SFB may outperform ERM.*
- 815 6)  $X_U$  contains some information about  $Y$  that is not included in  $X_S$   
816 (i.e.,  $X_U \not\perp\!\!\!\perp Y|X_S$ ), *and This is information we expect invariant*  
817 *risk minimization [IRM 2] to be unable to learn, and hence when*  
818 *we expect SFB to outperform IRM.*

819 *Then, as illustrated in Figure 3, three types of stable features are*  
820 *possible:*

- 821 1. Causal ancestors  $X_{S,C}$  of  $Y$ ,
- 822 2. Causal descendants  $X_{S,E}$  of  $Y$  that are not also descendants of  $E$ ,
- 823 3. Causal spouses  $X_{S,S}$  of  $Y$  (i.e., causal ancestors of  $X_{S,E}$ ), and
- 824 *while the only unstable features possible are descendants of  $Y$ .*

825 Notable special cases of the DAG in Figure 3 include:

- 826 1. the “cause-effect” settings, studied by Rojas-Carulla et al. [45], von Kügelgen et al. [59], where  
827  $X_S$  is a cause of  $Y$ ,  $X_U$  is an effect of  $Y$ , and  $E$  affects both  $X_S$  and  $X_U$  but affects  $Y$  only through  
828  $X_S$ . Note that this generalizes the commonly used “covariate shift” assumption, as not only the  
829 covariate distribution  $P_{X_S, X_U}$  but also the conditional distribution  $P_{Y|X_U}$  can change between  
830 environments.
- 831 2. the “anti-causal” setting, studied by Jiang and Veitch [30], where  $X_S$  and  $X_U$  are both effects of  
832  $Y$ , but  $X_S$  is unaffected by  $E$ .
- 833 3. the widely studied “covariate shift” setting [53, 22, 7, 52], which corresponds (see Sections 3 and  
834 5 of Schölkopf [49]) to a causal factorization  $P(X, Y) = P(X)P(Y|X)$  (i.e., in which the only

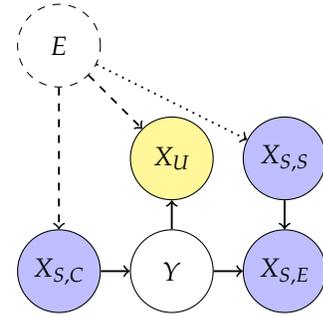


Figure 3: Causal DAGs over the environment  $E$ , three types of stable features (causes  $X_{S,C}$ , effects  $X_{S,E}$ , and spouses  $X_{S,S}$ ), unstable features  $X_U$ , and label  $Y$ , under conditions 1)-6). At least one, and possibly both, of the dashed edges  $E \rightarrow X_{S,C}$  and  $E \rightarrow X_U$  must be included. The dotted edge  $E \rightarrow X_{S,S}$  may or may not be included.

835 stable components  $X_S$  are causes  $X_{S,C}$  of  $Y$  or unconditionally independent (e.g., causal spouses  
 836  $X_{S,S}$ ) of  $Y$ .

837 However, this model is more general than these special cases. Also, for sake of simplicity, we assumed  
 838 causal sufficiency here; however, in the presence of unobserved confounders, other types of stable  
 839 features are also possible; for example, if we consider the possibility of unobserved confounders  $U$   
 840 influencing  $Y$  that are independent of  $E$  (i.e., invariant across domains), then our method can also  
 841 utilize stable features that are descendants of  $U$  (i.e., “siblings” of  $Y$ ).

## 842 E Datasets

843 In our experiments, we consider four datasets: Synthetic, ColorMNIST, PACS and Camelyon17.  
 844 While the first two offer controlled settings with a severe spurious-correlation shift, the latter two  
 845 offer real-world distribution shifts. Below, Fig. 4 depicts samples from the three image datasets.

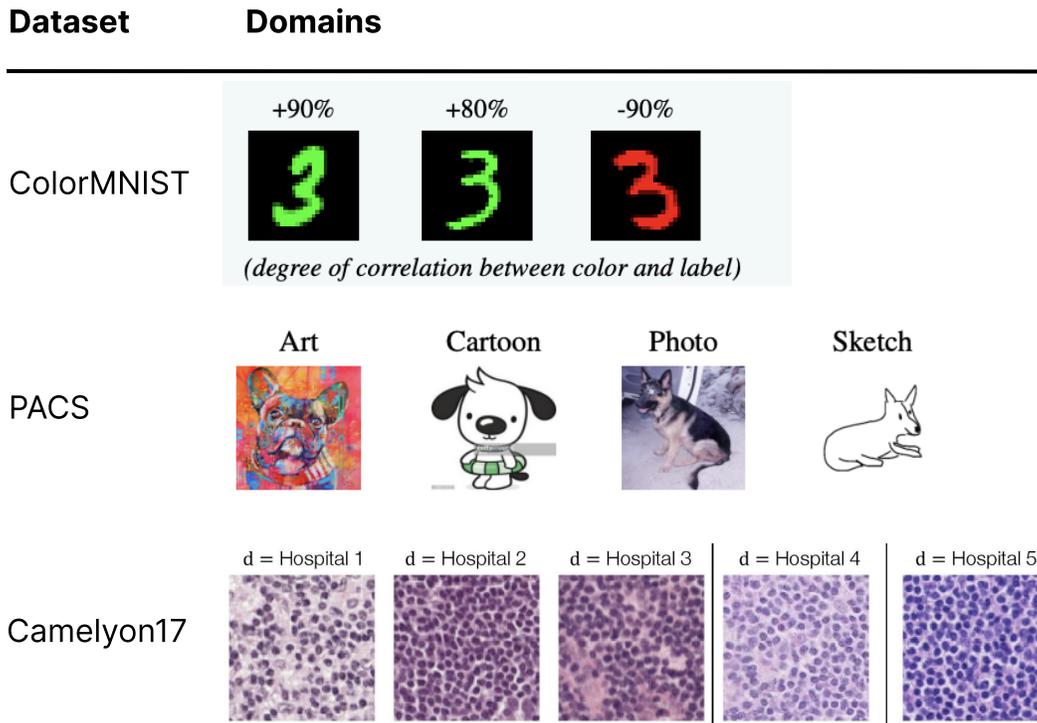


Figure 4: Examples from ColorMNIST [2], PACS [36] and Camelyon17 [5]. Figure and examples based on Gulrajani and Lopez-Paz [23, Table 3] and Koh et al. [32, Figure 4]. For ColorMNIST, we follow the standard approach [2] and use the first two domains for training and the final one for testing. For PACS [36], we follow the standard approach [23] and use each domain in turn for testing, using the remaining three domains for training. For Camelyon17 [5], we follow WILDS [32] and use the first three domains for training, the fourth for validation, and the fifth for testing.

## 846 F Further Experiments

847 This appendix provides further experiments which supplement those in the main text. In particular,  
 848 it provides: (i) experiments on a synthetic dataset where our assumption of complementarity (i.e.,  
 849 conditionally-independent unstable features) does not hold (Appendix F.1); and (ii) ablations on the  
 850 ColorMNIST dataset showing the effects of bias correction and post-hoc calibration (Appendix F.2).

### 851 F.1 Synthetic dataset

852 As depicted in Fig. 1 (right), our SFB approach assumes that the harnessed unstable features  $X_C \subseteq X_U$   
 853 are conditionally independent of the stable features  $X_S$ . If this assumption is violated, then adaptation  
 854 can fail as SFB is not guaranteed to learn an asymptotically-optimal predictor in the test domain.

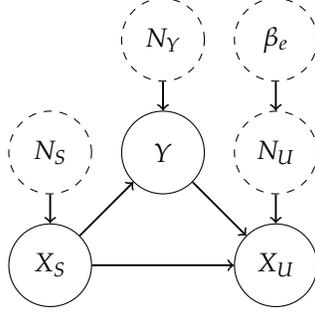


Figure 5: Causal DAG behind the synthetic dataset of Appendix F.1. Dashed circles indicate latent/unobserved variables, while solid circles indicate observed variables.

855 To investigate the adaptation performance of SFB when this assumption is violated, we conduct  
 856 experiments on a synthetic cause-effect dataset in which there is a direct dependence between  $X_S$  and  
 857  $X_U$ . In particular, similar to Jiang and Veitch [30, Appendix B], we generate synthetic data according  
 858 to the following structural equations (illustrated graphically in Fig. 5):

$$\begin{aligned}
 X_S &\leftarrow N_S, \text{ with } N_S \leftarrow \text{Bern}(0.5); \\
 Y &\leftarrow \text{XOR}(X_S, N_Y), \text{ with } N_Y \leftarrow \text{Bern}(0.75); \\
 X_U &\leftarrow \text{XOR}(\text{XOR}(Y, N_U), X_S), \text{ with } N_U \leftarrow \text{Bern}(\beta_e).
 \end{aligned}$$

859 Here, the input  $X = (X_S, X_U)$  and  $\text{Bern}(\beta)$  means that a random variable is 1 with probability  $\beta$   
 860 and 0 with probability  $1 - \beta$ . Following Jiang and Veitch [30, Appendix B], we create two training  
 861 domains with  $\beta_e \in \{0.95, 0.8\}$ , one validation domain with  $\beta_e = 0.2$ , and one test domain with  
 862  $\beta_e = 0.1$ . Like the anti-causal synthetic dataset of § 6, the idea is that prediction based on the  
 863 stable  $X_S$  results in lower accuracy (75%) than prediction based on the unstable  $X_U$ . Thus, models  
 864 optimizing for prediction accuracy only—and not stability—will use  $X_U$  and ultimately end up with  
 865 only 10% accuracy in the test domain. In addition, while the stable predictor achieves 75% accuracy  
 866 in the test domain, performance can be improved to 90% if  $X_U$  can be used correctly. However, unlike  
 867 the anti-causal synthetic dataset of § 6, the stable  $X_S$  and unstable  $X_U$  features are not conditionally  
 868 independent, i.e.,  $X_U \not\perp\!\!\!\perp X_S | Y$ , since  $X_S$  directly influences  $X_U$ . We use the same experimental  
 869 setup as for the anti-causal synthetic dataset in § 6: see Appendix G.4 for further details.

870 Looking at Table 4 we see that: (i) ACTIR has poor stable/invariant performance as its notion  
 871 of stability relies on the now-violated conditional-independence assumption; (ii) IRM has good  
 872 stable/invariant performance as its notion of stability does not rely on conditional independence; (iii)  
 873 SFB has good stable/invariant performance as its notion of stability does not rely on conditional  
 874 independence (IRM’s stability penalty is used); and (iv) surprisingly, SFB has near-optimal adapted  
 875 performance despite the conditional-independence assumption being violated. One explanation for  
 876 (iv) is that the conditional-independence assumption is only weakly violated in the test domain.  
 877 Another is that conditional independence isn’t necessary for SFB and some weaker, yet-to-be-  
 878 determined condition suffices.

Table 4: Test-domain accuracies on a synthetic cause-effect dataset with a *direct* dependence between  $X_S$  and  $X_U$ , meaning  $X_U \not\perp\!\!\!\perp X_S | Y$ . Means and standard errors are over 100 seeds.

Algorithm	Accuracy
ERM	11.57 ± 0.71
IRM	69.61 ± 1.26
ACTIR	43.51 ± 2.63
SFB w/o adapt	74.89 ± 3.64
SFB w. adapt	<b>88.56 ± 1.38</b>

879 **F.2 ColorMNIST**

880 We now provide ablations on the CoLoRmNIST dataset to illustrate the effectiveness of the different  
 881 components of SFB. In particular, we focus on bias correction and calibration, while also showing  
 882 how multiple rounds of pseudo-labelling can improve performance in practice.

883 **Bias correction.** To adapt the unstable classifier in the test domain, SFB employs the bias-corrected  
 884 adaptation algorithm of Alg. 1 (or Alg. 2 for the multi-class case) which corrects for biases caused by  
 885 possible disagreements between the stable-predictor pseudo-labels  $\hat{Y}$  and the true label  $Y$ . In this  
 886 (sub)section, we investigate the performance of SFB with and without bias correction (BC).

887 **Calibration.** As discussed in § 4.2, correctly combining the stable and unstable predictions post-  
 888 adaptation requires them to be properly calibrated. In particular, it requires the stable predictor  $f_S$  to be  
 889 calibrated with respect to the true labels  $Y$  and the unstable predictor  $f_U$  to be calibrated with respect  
 890 to the pseudo-labels  $\hat{Y}$ . In this (sub)section, we investigate the performance of SFB with and without  
 891 post-hoc calibration (in particular, simple temperature scaling [24]). More specifically, we investigate  
 892 the effect of calibrating the stable predictor (CS) and calibrating the unstable predictor (CU).

893 **Multiple rounds of pseudo-labelling.** While SFB learns the optimal unstable classifier  $h_U^e$  in  
 894 the test domain *given enough unlabelled data*, § 4.1 discussed how more accurate pseudo-labels  $\hat{Y}$   
 895 improve the sample efficiency of SFB. In particular, in a restricted-sample setting, more accurate  
 896 pseudo-labels result in an unstable classifier  $h_U^e$  which better harnesses  $X_U$  in the test domain. With  
 897 this in mind, note that, after adapting, we expect the joint predictions of SFB to be more accurate  
 898 than its stable-only predictions. This raises the question: can we use these improved predictions to  
 899 form more accurate pseudo-labels, and, in turn, an unstable classifier  $h_U^e$  that leads to even better  
 900 performance? Furthermore, can we repeat this process, using multiple rounds of pseudo-labelling to  
 901 refine our pseudo-labels and ultimately  $h_U^e$ ? While this multi-round approach loses the asymptotic  
 902 guarantees of § 4.2, we found it to work quite well in practice. In this (sub)section, we thus investigate  
 903 the performance of SFB with and without multiple rounds of pseudo-labelling (PL rounds).

Table 5: SFB ablations on CoLoRmNIST. Means and standard errors are over 3 random seeds. *BC*: bias correction. *CS*: post-hoc calibration of the stable classifier. *CU*: post-hoc calibration of the unstable classifier. *PL Rounds*: Number of pseudo-labelling rounds used. *GT adapt*: adapting using true labels in the test domain.

Model	Bias Correction	Calibration Stable	Calibration Unstable	PL Rounds	Test Acc.
SFB w/o adapt				1	70.6 ± 1.8
SFB with adapt				1	78.0 ± 2.9
+BC	✓			1	83.4 ± 2.8
+CS		✓		1	80.6 ± 3.4
+CU			✓	1	76.6 ± 2.4
+BC+CS+CU	✓	✓	✓	1	84.4 ± 2.2
+BC+CS	✓	✓		1	84.9 ± 2.6
+BC+CS	✓	✓		2	87.4 ± 1.9
+BC+CS	✓	✓		3	88.1 ± 1.8
+BC+CS	✓	✓		4	88.6 ± 1.3
+BC+CS	✓	✓		5	88.7 ± 1.3
SFB with GT adapt	✓	✓		1	89.0 ± 0.3

904 **Results.** Table 5 reports the ablations of SFB on CoLoRmNIST. Here we see that: (i) bias correction  
 905 significantly boosts performance (+BC); (ii) calibrating the stable predictor also boosts performance  
 906 without (+CS) and with (+BC+CS) bias correction, with the latter leading to the best performance;  
 907 (iii) calibrating the unstable predictor (with respect to the pseudo-labels) slightly hurts performance  
 908 without (+CU) and with (+BC+CS+CU) bias correction and stable-predictor calibration; (iv) multiple  
 909 rounds of pseudo-labelling boosts performance, while also reducing the performance variation across  
 910 random seeds; (v) using bias correction, stable-predictor calibration and 5 rounds of pseudo-labelling

911 results in near-optimal adaptation performance, as indicated by the similar performance of SFB when  
912 using true labels  $Y$  to adapt  $h_U^e$  (denoted “SFB with GT adapt” in Table 5).

## 913 G Implementation Details

914 Below we provide further implementation details for each of the experiments/datasets considered in  
915 this work. Code for reproducing all experimental results will be made available upon acceptance.

### 916 G.1 ColorMNIST

917 **Training details.** We follow the setup of Eastwood et al. [14, §6.1] and build on their open-source  
918 code<sup>5</sup>. In particular, we use the original MNIST training set to create training and validation sets  
919 for each domain, and the original MNIST test set for the test sets of each domain. For all methods,  
920 we use a 2-hidden-layer MLP with 390 hidden units, the Adam optimizer, a learning rate of 0.0001  
921 with cosine scheduling, and dropout with  $p = 0.2$ . In addition, we use full batches (size 25000),  
922 400 steps for ERM pertaining (which directly corresponds to the delicate penalty “annealing” or  
923 warm-up periods used by penalty-based methods on CoLoRMNIST [2, 34, 14]), and 600 total steps.  
924 We sweep over stability-penalty weights in  $\{50, 100, 500, 1000, 5000\}$  for IRM, VREx and SFB  
925 and  $\alpha$ ’s in  $1 - \{e^{-100}, e^{-250}, e^{-500}, e^{-750}, e^{-1000}\}$  for EQRM. As the stable (shape) and unstable  
926 (color) features are conditionally independent given the label, we fix SFB’s conditional-independence  
927 penalty weight  $\lambda_C = 0$ . As is the standard for CoLoRMNIST, we use a test-domain validation set to  
928 select the best settings (after the total number of steps), and then report the mean and standard error  
929 over 10 random seeds on a test-domain test set. As in previous works, the hyperparameter ranges of  
930 all methods are selected by peeking at test-domain performance. While far from ideal, this is quite  
931 difficult to avoid with CoLoRMNIST and highlights a core problem with hyperparameter selection in  
932 DG—as discussed by many previous works [2, 34, 23, 64, 14].

933 **Adaptation details.** For SFB’s unsupervised adaptation in the test domain, we use a batch size of  
934 2048 and employ the bias correction of Alg. 1. In addition, we calibrate the stable predictor using post-  
935 hoc temperature scaling, choosing the temperature to minimize the expected calibration error (ECE,  
936 [24]) across the two training domains. Again using the two training domains for hyperparameter  
937 selection, we sweep over adaptation learning rates in  $\{0.1, 0.01\}$ , choose the best adaptation step in  
938  $[5, 20]$  (via early stopping), and sweep over the number of pseudo-labelling rounds in  $[1, 5]$ . Finally,  
939 we report the mean and standard error over 3 random seeds for adaptation.

### 940 G.2 PACS

941 We follow the experimental setup of Jiang and Veitch [30, Section 6.4] and build on their open-source  
942 implementation<sup>6</sup>. This means using an ImageNet-pretrained ResNet-18, the Adam optimizer with a  
943 learning rate of  $10^{-4}$ , and, following [23], choosing hyperparameters using leave-one-domain-out  
944 cross-validation. This is akin to K-fold cross-validation except with domains, meaning that we  
945 train 3 models—each time leaving out 1 of the 3 training domains for validation—and then select  
946 hyperparameters based on the best average performance across the held-out validation domains.  
947 Finally, we use the selected hyperparameters to retrain the model using all 3 training domains.

948 For SFB, we sweep over  $\lambda_S$  in  $\{0.01, 0.1, 1, 5, 10, 20\}$ ,  $\lambda_C$  in  $\{0.01, 0.1, 1\}$ , and learning rates in  
949  $\{10^{-4}, 50^{-4}\}$ . For SFB’s unsupervised adaptation, we employ the multi-class bias correction of  
950 Alg. 2 and calibrate the stable predictor using post-hoc temperature scaling, choosing the temperature  
951 to minimize the expected calibration error (ECE, [24]) across the three training domains. In addition,  
952 we use the Adam optimizer with an adaptation learning rate of 0.01, choosing the number of adaptation  
953 steps in  $[1, 20]$  (via early stopping) using the training domains. Finally, we report the mean and  
954 standard error over 3 random seeds.

### 955 G.3 Camelyon17

956 We follow the experimental setup of Jiang and Veitch [30, Section 6.3] and build on their open-source  
957 implementation<sup>7</sup>. This means using an ImageNet-pretrained ResNet-18, the Adam optimizer, and,  
958 following [32], choosing hyperparameters using the validation domain (hospital 4). In contrast to

<sup>5</sup><https://github.com/cianeastwood/qrm/tree/main/CMNIST>

<sup>6</sup><https://github.com/ybjiaang/ACTIR>.

<sup>7</sup>See Footnote 6.

959 [30], we use a learning rate of  $10^{-5}$  for all methods, rather than  $10^{-4}$ , and employ early stopping  
960 using the validation domain. We found this to significantly improve all methods. E.g., the baselines  
961 of ERM and IRM improve by approximately 20 percentage points, jumping from  $\approx 70\%$  to  $\approx 90\%$ .

962 For SFB, we sweep over  $\lambda_S$  in  $\{0.01, 0.1, 1, 5, 10, 20\}$  and  $\lambda_C$  in  $\{0.01, 0.1, 1\}$ . For SFB’s unsuper-  
963 vised adaptation, we employ the bias correction of Alg. 1 and calibrate the stable predictor using  
964 post-hoc temperature scaling, choosing the temperature to minimize the expected calibration er-  
965 ror (ECE, [24]) on the validation domain. In addition, we use the Adam optimizer with an adaptation  
966 learning rate of 0.01, choosing the number of adaptation steps in  $[1, 20]$  (via early stopping) using the  
967 validation domain. Finally, we report the mean and standard error over 3 random seeds.

#### 968 G.4 Synthetic

969 Following Jiang and Veitch [30], we use a simple three-layer network with 8 units in each hidden  
970 layer and the Adam optimizer, choosing hyperparameters using the validation domain.

971 For SFB, we sweep over  $\lambda_S$  in  $\{0.01, 0.1, 1, 5, 10, 20\}$  and  $\lambda_C$  in  $\{0.01, 0.1, 1\}$ . For SFB’s unsuper-  
972 vised adaptation, we employ the bias correction of Alg. 1 and calibrate the stable predictor using  
973 post-hoc temperature scaling, choosing the temperature to minimize the expected calibration er-  
974 ror (ECE, [24]) on the validation domain. In addition, we use the Adam optimizer with an adaptation  
975 learning rate of 0.01, choosing the number of adaptation steps in  $[1, 20]$  (via early stopping) using the  
976 validation domain. Finally, we report the mean and standard error over 100 random seeds.

## 977 H Further Related Work

978 **Using spurious or unstable features without labels.** Bui et al. [12] exploit-domain specific or  
979 unstable features with a meta-learning approach. However, they use the unstable features *in the same*  
980 *way* in the test domain, which, by their very definition, can lead to degraded performance. In contrast,  
981 we seek a *robust* approach to safely harness the unstable features in the test domain, as summarised  
982 in Table 1. Sun et al. [54] share the goal of exploiting spurious or unstable features to go “beyond  
983 invariance”. However, their approach requires labels for the spurious features at training time and  
984 only applies to label shifts. In contrast, we do not require labels for the spurious features and are not  
985 restricted to label shifts.

986 **Self-learning via pseudo-labelling.** In the source-free and test-time domain adaptation literature,  
987 adapting to the test domain using a model’s own pseudo-labels is a common approach [35, 38, 61, 29]—  
988 see Rusak et al. [48] for a recent review. In contrast to these approaches, we use one model to provide  
989 the pseudo-labels (the stable model) and the other to use/adapt to the pseudo-labels (the unstable  
990 model). In addition, while the majority of this pseudo-labelling work is purely empirical, we provide  
991 theoretical justification and guarantees for our SFB approach.

## 992 I Limitations

993 In our view, the most significant limitation of this work is the assumption of complementarity (i.e.,  
994 that the spurious features are conditionally independent of the stable features, given the label).  
995 Complementarity is implicit in the causal generative models assumed by existing related work [45,  
996 60, 30], and, as Example A.1 in Appendix A.1 demonstrates, is cannot simply be dropped from our  
997 theoretical motivation. In the related context of co-training, this condition was initially assumed and  
998 then weakened in subsequent work [10, 4, 1, 62]; similarly, we hope future work will identify weaker  
999 conditions that are sufficient for SFB to succeed. On the other hand, our experimental results on the  
1000 synthetic dataset of Appendix F.1, as well as the real datasets of PACS and Camelyon17, suggest  
1001 that SFB may be robust to violations of complementarity—perhaps mirroring the surprisingly good  
1002 practical performance of methods such as naive Bayes classification which are justified under similar  
1003 assumptions [44].

1004