
Don't fear the unlabelled: Safe semi-supervised learning via simple debiasing

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Semi-supervised learning (SSL) provides an effective means of leveraging un-
2 labelled data to improve a model's performance. Even though the domain has
3 received a considerable amount of attention in the past years, most methods present
4 the common drawback of lacking theoretical guarantees. Our starting point is to
5 notice that the estimate of the risk that most discriminative SSL methods minimise
6 is biased, even asymptotically. This bias impedes the use of standard statistical
7 learning theory and can hurt empirical performance. We propose a simple way of
8 removing the bias. Our debiasing approach is straightforward to implement and
9 applicable to most deep SSL methods. We provide simple theoretical guarantees on
10 the trustworthiness of these modified methods, without having to rely on the strong
11 assumptions on the data distribution that SSL theory usually requires. In particular,
12 we provide generalisation error bounds for the proposed methods. We evaluate
13 debiased versions of different existing SSL methods, such as the Pseudo-label
14 method and Fixmatch, and show that debiasing can compete with classic deep SSL
15 techniques in various settings by providing better calibrated models. Additionally,
16 we provide a theoretical explanation of the intuition of the popular SSL methods.

17 1 Introduction

18 The promise of semi-supervised learning (SSL) is to be able to learn powerful predictive models
19 using partially labelled data. In turn, this would allow machine learning to be less dependent on
20 the often costly and sometimes dangerously biased task of labelling data. Early SSL approaches—
21 e.g. Scudder's (1965) untaught pattern recognition machine—simply replaced unknown labels with
22 predictions made by some estimate of the predictive model and used the obtained *pseudo-labels* to
23 refine their initial estimate. Other more complex branches of SSL have been explored since, notably
24 using generative models (from McLachlan, 1977, to Kingma et al., 2014) or graphs (notably following
25 Zhu et al., 2003). Deep neural networks, which are state-of-the art supervised predictors, have been
26 trained successfully using SSL. Somewhat surprisingly, the main ingredient of their success is still the
27 notion of pseudo-labels (or one of its variants), combined with systematic use of data augmentation
28 (e.g. Xie et al., 2019; Sohn et al., 2020; Rizve et al., 2021).

29 An obvious SSL baseline is simply throwing away the unlabelled data. We will call such a baseline the
30 *complete case*, following the missing data literature (e.g. Tsiatis, 2006). As reported in van Engelen &
31 Hoos (2020), the main risk of SSL is the potential degradation caused by the introduction of unlabelled
32 data. Indeed, semi-supervised learning outperforms the complete case baseline only in specific cases
33 (Singh et al., 2008; Schölkopf et al., 2012; Li & Zhou, 2014). This degradation risk for generative
34 models has been analysed in Chapelle et al. (2006, Chapter 4). To overcome this issue, previous works
35 introduced the notion of *safe* semi-supervised learning for techniques which never reduce predictive
36 performance by introducing unlabelled data (Li & Zhou, 2014; Guo et al., 2020). Our loose definition

37 of safeness is as follows: *an SSL algorithm is safe if it has theoretical guarantees that are similar*
 38 *or stronger to the complete case baseline.* The “theoretical” part of the definition is motivated by the
 39 fact that any empirical assessment of generalisation performances of an SSL algorithm is jeopardised
 40 by the scarcity of labels. Unfortunately, popular deep SSL techniques generally do not benefit from
 41 theoretical guarantees without strong and essentially untestable assumptions on the data distribution
 42 (Mey & Loog, 2019) such as the smoothness assumption (small perturbations on the features x do not
 43 cause large modification in the labels, $p(y|pert(x)) \approx p(y|x)$) or the cluster assumption (data points
 44 are distributed on discrete clusters and points in the same cluster are likely to share the same label).

45 Most semi-supervised methods rely on these dis-
 46 tributional assumptions to ensure performance
 47 in entropy minimisation, pseudo-labelling and
 48 consistency-based methods. However, no proof
 49 is given that guarantees the effectiveness of state-
 50 of-the-art methods (Tarvainen & Valpola, 2017;
 51 Miyato et al., 2018; Sohn et al., 2020; Pham
 52 et al., 2021). To illustrate that SSL requires spe-
 53 cific assumptions, we show in a toy example that
 54 pseudo-labelling fails at learning. To do so, we
 55 draw samples from two uniform distributions
 56 with a small overlap. Both supervised and semi-
 57 supervised neural networks are trained using the
 58 same labelled dataset. While the supervised algo-
 59 rithm learns perfectly the true distribution of
 60 $p(1|x)$, the semi-supervised learning methods
 61 (both entropy minimisation and pseudo-label)
 62 underestimate $p(1|x)$ for $x \in [1, 3]$ (see Figure
 63 1). We also test our proposed method (DeSSL)
 64 on this dataset and show that the unbiased ver-
 65 sion of each SSL technique learns the true dis-
 66 tribution accurately. See Appendix A for the
 67 results with Entropy Minimisation.

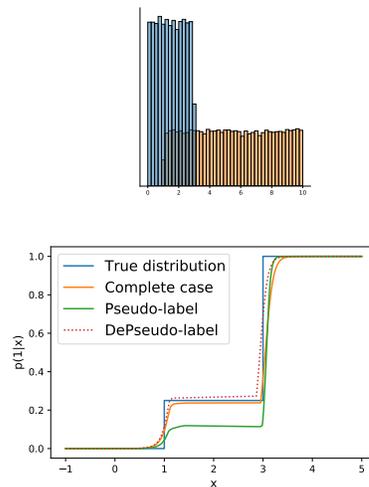


Figure 1: (Left) Data histogram. (Right) Posterior probabilities $p(1|x)$ of the same model trained following either complete case (only labelled data), Pseudo-label or our DePseudo-label.

68 1.1 Contributions

69 Rather than relying on the strong geometric assumptions usually used in SSL theory, we simply use
 70 the *missing completely at random (MCAR)* assumption, a standard assumption from the missing data
 71 literature (see e.g. Little & Rubin, 2019). With this only assumption on the data distribution, we
 72 propose a new safe SSL method derived from simply debiasing common SSL risk estimates. Our
 73 main contributions are:

- 74 • We introduce debiased SSL (DeSSL), a safe method that can be applied to most deep SSL
 75 algorithms without assumptions on the data distribution;
- 76 • We propose a theoretical explanation of the intuition of popular SSL methods. We provide
 77 theoretical guarantees on the safeness of using DeSSL both on consistency and calibration
 78 of the method. We also provide a generalisation error bound;
- 79 • We show how simple it is to apply DeSSL to the most popular methods such as Pseudo-label
 80 and Fixmatch, and show empirically that DeSSL leads to models that are never worse than
 81 their classical counterparts, generally better calibrated and sometimes much more accurate.

82 2 Semi-supervised learning

83 2.1 Learning with labelled data

84 The ultimate objective of most of the learning frameworks is to minimise a risk \mathcal{R} , defined as
 85 the expectation of a particular loss function L over a data distribution $p(x, y)$, on a set of models
 86 $f_\theta(x)$, parametrised by $\theta \in \Theta$. Thus, the learning task is finding θ^* that minimises the risk:
 87 $\mathcal{R}(\theta) = \mathbb{E}_{(X,Y) \sim p(x,y)} [L(\theta; X, Y)]$. The distribution $p(x, y)$ being unknown, we generally minimise

88 an approximation of the risk, the empirical risk $\hat{\mathcal{R}}(\theta)$ computed on a sample of n i.i.d points drawn
 89 from $p(x, y)$. $\hat{\mathcal{R}}(\theta)$ is an unbiased and consistent estimate of $\mathcal{R}(\theta)$ under mild assumptions. Its
 90 unbiased nature is one of the basic properties that is used for the development of traditional learning
 91 theory and asymptotic statistics (van der Vaart, 2000; Shalev-Shwartz & Ben-David, 2014).

92 2.2 Learning with both labelled and unlabelled data

93 Semi-supervised learning leverages both labelled and unlabelled data to improve the model’s per-
 94 formance and generalisation. Further information on the distribution $p(x)$ provides a better under-
 95 standing of the distributions $p(x, y)$ and also $p(y|x)$. Indeed, $p(x)$ may contain information on $p(y|x)$
 96 (Schölkopf et al., 2012, Goodfellow et al., 2016, Chapter 7.6, van Engelen & Hoos, 2020).

97 In the following, we have access to n samples drawn from the distribution $p(x, y)$ where some of the
 98 labels are missing. We introduce a new random variable $r \in \{0, 1\}$ that governs whether or not a
 99 data point is labelled ($r = 0$ missing, $r = 1$ observed). We note n_l the number of labelled and n_u the
 100 number of unlabelled datapoints. The MCAR assumption states that the missingness of a label y is
 101 independent of its features and the value of the label: $p(x, y, r) = p(x, y)p(r)$, then $r \sim \mathcal{B}(\pi)$. This
 102 is the case when nor features nor label carry information about the potential missingness of the labels.
 103 This description of semi-supervised learning as a missing data problem has already been done in
 104 multiple works –e.g. Seeger, 2000; Ahfock & McLachlan, 2019. Moreover, the MCAR assumption
 105 is implicitly made in most of the SSL works to design the experiments, indeed, missing labels are
 106 drawn completely as random in datasets such as MNIST, CIFAR or SVHN (Tarvainen & Valpola,
 107 2017; Miyato et al., 2018; Xie et al., 2019; Sohn et al., 2020).

108 2.2.1 Complete case: throwing the unlabelled data away

109 In missing data theory, the complete case is the learning scheme that only uses fully observed
 110 instances, namely labelled data. The natural estimator of the risk is then simply the empirical risk
 111 computed on the labelled data. Fortunately, in the MCAR setting, the complete case risk estimate
 112 keeps the same good properties of the traditional supervised one: it is unbiased and converges
 113 pointwisely to $\mathcal{R}(\theta)$. Therefore, traditional learning theory holds for the complete case under MCAR.
 114 While these observations are hardly new (see e.g. Liu & Goldberg, 2020), they can be seen as
 115 particular cases of the theory that we develop below. The risk to minimise is

$$\hat{\mathcal{R}}_{CC}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i). \quad (1)$$

116 2.2.2 Incorporating unlabelled data

117 A major drawback of the complete case framework is that a lot of data ends up not being exploited. A
 118 class of SSL approaches, mainly inductive methods with respect to the taxonomy of van Engelen &
 119 Hoos (2020), generally aim to minimise a modified estimator of the risk by including unlabelled data.
 120 Therefore, the optimisation problem generally becomes finding $\hat{\theta}$ that minimises the SSL risk,

$$\hat{\mathcal{R}}_{SSL}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i). \quad (2)$$

121 where H is a term that does not depend on the labels and λ is a scalar weight which balances the
 122 labelled and unlabelled terms. In the literature, H can generally be seen as a surrogate of L . Indeed,
 123 it looks like the intuitive choices of H are equal or equivalent to a form of expectation of L on a
 124 distribution given by the model.

125 2.2.3 Some examples of surrogates

126 A recent overview of the recent SSL techniques has been proposed by van Engelen & Hoos (2020).
 127 In this work, we focus on methods suited for a discriminative probabilistic model $p_\theta(y|x)$ that
 128 approximates the conditional $p(y|x)$. We categorised methods into two distinct sections, entropy
 129 and consistency-based.

130 **Entropy-based methods** Entropy-based methods aim to minimise a term of entropy of the predic-
 131 tions computed on unlabelled data. Thus, they encourage the model to be confident on unlabelled
 132 data, implicitly using the cluster assumption. Entropy-based methods can all be described as an
 133 expectation of L under a distribution π_x computed at the datapoint x :

$$H(\theta; x) = \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})}[L(\theta; \tilde{x}, \tilde{y})]. \quad (3)$$

134 For instance, Grandvalet & Bengio (2004) simply use the Shannon entropy as $H(\theta; x)$ which can be
 135 rewritten as equation (3) with $\pi_x(\tilde{x}, \tilde{y}) = \delta_x(\tilde{x})p_\theta(\tilde{y}|\tilde{x})$, where δ_x is the dirac distribution in x . Also,
 136 pseudo-label methods, which consist in picking the class with the maximum predicted probability
 137 as a pseudo-label for the unlabelled data (Scudder, 1965), can also be described as Equation 3. See
 138 Appendix B for complete description of the entropy-based literature (Berthelot et al., 2019; 2020;
 139 Xie et al., 2019; Sohn et al., 2020; Rizve et al., 2021; Zhang et al., 2021a) and further details.

140 **Consistency-based methods** Another range of SSL methods minimise a consistency objective
 141 that encourages invariant prediction for perturbations either on the data either on the model in order
 142 to enforce stability on model predictions. These methods rely on the smoothness assumption. In
 143 this category, we cite Π -model from (Sajjadi et al., 2016), temporal ensembling from (Laine & Aila,
 144 2017), Mean-teacher proposed by (Tarvainen & Valpola, 2017), virtual adversarial training (VAT)
 145 from (Miyato et al., 2018) and interpolation consistent training (ICT) from (Verma et al., 2019). We
 146 remark that these objectives H are equivalent to an expectation of L (see Appendix B). The general
 147 form of the unsupervised objective can be written as

$$C_1 \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})}[L(\theta; \tilde{x}, \tilde{y})] \leq H(\theta; x) = \mathbf{Div}(f_{\hat{\theta}}(x, \cdot), \text{pert}(f_{\hat{\theta}}(x, \cdot))) \leq C_2 \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})}[L(\theta; \tilde{x}, \tilde{y})], \quad (4)$$

148 where $f_{\hat{\theta}}$ is the predictions of the model, the \mathbf{Div} is a non-negative function that measures the
 149 divergence between two distributions, $\hat{\theta}$ is a fixed copy of the current parameter θ (the gradient is not
 150 propagated through $\hat{\theta}$), pert is a perturbation applied to the model or the data and $0 \leq C_1 \leq C_2$.

151 Previous works also remarked that H is an expectation of L for entropy-minimisation and pseudo-
 152 label (Zhu et al., 2022; Aminian et al., 2022). We describe a more general framework covering further
 153 methods and provide with our theory an intuition on the choice of H .

154 2.3 Theoretical guarantees

155 The main risk of SSL is the potential degradation caused by the introduction of unlabelled data when
 156 distributional assumptions are not satisfied (Singh et al., 2008; Schölkopf et al., 2012; Li & Zhou,
 157 2014), specifically in settings where the MCAR assumption does not hold anymore (Oliver et al.,
 158 2018; Guo et al., 2020). Additionally, in (Zhu et al., 2022), the authors show disparate impacts of
 159 pseudo-labelling on the different sub-classes of the population. To mitigate these problems, previous
 160 works introduced the notion *safe* semi-supervised learning for techniques which never reduce learning
 161 performance by introducing unlabelled data (Li & Zhou, 2014; Kawakita & Takeuchi, 2014; Li et al.,
 162 2016; Gan et al., 2017; Trapp et al., 2017; Guo et al., 2020). As remark by Oliver et al. (2018),
 163 SSL performances are enabled by leveraging large validation sets which is not suited for real-world
 164 applications. Then, theoretical guarantees are required to use safely SSL algorithms. For this reason,
 165 in our work, we consider as *safe* an SSL algorithm that has theoretical guarantees that are similar
 166 or stronger than those of the complete case baseline. Even though the methods presented above
 167 produce good performances in a variety of SSL benchmarks, they generally do not benefit from
 168 theoretical guarantees, even elementary. More over, Schölkopf et al. (2012) identify settings on the
 169 causal relation between the features x and the target y where SSL may systematically fail, even if
 170 classic SSL assumptions hold. Our example of Figure 1 also shows that classic SSL may fail to
 171 generalise in a very benign setting with a large number of labelled data.

172 Presented methods minimise a biased version of the risk under the MCAR assumption and therefore
 173 classical learning theory cannot be applied anymore, as we argue more precisely in Appendix C.
 174 Learning over a biased estimate of the risk is not necessarily unsafe but it is difficult to provide
 175 theoretical guarantees on such methods even if some works try to do so with strong assumptions
 176 on the data distribution (Mey & Loog 2019, Section 4 and 5). Additionally, we remark that the
 177 choice of H can be confusing as seen in the literature. For instance, Grandvalet & Bengio (2004) and
 178 Corduneanu & Jaakkola (2003) perform respectively entropy and mutual information *minimisation*
 179 whereas Pereyra et al. (2017) and Krause et al. (2010) perform *maximisation* of the same quantities.

180 **2.4 Related works**

181 Previous works already proposed safe SSL methods with theoretical guarantees. Unfortunately,
 182 so far these methods come with either strong assumptions or important computational burdens.
 183 Li & Zhou (2014) introduced a safe semi-supervised SVM and showed that the accuracy of their
 184 method is never worse than SVMs trained with only labelled data with the assumption that the true
 185 model is accessible. However, if the distributional assumptions are not satisfied, no improvement or
 186 degeneration is expected. Sakai et al. (2017) proposed an unbiased estimate of the risk for binary
 187 classification by including unlabelled data. The key idea is to use unlabelled data to better evaluate
 188 on the one hand the risk of positive class samples and on the other the risk of negative samples.
 189 They provided theoretical guarantees on its variance and a generalisation error bound. The method
 190 is designed only for binary classification and has not been tested in a deep learning setting. It has
 191 been extended to ordinal regression in follow-up work (Tsuchiya et al., 2021). In the context of
 192 kernel machines, Liu & Goldberg (2020) used an unbiased estimate of risk, like ours, for a specific
 193 choice of H . Guo et al. (2020) proposed DS^3L , a safe method that needs to approximately solve
 194 a bi-level optimisation problem. In particular, the method is designed for a different setting, not
 195 under the MCAR assumption, where there is a class mismatch between labelled and unlabelled data.
 196 The resolution of the optimisation problem provides a solution not worse than the complete case but
 197 comes with approximations. They provide a generalisation error bound. Also, the method does not
 198 outperform classic SSL methods in the MCAR setting as it is designed for non-MCAR situations.
 199 Sokolovska et al. (2008) proposed a safe method with strong assumptions such that the feature space
 200 is finite and the marginal probability distribution of x is fully known. Fox-Roberts & Rosten (2014)
 201 proposed an unbiased estimator in the generative setting applicable to a large range of models and
 202 they prove that this estimator has a lower variance than the one of the complete case.

203 **3 DeSSL: Unbiased semi-supervised learning**

204 To overcome the issues introduced by the second term in the approximation of the risk for the semi-
 205 supervised learning approach, we propose DeSSL, an unbiased version of the SSL estimator using
 206 labelled data to annul the bias. The idea here is to retrieve the properties of classical learning theory.
 207 Fortunately, we will see that the proposed method can eventually have better properties than the complete
 208 case, in particular with regard to the variance of the estimate. The proposed DeSSL objective is

$$\hat{\mathcal{R}}_{DeSSL}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i). \quad (5)$$

209

210 Under the MCAR assumption, this estimator is unbiased for any value of the parameter λ . For proof
 211 of this result see Appendix D. **We prove the optimality of debiasing with the labelled data in Appendix**
 212 **F.**

213 Intuitively, for entropy-based methods, H should be applied only on unlabelled data to enforce the
 214 confidence of the model only on unlabelled datapoints. Whereas, for consistency-based methods,
 215 H can be applied to any subset of data points. Our theory and proposed method remain the same
 216 whether H is applied to all the available data or not (see Appendix K).

217 **3.1 Does the DeSSL risk estimator make sense?**

218 The most intuitive interpretation is that by debiasing the risk estimator, we get back to the basics of
 219 learning theory. This way of debiasing is closely related to the method of control variates (Owen,
 220 2013, Chapter 8) which is a common variance reduction technique. The idea is to add an additional
 221 term to a Monte-Carlo estimator with a null expectation in order to reduce the variance of the
 222 estimator without modifying the expectation. Here, DeSSL can also be interpreted as a control variate
 223 on the risk’s gradient itself and should improve the optimisation scheme. This idea is close to the
 224 optimisation schemes introduced by Johnson & Zhang (2013) and Defazio et al. (2014) which reduce
 225 the variance of the gradients’ estimate to improve optimisation performance.

226 Another interesting way to interpret DeSSL is as a constrained optimisation problem. Indeed, min-
 227 imising $\hat{\mathcal{R}}_{DeSSL}$ is equivalent to minimising the Lagrangian of the following optimisation problem:

$$\begin{aligned} \min_{\theta} \quad & \hat{\mathcal{R}}_{CC}(\theta) \\ \text{s.t.} \quad & \frac{1}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) = \frac{1}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i). \end{aligned} \quad (6)$$

228 The idea of this optimisation problem is to minimise the complete case risk estimator by assessing
 229 that some properties represented by H are on average equal for the labelled data and the unlabelled
 230 data. For example, if we consider entropy-minimisation, this program encourages the model to have
 231 the same confidence on the unlabelled examples as on the labelled ones.

232 The debiasing term of our objective will penalise the confidence of the model on the labelled data.
 233 Pereyra et al. (2017) show that penalising the entropy in a supervised context acts as a strong
 234 regulator for supervised models and improves on the state-of-the-art on common benchmarks. This
 235 comforts us in the idea of debiasing using labelled data in the case of entropy-minimisation. Similarly,
 236 the debiasing term in pseudo-label turns the problem into plausibility inference as described by
 237 Barndorff-Nielsen (1976). Our objective also resembles doubly-robust risk estimates used for SSL in
 238 the context of kernel machines by Liu & Goldberg (2020) and for deep learning in a recent preprint
 239 (Hu et al., 2022). In both cases, their focus is quite different, as they consider weaker conditions
 240 than MCAR, but very specific choices of H .

241 3.2 Is $\hat{\mathcal{R}}_{DeSSL}(\theta)$ an accurate risk estimate?

242 Because of the connections between our debiased estimate and variance reduction techniques, we
 243 have a natural interest in the variance of the estimate. Having a lower-variance estimate of the risk
 244 would mean estimating it more accurately, leading to better models. Similarly to traditional control
 245 variates (Owen, 2013), the variance can be computed, and optimised in λ :

246 **Theorem 3.1.** *The function $\lambda \mapsto \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta))$ reaches its minimum for:*

$$\lambda_{opt} = \frac{n_u}{n} \frac{\text{Cov}(L(\theta; x, y), H(\theta; x))}{\mathbb{V}(H(\theta; x))}, \quad (7)$$

247 and at λ_{opt} :

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta))|_{\lambda_{opt}} = \left(1 - \frac{n_u}{n} \rho_{L,H}^2\right) \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)) \leq \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)), \quad (8)$$

248 where $\rho_{L,H} = \text{Corr}(L(\theta; x, y), H(\theta; x))$.

249 **Additionally, $\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)) \leq \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta))$ for all λ between 0 and $2\lambda_{opt}$.** A proof of this theorem
 250 is available as Appendix E. This theorem provides a formal justification to the heuristic idea that
 251 H should be a surrogate of L . Indeed, DeSSL is a more accurate risk estimate when H is strongly
 252 positively correlated with L , which is likely to be the case when H is equal or equivalent to an
 253 expectation of L . Then, choosing λ positive is a coherent choice. We also demonstrate in Appendix E
 254 that L and H are positively correlated when L is the negative likelihood and H is the entropy. Other
 255 SSL methods have variance reduction guarantees and already has shown great promises in SSL, see
 256 Fox-Roberts & Rosten (2014) and Sakai et al. (2017). In a purely supervised context, Chen et al.
 257 (2020) show that the effectiveness of data augmentation techniques lays partially on the variance
 258 reduction of the risk estimate. A natural application of this theorem would be to tune λ automatically
 259 by estimating λ_{opt} . In our case however, the estimation of $\text{Cov}(L(\theta; x, y), H(\theta; x))$ with few labels
 260 led to extremely unstable unsatisfactory results. **However, we estimate it more accurately using the
 261 test set (which is of course impossible in practice) on different datasets and methods to provide
 262 intuition on the order of λ_{opt} and the range of the variance reduction regime in Appendix M.2.**

263 3.3 Calibration

264 The calibration of a model is its capacity of predicting probability estimates that are representative
 265 of the true distribution. This property is determinant in real-world application when we need

266 reliable predictions. A scoring rule \mathcal{S} is a function assigning a score to the predictive distribution
 267 $p_\theta(y|x)$ relative to the event $y|x \sim p(y|x)$, $\mathcal{S}(p_\theta, (x, y))$, where $p(x, y)$ is the true distribution (see
 268 e.g. Gneiting & Raftery, 2007). A scoring rule measures both the accuracy and the quality of
 269 predictive uncertainty, meaning that better calibration is rewarded. The expected scoring rule is
 270 defined as $\mathcal{S}(p_\theta, p) = \mathbb{E}_p[\mathcal{S}(p_\theta, (x, y))]$. A proper scoring rule is defined as a scoring rule such
 271 that $\mathcal{S}(p_\theta, p) \leq \mathcal{S}(p, p)$ (Gneiting & Raftery, 2007). The motivation behind having proper scoring
 272 rules comes from the following: suppose that the true data distribution p is accessible by our
 273 set of models. Then, the scoring rule encourages to predict $p_\theta = p$. The opposite of a proper
 274 scoring rule can then be used to train a model to encourage the calibration of predictive uncertainty:
 275 $L(\theta; x, y) = -\mathcal{S}(p_\theta, (x, y))$. Most common losses used to train models are proper scorings rule such
 276 as log-likelihood.

277 **Theorem 3.2.** *If $\mathcal{S}(p_\theta, (x, y)) = -L(\theta; x, y)$ is a proper scoring rule, then $\mathcal{S}'(p_\theta, (x, y, r)) =$
 278 $-(\frac{rn}{n_l}L(\theta; x, y) + \lambda n(\frac{1-r}{n_u} - \frac{r}{n_l})H(\theta; x))$ is also a proper scoring rule.*

279 The proof is available in Appendix G, and follows directly from unbiasedness and the MCAR
 280 assumption. The main interpretation of this theorem is that we can expect DeSSL to be as well-
 281 calibrated as the complete case.

282 3.4 Consistency

283 We say that $\hat{\theta}$ is consistent if $d(\hat{\theta}, \theta^*) \xrightarrow{P} 0$ when $n \rightarrow \infty$, where d is a distance on Θ . The asymptotic
 284 properties of $\hat{\theta}$ depend on the behaviours of the functions L and H . We will thus require the following
 285 standard assumptions.

286 **Assumption 3.3.** The minimum θ^* of \mathcal{R} is well-separated: $\inf_{\theta: d(\theta^*, \theta) \geq \epsilon} \mathcal{R}(\theta) > \mathcal{R}(\theta^*)$.

287 **Assumption 3.4.** The uniform weak law of large number holds for both L and H .

288 **Theorem 3.5.** *Under the MCAR assumption, Assumption 3.3 and Assumption 3.4, $\hat{\theta} =$
 289 $\arg \min \hat{\mathcal{R}}_{DeSSL}$ is consistent.*

290 For proof of this theorem see Appendix G. This theorem is a simple application of van der Vaart's
 291 (2000) Theorem 5.7 proving the consistency of an M-estimator. Also, this result holds for the
 292 complete case, with $\lambda = 0$ which proves that the complete case is a solid baseline under the MCAR
 293 assumption. **Going further, we prove the asymptotic normality of $\hat{\theta}_{DeSSL}$ and showed that the
 294 asymptotic variance can be optimised with respect to λ .**

295 **Coupling of n_l and n_u under the MCAR assumption** Under the MCAR assumption, n_l and n_u
 296 are random variables. We have that $r \sim \mathcal{B}(\pi)$ (i.e. any x has the probability π of being labelled).
 297 Then, with n growing to infinity, we have $\frac{n_l}{n} = \frac{n_l}{n_l + n_u} \rightarrow \pi$. Therefore, both n_l and n_u grow to
 298 infinity and $\frac{n_l}{n_u} \rightarrow \frac{\pi-1}{\pi}$. This implies $n_u = \mathcal{O}(n_l)$ and then when n goes to infinity, both n_u and n_l
 299 go to infinity too and even if $n_u \gg n_l$.

300 3.5 Rademacher complexity and generalisation bounds

301 In this section, we prove an upper bound for the generalisation error of DeSSL. The unbiasedness of
 302 $\hat{\mathcal{R}}_{DeSSL}$ can directly be used to derive generalisation bounds based on the Rademacher complexity
 303 (Bartlett & Mendelson, 2002), defined in our case as

$$R_n = \mathbb{E}_{(\varepsilon_i)_{i \leq n}} \left[\sup_{\theta \in \Theta} \left(\frac{1}{n_l} \sum_{i=1}^{n_l} \varepsilon_i L(\theta; x_i, y_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} \varepsilon_i H(\theta; x_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} \varepsilon_i H(\theta; x_i) \right) \right], \quad (9)$$

304 where ε_i are i.i.d. Rademacher variables independent of the data. In the particular case of $\lambda = 0$,
 305 we recover the standard Rademacher complexity of the complete case. We can then now bound the
 306 generalisation error of a model trained using our new loss function.

307 **Theorem 3.6.** *We assume that labels are MCAR and that both L and H are bounded. Then, there
 308 exists a constant $\kappa > 0$, that depends on λ , L , H , and the ratio of observed labels, such that, with
 309 probability at least $1 - \delta$, for all $\theta \in \Theta$,*

$$\mathcal{R}(\theta) \leq \hat{\mathcal{R}}_{DeSSL}(\theta) + 2R_n + \kappa \sqrt{\frac{\log(4/\delta)}{n}}. \quad (10)$$

310 The proof follows Shalev-Shwartz & Ben-David (2014, Chapter 26), and is available in Appendix J.

311 4 Experiments

312 We evaluate the performance of DeSSL against different classic methods. The goal here is to compare
313 DeSSL methods and their original counterparts. In particular, we perform experiments with simple
314 SSL methods such as pseudo-label (PseudoLabel) and entropy minimisation (EntMIN) with varying
315 λ on MNIST (LeCun & Cortes, 2010) and CIFAR-10 and CIFAR-100 (Krizhevsky, 2009) and
316 compare them to the debiased method, respectively DeEntMin and DePseudoLabel. We also compare
317 PseudoLabel and DePseudoLabel on five small datasets of MedMNIST (Yang et al., 2021a;b) with a
318 fixed λ . The results of these experiments are reported below. In our figures, the error bars represent
319 the size of the 95% confidence interval (CI). Finally, we modified the implementation of Fixmatch
320 (Sohn et al., 2020) and compare it with its debiased version on CIFAR-10.

321 We also compare DeEntMin and DePseudoLabel to the biased version on a large range of tabular
322 datasets commonly used in SSL benchmarks (Chapelle et al., 2006; Guo et al., 2010). We do not
323 observe differences between the performance, see Appendix P. Finally, we show how simple it is to
324 debias an existing implementation, by demonstrating it on the consistency-based models benchmarked
325 by (Oliver et al., 2018), namely VAT, Π -model and MeanTeacher on CIFAR-10 and SVHN (Netzer
326 et al., 2011). We observe similar performances between the debiased and biased versions for the differ-
327 ent methods, both in terms of cross-entropy and accuracy. Moreover, these results have been obtained
328 using the hyperparameters finetuned for the biased versions. Therefore, it is likely that optimising the
329 hyperparameters for DeSSL will yield even better with the right hyperparameters, see Appendix O.

330 4.1 MNIST

331 MNIST is an advantageous dataset for SSL since classes are
332 well-separated. We compare PseudoLabel and DePseudoLabel
333 for a LeNet-like architecture using $n_l = 1000$ labelled data on
334 10 different splits of the training dataset into a labelled and unla-
335 belled set. Models are then evaluated using the standard 10,000
336 test samples. We used 10% of n_l as the validation set. We test
337 the influence of the hyperparameter λ and report the accuracy,
338 the cross-entropy and the expected calibration error (ECE, Guo
339 et al., 2017) at the epoch of best validation accuracy, see Fig-
340 ure 2 and Appendix L. In this example SSL and DeSSL have
341 almost the same accuracy for all λ , however, DeSSL seems to
342 be always better calibrated. To break the cluster assumption, we
343 reproduced the same experiment on a modified MNIST. Indeed,
344 we had label noise by replacing the true label for 20% of the
345 dataset with a randomly sampled label, see Appendix L. In this
346 setting, DeSSL performs better for large λ in terms of accuracy
347 and also provides a better calibration.

348 4.2 MedMNIST

349 We compare PseudoLabel and DePseudoLabel on different datasets of MedMNIST, a large-scale
350 MNIST-like collection of biomedical images. We selected the five smallest 2D datasets of the
351 collection, for these datasets it is likely that the cluster assumption no longer holds. We trained
352 a 5-layer CNN with a fixed $\lambda = 1$ and n_l at 10% of the training data. We report in Table 1 the
353 mean accuracy and cross-entropy on 5 different splits of the labelled and unlabelled data and the
354 number of labelled data used. We report the AUC in Appendix L. DePseudoLabel competes with
355 PseudoLabel in terms of accuracy and even success when PseudoLabel’s accuracy is less than the
356 complete case. Moreover, DePseudoLabel is always better in terms of cross-entropy, so calibration,
357 whereas PseudoLabel is always worse than the complete case.

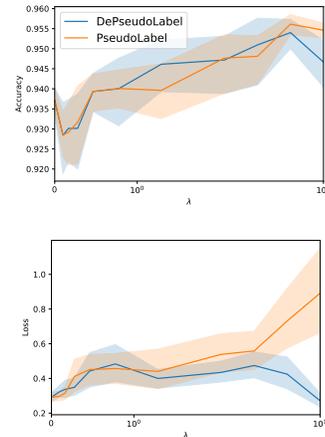


Figure 2: The influence of λ on Pseudo-label and DePseudo-label for a LeNet trained on MNIST with $n_l = 1000$: (Top) Mean test accuracy; (Bottom) Mean test cross-entropy, with 95% CI.

Table 1: Test accuracy and cross-entropy of Complete Case (CC), PseudoLabel (PL) and DePseudoLabel (DePL) on five datasets of MedMNIST.

DATASET	NL	CC		PL		DEPL	
		CROSS-ENTROPY	ACCURACY	CROSS-ENTROPY	ACCURACY	CROSS-ENTROPY	ACCURACY
DERMA	1000	1.95 ± 0.09	68.99 ± 1.20	2.51 ± 0.20	68.88 ± 1.03	1.88 ± 0.12	69.30 ± 0.85
PNEUMONIA	585	1.47 ± 0.04	83.94 ± 2.40	2.04 ± 0.04	85.83 ± 2.13	1.40 ± 0.06	84.36 ± 3.79
RETINA	160	1.68 ± 0.03	48.30 ± 3.06	1.80 ± 0.18	47.75 ± 2.50	1.67 ± 0.06	49.40 ± 2.62
BREAST	78	0.80 ± 0.04	76.15 ± 0.75	1.00 ± 0.26	74.74 ± 1.04	0.70 ± 0.03	76.67 ± 1.32
BLOOD	1700	6.11 ± 0.17	84.13 ± 0.83	6.61 ± 0.22	84.09 ± 1.17	6.53 ± 0.30	83.68 ± 0.59

358 4.3 CIFAR

359 We compare PseudoLabel and DePseudoLabel on CIFAR-10
360 and CIFAR-100. We trained a CNN-13 from Tarvainen &
361 Valpola (2017) on 5 different splits. For this experiment, we use
362 $n_l = 4000$ and use the rest of the dataset as unlabelled. Models
363 are then evaluated using the standard 10,000 test samples. For a
364 more realistic validation set, we used 10% of n_l as the validation
365 set. We test the influence of the hyperparameter λ and report the
366 accuracy and the cross-entropy at the epoch of best validation
367 accuracy, see Figure 3. We report the ECE in Appendix M. The
368 performance of both methods on CIFAR-100 with $n_l = 10000$
369 are reported in Appendix M. We observe DeSSL provides both
370 a better cross-entropy and ECE with the same accuracy for
371 small λ . For larger λ , DeSSL performs better in all the reported
372 metrics. We performed a paired Student’s t-test to ensure that
373 our results are significant and reported the p-values in Appendix
374 M. The p-values indicate that for λ close to 10, DeSSL is often
375 significantly better in all the metrics. Moreover, DeSSL for large
376 λ provides a better cross-entropy and ECE than the complete
377 case whereas SSL never does.

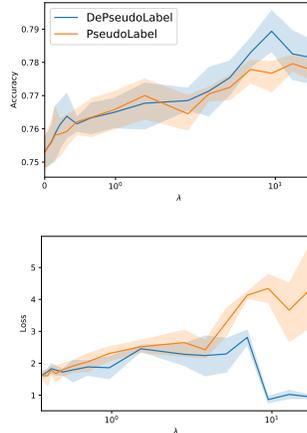


Figure 3: Influence of λ on Pseudo-label and DePseudo-label for a CNN trained on CIFAR with $n_l = 4000$: (Left) Mean test accuracy; (Right) Mean test cross-entropy, with 95% CI.

378 4.4 Fixmatch (Sohn et al., 2020)

379 We debiased a version of Fixmatch, see Appendix N for
380 further details. For this experiment, we use $n_l = 4000$ on 5
381 different folds. First, we report that a strong baseline using
382 data augmentation reach 87.27% accuracy. Then, we ob-
383 serve that on the debiasing method improve both accuracy
384 and cross-entropy of this modified version of Fixmatch.
385 Inspired by Zhu et al. (2022), we show that our method
386 improved performance on “poor” classes more equally than the biased version. Indeed, DeFixmatch
387 improves Fixmatch by 1.57% overall but by 4.91% on the worst class. We report in Appendix N
388 the accuracy per class of the different methods and the *benefit ratio* as defined by Zhu et al. (2022).

Table 2: 1st line: Accuracy, 2nd line: Worst class accuracy, 3rd line: Cross-entropy.

COMPLETE CASE	FIXMATCH	DEFIXMATCH
87.27 ± 0.25	93.87 ± 0.13	95.44 ± 0.10
70.08 ± 0.93	82.25 ± 2.27	87.16 ± 0.46
0.60 ± 0.01	0.27 ± 0.01	0.20 ± 0.01

389 5 Conclusion

390 Motivated by the remarks of van Engelen & Hoos (2020) and Oliver et al. (2018) on the missingness
391 of theoretical guarantees in SSL, we proposed a simple modification of SSL frameworks. We consider
392 frameworks based on the inclusion of unlabelled data in the computation of the risk estimator and
393 debias them using labelled data. We show theoretically that this debiasing comes with several theo-
394 retical guarantees. We demonstrate these theoretical results experimentally on several common SSL
395 datasets and some more challenging ones such as MNIST with label noise. DeSSL shows competitive
396 performance in terms of accuracy compared to its biased version but improves significantly the
397 calibration. There are several future directions open to us. We showed that λ_{opt} exists (Theorem 3.1)
398 and therefore our formula provides guidelines for the optimisation of λ . Finally, an interesting im-
399 provement would be to go beyond the MCAR assumption by considering settings with a distribution
400 mismatch between labelled and unlabelled data (Guo et al., 2020; Cao et al., 2021; Hu et al., 2022).

401 References

- 402 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean,
403 J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz,
404 R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah,
405 C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V.,
406 Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and
407 Zheng, X. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL
408 <https://www.tensorflow.org/>. Software available from tensorflow.org.
- 409 Ahfock, D. and McLachlan, G. J. On missing label patterns in semi-supervised learning. *arXiv*
410 *preprint arXiv:1904.02883*, 2019.
- 411 Aminian, G., Abroshan, M., Khalili, M. M., Toni, L., and Rodrigues, M. An information-theoretical
412 approach to semi-supervised learning under covariate-shift. In *International Conference on*
413 *Artificial Intelligence and Statistics*, pp. 7433–7449. PMLR, 2022.
- 414 Avramidis, A. N. and Wilson, J. R. A splitting scheme for control variates. *Operations Research*
415 *Letters*, 1993.
- 416 Barndorff-Nielsen, O. Plausibility inference. *Journal of the Royal Statistical Society: Series B*
417 *(Methodological)*, 38(2):103–123, 1976.
- 418 Bartlett, P. L. and Mendelson, S. Rademacher and Gaussian complexities: Risk bounds and structural
419 results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 420 Berthelot, D., Carlini, N., Goodfellow, I., Papernot, N., Oliver, A., and Raffel, C. A. Mixmatch:
421 A holistic approach to semi-supervised learning. *Advances in Neural Information Processing*
422 *Systems*, 2019.
- 423 Berthelot, D., Carlini, N., Cubuk, E. D., Kurakin, A., Sohn, K., Zhang, H., and Raffel, C. ReMix-
424 Match: Semi-supervised learning with distribution matching and augmentation anchoring. *Internat-*
425 *ional conference on Learning Representations*, 2020.
- 426 Cao, K., Brbic, M., and Leskovec, J. Open-world semi-supervised learning, 2021.
- 427 Chapelle, O., Schölkopf, B., and Zien, A. Semi-supervised learning. *MIT Press*, 2006.
- 428 Chen, S., Dobriban, E., and Lee, J. H. A group-theoretic framework for data augmentation. *Journal*
429 *of Machine Learning Research*, 21(245):1–71, 2020.
- 430 Corduneanu, A. and Jaakkola, T. On information regularization. In *UAI*. UAI, 2003.
- 431 Cubuk, E. D., Zoph, B., Shlens, J., and Le, Q. V. Randaugment: Practical automated data augmen-
432 tation with a reduced search space. In *Conference on Computer Vision and Pattern Recognition*
433 *Workshops*, 2020.
- 434 Defazio, A., Bach, F., and Lacoste-Julien, S. SAGA: A fast incremental gradient method with support
435 for non-strongly convex composite objectives. *Advances in Neural Information Processing Systems*,
436 2014.
- 437 Fox-Roberts, P. and Rosten, E. Unbiased generative semi-supervised learning. *The Journal of*
438 *Machine Learning Research*, 15(1):367–443, 2014.
- 439 Gan, H., Li, Z., Fan, Y., and Luo, Z. Dual learning-based safe semi-supervised learning. *IEEE Access*,
440 6:2615–2621, 2017.
- 441 Gneiting, T. and Raftery, A. E. Strictly proper scoring rules, prediction, and estimation. *Journal of*
442 *the American statistical Association*, 102(477):359–378, 2007.
- 443 Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and
444 Bengio, Y. Generative adversarial nets. *Advances in Neural Information Processing Systems*, 2014.
- 445 Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. MIT Press, 2016.

- 446 Grandvalet, Y. and Bengio, Y. Semi-supervised learning by entropy minimization. *Advances in*
447 *Neural Information Processing Systems*, 2004.
- 448 Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks.
449 *International Conference on Machine Learning*, 2017.
- 450 Guo, L.-Z., Zhang, Z.-Y., Jiang, Y., Li, Y.-F., and Zhou, Z.-H. Safe deep semi-supervised learning for
451 unseen-class unlabeled data. *International Conference on Machine Learning*, 2020.
- 452 Guo, Y., Niu, X., and Zhang, H. An extensive empirical study on semi-supervised learning. *IEEE*
453 *International Conference on Data Mining*, 2010.
- 454 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser,
455 E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M.,
456 Haldane, A., Fernández del Rfo, J., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K.,
457 Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. Array programming with
458 NumPy. *Nature*, 585:357–362, 2020. doi: 10.1038/s41586-020-2649-2.
- 459 Hu, X., Niu, Y., Miao, C., Hua, X.-S., and Zhang, H. On non-random missing labels in semi-
460 supervised learning. In *International Conference on Learning Representations*, 2022.
- 461 Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance
462 reduction. *Advances in Neural Information Processing Systems*, 2013.
- 463 Kawakita, M. and Takeuchi, J. Safe semi-supervised learning based on weighted likelihood. *Neural*
464 *Networks*, 53:146–164, 2014.
- 465 Kingma, D. P., Mohamed, S., Rezende, D. J., and Welling, M. Semi-supervised learning with deep
466 generative models. In *Advances in Neural Information Processing Systems*, pp. 3581–3589, 2014.
- 467 Krause, A., Perona, P., and Gomes, R. Discriminative clustering by regularized information maxi-
468 mization. *Advances in neural information processing systems*, 23, 2010.
- 469 Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, MIT, NYU,
470 2009.
- 471 Laine, S. and Aila, T. Temporal ensembling for semi-supervised learning. *International Conference*
472 *on Learning Representations*, 2017.
- 473 LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- 475 Lee, D.-H. Pseudo-Label : The simple and efficient semi-supervised learning method for deep
476 neural networks. *Workshop on challenges in representation learning, International conference on*
477 *machine learning*, 2013.
- 478 Li, Y.-F. and Zhou, Z.-H. Towards making unlabeled data never hurt. *IEEE transactions on pattern*
479 *analysis and machine intelligence*, 37:175–188, 2014.
- 480 Li, Y.-F., Kwok, J. T., and Zhou, Z.-H. Towards safe semi-supervised learning for multivariate
481 performance measures. *AAAI Conference on Artificial Intelligence*, 2016.
- 482 Little, R. J. and Rubin, D. B. *Statistical Analysis with Missing Data*. John Wiley & Sons, 2019.
- 483 Liu, T. and Goldberg, Y. Kernel machines with missing responses. *Electronic Journal of Statistics*,
484 14:3766–3820, 2020.
- 485 Lundh, F., Ellis, M., et al. Python imaging library (pil), 2012.
- 486 McKinney, W. et al. Data structures for statistical computing in python. In *Proceedings of the 9th*
487 *Python in Science Conference*, volume 445, pp. 51–56. Austin, TX, 2010.
- 488 McLachlan, G. J. Estimating the linear discriminant function from initial samples containing a small
489 number of unclassified observations. *Journal of the American statistical association*, 72:403–406,
490 1977.

- 491 Mey, A. and Loog, M. Improvability through semi-supervised learning: A survey of theoretical
492 results. *arXiv preprint arXiv:1908.09574*, 2019.
- 493 Miyato, T., Maeda, S.-i., Koyama, M., and Ishii, S. Virtual adversarial training: A regularization
494 method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and*
495 *machine intelligence*, 41:1979–1993, 2018.
- 496 Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images
497 with unsupervised feature learning. 2011.
- 498 Newey, W. K. and McFadden, D. Large sample estimation and hypothesis testing. *Handbook of*
499 *econometrics*, 4:2111–2245, 1994.
- 500 Oliver, A., Odena, A., Raffel, C., Cubuk, E. D., and Goodfellow, I. J. Realistic evaluation of deep
501 semi-supervised learning algorithms. *Advances in Neural Information Processing Systems*, 2018.
- 502 Owen, A. B. *Monte Carlo theory, methods and examples*. 2013.
- 503 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
504 Gimselshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani,
505 A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. PyTorch: An imperative
506 style, high-performance deep learning library. In Wallach, H., Larochelle, H., Beygelzimer, A.,
507 d'Alché-Buc, F., Fox, E., and Garnett, R. (eds.), *Advances in Neural Information Processing*
508 *Systems 32*. Curran Associates, Inc., 2019. URL [http://papers.neurips.cc/paper/](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
509 [9015-pytorch-an-imperative-style-high-performance-deep-learning-library.](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
510 [pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 511 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
512 Prettenhofer, P., Weiss, R., Dubourg, V., et al. Scikit-learn: Machine learning in python. *Journal of*
513 *machine learning research*, 12(Oct):2825–2830, 2011.
- 514 Pereyra, G., Tucker, G., Chorowski, J., Kaiser, Ł., and Hinton, G. Regularizing neural networks by
515 penalizing confident output distributions. *Workshop track, International Conference on Learning*
516 *Representations*, 2017.
- 517 Pham, H., Dai, Z., Xie, Q., and Le, Q. V. Meta pseudo labels. *Conference on Computer Vision and*
518 *Pattern Recognition*, 2021.
- 519 Rizve, M. N., Duarte, K., Rawat, Y. S., and Shah, M. In defense of pseudo-labeling: An uncertainty-
520 aware pseudo-label selection framework for semi-supervised learning. *International Conference*
521 *on Learning Representations*, 2021.
- 522 Sajjadi, M., Javanmardi, M., and Tasdizen, T. Regularization with stochastic transformations and
523 perturbations for deep semi-supervised learning. *Advances in Neural Information Processing*
524 *Systems*, 2016.
- 525 Sakai, T., Plessis, M. C., Niu, G., and Sugiyama, M. Semi-supervised classification based on
526 classification from positive and unlabeled data. *International conference on machine learning*,
527 2017.
- 528 Schölkopf, B., Janzing, D., Peters, J., Sgouritsa, E., Zhang, K., and Mooij, J. On causal and anticausal
529 learning. *International conference on machine learning*, 2012.
- 530 Scudder, H. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions*
531 *on Information Theory*, 11:363–371, 1965.
- 532 Seeger, M. Learning with labeled and unlabeled data. *Technical report*, 2000.
- 533 Serre, D. *Matrices. Theory and Applications (Second edition)*. Springer, 2010.
- 534 Shalev-Shwartz, S. and Ben-David, S. *Understanding Machine Learning: From Theory to Algorithms*.
535 Cambridge university press, 2014.
- 536 Singh, A., Nowak, R., and Zhu, J. Unlabeled data: Now it helps, now it doesn't. *Advances in Neural*
537 *Information Processing Systems*, 2008.

- 538 Sohn, K., Berthelot, D., Li, C.-L., Zhang, Z., Carlini, N., Cubuk, E. D., Kurakin, A., Zhang, H.,
539 and Raffel, C. FixMatch: Simplifying semi-supervised learning with consistency and confidence.
540 *Avances in Neural Information Processing Systems*, 2020.
- 541 Sokolovska, N., Cappé, O., and Yvon, F. The asymptotics of semi-supervised learning in discrimina-
542 tive probabilistic models. In *International Conference on Machine Learning*, 2008.
- 543 Tarvainen, A. and Valpola, H. Mean teachers are better role models: Weight-averaged consistency
544 targets improve semi-supervised deep learning results. *Advancer in Neural Information Processing*
545 *Systems*, 2017.
- 546 Trapp, M., Madl, T., Peharz, R., Pernkopf, F., and Trappl, R. Safe semi-supervised learning of
547 sum-product networks. *Conference on Uncertainty in Artificial Intelligence*, 2017.
- 548 Tsiatis, A. A. *Semiparametric theory and missing data*. Springer, 2006.
- 549 Tsuchiya, T., Charoenphakdee, N., Sato, I., and Sugiyama, M. Semisupervised ordinal regression
550 based on empirical risk minimization. *Neural Computation*, 33:3361–3412, 2021.
- 551 van der Vaart, A. W. *Asymptotic statistics*. Cambridge university press, 2000.
- 552 van Engelen, J. E. and Hoos, H. H. A survey on semi-supervised learning. *Machine Learning*, 109:
553 373–440, 2020.
- 554 Van Rossum, G. and Drake Jr, F. L. *Python reference manual*. Centrum voor Wiskunde en Informatica
555 Amsterdam, 1995.
- 556 Verma, V., Kawaguchi, K., Lamb, A., Kannala, J., Bengio, Y., and Lopez-Paz, D. Interpolation
557 consistency training for semi-supervised learning. *International Joint Conference on Artificial*
558 *Intelligence*, 2019.
- 559 Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Lukauskas, S., Gemperline, D. C., Augspurger,
560 T., Halchenko, Y., Cole, J. B., Warmenhoven, J., de Ruiter, J., Pye, C., Hoyer, S., Vanderplas,
561 J., Villalba, S., Kunter, G., Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram,
562 Y., Yarkoni, T., Williams, M. L., Evans, C., Fitzgerald, C., Brian, Fonnesbeck, C., Lee, A.,
563 and Qalieh, A. mwaskom/seaborn: v0.8.1 (september 2017), September 2017. URL <https://doi.org/10.5281/zenodo.883859>.
- 565 Wei, C., Shen, K., Chen, Y., and Ma, T. Theoretical analysis of self-training with deep networks on
566 unlabeled data. In *International Conference on Learning Representations*, 2021.
- 567 Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. Unsupervised data augmentation for
568 consistency training. *Advances in Neural Information Processing Systems*, 2019.
- 569 Yang, J., Shi, R., and Ni, B. MedMNIST classification decathlon: A lightweight AutoML benchmark
570 for medical image analysis. In *IEEE 18th International Symposium on Biomedical Imaging (ISBI)*,
571 pp. 191–195, 2021a.
- 572 Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., and Ni, B. MedMNIST v2: A
573 large-scale lightweight benchmark for 2D and 3D biomedical image classification. *arXiv preprint*
574 *arXiv:2110.14795*, 2021b.
- 575 Zhang, B., Wang, Y., Hou, W., Wu, H., Wang, J., Okumura, M., and Shinozaki, T. FlexMatch: Boost-
576 ing semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information*
577 *Processing Systems*, 2021a.
- 578 Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization.
579 *Internation Conference on Learning Representations*, 2017.
- 580 Zhang, S., Wang, M., Liu, S., Chen, P.-Y., and Xiong, J. How unlabeled data improve generalization
581 in self-training? A one-hidden-layer theoretical analysis. In *International Conference on Learning*
582 *Representations*, 2021b.
- 583 Zhu, X., Ghahramani, Z., and Lafferty, J. D. Semi-supervised learning using Gaussian fields and
584 harmonic functions. *International conference on machine learning*, 2003.
- 585 Zhu, Z., Luo, T., and Liu, Y. The rich get richer: Disparate impact of semi-supervised learning. In
586 *International Conference on Learning Representations*, 2022.

587 **Checklist**

- 588 1. For all authors...
- 589 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
590 contributions and scope? [Yes]
- 591 (b) Did you describe the limitations of your work? [Yes] See section 5
- 592 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Theoretical
593 work
- 594 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
595 them? [Yes]
- 596 2. If you are including theoretical results...
- 597 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 598 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices
- 599 3. If you ran experiments...
- 600 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
601 perimental results (either in the supplemental material or as a URL)? [Yes] code and
602 instructions to run Fixmatch.
- 603 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
604 were chosen)? [Yes] See Appendices
- 605 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
606 ments multiple times)? [Yes] Confidence intervals on all figures.
- 607 (d) Did you include the total amount of compute and the type of resources used (e.g., type
608 of GPUs, internal cluster, or cloud provider)? [Yes] An estimation in Appendix
- 609 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 610 (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix
- 611 (b) Did you mention the license of the assets? [Yes]
- 612 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
613 Code
- 614 (d) Did you discuss whether and how consent was obtained from people whose data you're
615 using/curating? [N/A]
- 616 (e) Did you discuss whether the data you are using/curating contains personally identifiable
617 information or offensive content? [N/A]
- 618 5. If you used crowdsourcing or conducted research with human subjects...
- 619 (a) Did you include the full text of instructions given to participants and screenshots, if
620 applicable? [N/A]
- 621 (b) Did you describe any potential participant risks, with links to Institutional Review
622 Board (IRB) approvals, if applicable? [N/A]
- 623 (c) Did you include the estimated hourly wage paid to participants and the total amount
624 spent on participant compensation? [N/A]

625 **A Toy example**

626 We trained a 4-layer neural network (1/20/100/20/1) with ReLU activation function using 25,000
 627 labelled and 25,000 unlabelled points drawn from two 1D uniform laws with an overlap. We used
 628 $\lambda = 1$ and a confidence threshold for Pseudo-label $\tau = 0.70$. We optimised the model’s weights
 629 using a stochastic gradient descent (SGD) optimiser with a learning rate of 0.1.

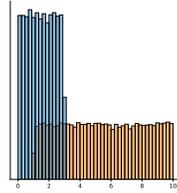


Figure 4: Data histogram

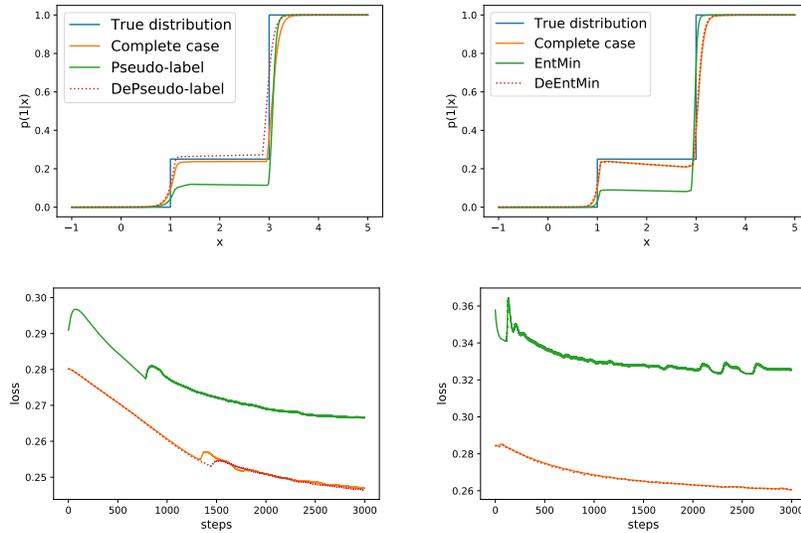


Figure 5: 4-layer neural net trained using SSL methods on a 1D dataset drawn from two uniform laws. (Top-left) Posterior probabilities $p(1|x)$ of the same model trained following either complete case (only labelled data), Pseudo-label or our DePseudo-label. (Top-right) Same for EntMin and DeEntMin (Bottom-left) Training cross-entropy for Pseudo-label and DePseudo-label (Bottom-right) Training cross-entropy for EntMin and DeEntMin

630 **B Details on surrogates and more examples**

631 We provide in this appendix further details on our classification of SSL methods between entropy-
 632 based and consistency-based (see Section 2.2.3). We detail a general framework for both of these
 633 methods’ classes. We also show how popular SSL methods are related to our framework.

634 **B.1 Entropy-based**

635 We class as entropy-based, methods that aim to minimise a term of entropy such as Grandvalet &
 636 Bengio (2004) which minimises Shannon’s entropy or pseudo-label which is a form of entropy, see
 637 Remark E.5. These methods encourage the model to be confident on unlabelled data, implicitly using
 638 the cluster assumption. We recall those entropy-based methods can all be described as an expectation
 639 of L under a distribution π_x computed at the datapoint x :

$$H(\theta; x) = \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})}[L(\theta; \tilde{x}, \tilde{y})]. \quad (11)$$

640 **Pseudo-label:** As presented in the core article, the unsupervised objective of pseudo-label can be
 641 written as an expectation of L on the distribution $\pi_x(\tilde{x}, \tilde{y}) = \delta_x(\tilde{x})p_\theta(\tilde{y}|\tilde{x})$. Recently, Lee (2013)
 642 encouraged the pseudo-labels method for deep semi-supervised learning. Then, Rizve et al. (2021)
 643 recently improved the pseudo-label selection by introducing an uncertainty-aware mechanism on
 644 the confidence of the model concerning the predicted probabilities. Pham et al. (2021) reaches
 645 state-of-the-art on the Imagenet challenge using pseudo-labels on a large dataset of additional images.

646 **B.2 Pseudo-label and data augmentation**

647 Recently, several methods based on data augmentation have been proposed and proven to perform well
 648 on a large spectrum of SSL tasks. The idea is to have a model resilient to strong data-augmentation of
 649 the input (Berthelot et al., 2019; 2020; Sohn et al., 2020; Xie et al., 2019; Zhang et al., 2021a). These
 650 methods rely both on the cluster assumption and the smoothness assumption and are at the border
 651 between entropy-based and consistency-based methods. The idea is to have the same prediction
 652 for an input and an augmented version of it. For instance, in Sohn et al. (2020), we first compute
 653 pseudo-labels predicted using a weakly-augmented version of x (flip-and-shift data augmentation)
 654 and then minimise the likelihood with the predictions of the model on a strongly augmented version
 655 of x . In Xie et al. (2019), the method is a little bit different as we minimise the cross entropy between
 656 the prediction of the model on x and the predictions of an augmented version. In both cases, the
 657 unsupervised part of the risk estimator can be reformulated as Equation 11.

658 **Fixmatch:** In Fixmatch, Sohn et al. (2020), the unsupervised objective can be written as:

$$H(\theta; x) = \mathbb{1}[\max_y p_{\hat{\theta}}(y|x_1) > \tau]L(\theta; x_2, \arg \max_y p_{\hat{\theta}}(y|x_1)) \quad (12)$$

659 where $\hat{\theta}$ is a fixed copy of the current parameters θ indicating that the gradient is not prop-
 660 agated through it, x_1 is a weakly-augmented version of x and x_2 a strongly-augmented
 661 one. Therefore, we write H as an expectation of L on the distribution $\pi_x(\tilde{x}, \tilde{y}) =$
 662 $\delta_{x_2}(\tilde{x})\delta_{\arg \max_y p_{\hat{\theta}}(y|x_1)}(\tilde{y})\mathbb{1}[\max_y p_{\hat{\theta}}(y|x_1) > \tau]$.

663 **UDA:** In UDA, Xie et al. (2019), the unsupervised objective can be written as:

$$H(\theta; x) = \sum_y p_{\hat{\theta}}(y|x)L(\theta; x_1, y) \quad (13)$$

664 where $\hat{\theta}$ is a fixed copy of the current parameters θ indicating that the gradient is not propagated
 665 through it and x_1 is an augmented version of x . Therefore, we write H as an expectation of L on the
 666 distribution $\pi_x(\tilde{x}, \tilde{y}) = \delta_{x_1}(\tilde{x})p_{\hat{\theta}}(\tilde{y}|\tilde{x})$.

667 **Others:** Recently, have been proposed in the literature Zhang et al. (2021a) and Rizve et al. (2021).
668 The former is an improved version of Fixmatch with a variable threshold τ with respect to the class
669 and the training stage. The latter introduces a measurement of uncertainty in the pseudo-labelling
670 step to improve the selection. They also introduce negative pseudo-labels to improve the single-label
671 classification.

672 B.3 Consistency-based

673 Consistency-based methods aim to smooth the decision function of the models or have more stable
674 predictions. These objectives H are not directly a form of expectation of L but are equivalent to an
675 expectation of L . For all the following methods we can write the unsupervised objective H such that:
676

$$C_1 \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})} [L(\theta; \tilde{x}, \tilde{y})] \leq H(\theta; x) \leq C_2 \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})} [L(\theta; \tilde{x}, \tilde{y})], \quad (14)$$

677 with $0 \leq C_1 \leq C_2$.

678 Indeed, consistency-based methods minimise an unsupervised objective that is a divergence between
679 the model predictions and a modified version of the input (data augmentation) or a perturbation of the
680 model. Using the fact that all norms are equivalent in a finite-dimensional space such as the space of
681 the labels, we have the equivalence between a consistency-based H and an expectation of L .

682 **VAT** The virtual adversarial training method proposed by (Miyato et al., 2018) generates the most
683 impactful perturbation r_{adv} to add to x . The objective is to train a model robust to input perturbations.
684 This method is closely related to adversarial training introduced by Goodfellow et al. (2014).

$$H(\theta; x) = \mathbf{Div}(f_{\hat{\theta}}(x, \cdot), f_{\theta}(x + r_{adv}, \cdot))$$

685 where the **Div** is a non-negative function that measures the divergence between two distributions, the
686 cross-entropy or the KL divergence for instance. If the divergence function is the cross-entropy, it is
687 straightforward to write the unlabelled objective as Equation 3. If the objective function is the KL
688 divergence, we can write the objective as

$$H(\theta; x) = \mathbb{E}_{\pi_x(\tilde{x}+r, \tilde{y})} [L(\theta; \tilde{x}, \tilde{y})] - \mathbb{E}_{\pi_x(\tilde{x}, \tilde{y})} [L(\hat{\theta}; \tilde{x}, \tilde{y})] \quad (15)$$

689 with $\pi_x(\tilde{x}, \tilde{y}) = \delta_x(\tilde{x})p_{\hat{\theta}}(y|x)$. Therefore, variation of H with respect to θ are the same as
690 $\mathbb{E}_{\pi_x(\tilde{x}+r, \tilde{y})} [L(\theta; \tilde{x}, \tilde{y})]$. VAT is also a method between consistency-based and entropy-based methods
691 as long as we use the KL-divergence or the cross-entropy as the measure of divergence.

692 **Mean-Teacher** A different form of pseudo-labelling is the Mean-Teacher approach proposed by
693 (Tarvainen & Valpola, 2017) where pseudo-labels are generated by a teacher model for a student
694 model. The parameters of the student model are updated, while the teacher’s are a moving average
695 of the student’s parameters from the previous training steps. The idea is to have a more stable
696 pseudo-labelling using the teacher than in the classic Pseudo-label. Final predictions are made by the
697 student model. A generic form of the unsupervised part of the risk estimator is then

$$H(\theta; x) = \sum_y (p_{\theta}(y|x) - p_{\hat{\theta}}(y|x))^2,$$

698 where $\hat{\theta}$ are the fixed parameters of the teacher.

699 **II-Model** The II-Models are intrinsically stochastic models (for example a model with dropout)
700 encouraged to make consistent predictions through several passes of the same x in the model. The
701 SSL loss is using the stochastic behaviour of the model where the model f_{θ} and penalises different
702 predictions for the same x (Sajjadi et al., 2016). Let’s note $f_{\theta}(x, \cdot)_1$ and $f_{\theta}(x, \cdot)_2$ two passes of x
703 through the model f_{θ} . A generic form of the unsupervised part of the risk estimator is then

$$H(\theta; x) = \mathbf{Div}(f_{\theta}(x, \cdot)_1, f_{\theta}(x, \cdot)_2), \quad (16)$$

704 where **Div** is a measure of divergence between two distributions (often the Kullback-Leibler diver-
705 gence).

706 **Temporal ensembling** Temporal ensembling (Laine & Aila, 2017) is a form of Π -Model where
707 we compare the current prediction of the model on the input x with an accumulation of the previous
708 passes through the model. Then, the training is faster as the network is evaluated only once per input
709 on each epoch and the perturbation is expected to be less noisy than for Π -models.

710 **ICT** Interpolation consistency training (Verma et al., 2019) is an SSL method based on the mixup
711 operation (Zhang et al., 2017). The model trained is then consistent to predictions at interpolations.
712 The unsupervised term of the objective is then computed on two terms:

$$H(\theta; x_1, x_2) = \mathbf{Div} (f_\theta(\alpha x_1 + (1 - \alpha)x_2, \cdot), \alpha f_{\hat{\theta}}(x_1, \cdot) + (1 - \alpha)f_{\hat{\theta}}(x_2, \cdot)), \quad (17)$$

713 with α drawn with from a distribution $\mathcal{B}(a, a)$. With the exact same transformation, we will be able
714 to show that this objective is equivalent to a form of expectation of L .

715 C On the semi-supervised bias

716 We provide in this appendix a further explanation of the risk induced by the SSL bias as introduced
717 in Section 2.3.

718 Presented methods minimise a biased version of the risk under the MCAR assumption and therefore
719 classical learning theory does not apply anymore,

$$\mathbb{E}[\hat{\mathcal{R}}_{SSL}(\theta)] = \mathbb{E}[L(\theta; x, y)] + \lambda \mathbb{E}[H(\theta; x, y)] \neq \mathcal{R}(\theta). \quad (18)$$

720 Learning over a biased estimate of the risk is not necessarily unsafe but it is difficult to provide
721 theoretical guarantees on such methods even if some works try to do so with strong assumptions on the
722 data distribution (Mey & Loog 2019, Section 4 and 5, Zhang et al. 2021b). Previous works proposed
723 generalisation error bounds of SSL methods under strong assumptions on the data distribution or the
724 true model. We refer to the survey by Mey & Loog (2019). More recently, Wei et al. (2021) proves an
725 upper bound for training deep models with the pseudo-label method under strong assumption. Under
726 soft assumptions, Aminian et al. (2022) provides an error bound showing that the choice of H is
727 crucial to provide good performances.

728 Indeed, the unbiased nature of the risk estimate is crucial in the development of learning theory. This
729 bias on the risk estimate may look like the one of a regularisation, such as the ridge regularisation.
730 However, SSL and regularisation are intrinsically different for several reasons:

- 731 • Regularisers have a vanishing impact in the limit of infinite data whereas SSL usually do
732 not in the proposed methods, see Equation 18. A solution would be to choose λ with respect
733 of the number of data points and make it vanish when n goes to infinity. However, in most
734 works, the choice of λ is independent of the number of n or n_l (Oliver et al., 2018; Sohn
735 et al., 2020).
- 736 • One of the main advantages of regularisation is to turn the learning problem into a “more
737 convex” problem, see Shalev-Shwartz & Ben-David (2014, Chapter 13). Indeed, ridge
738 regularisation will often turn a convex problem into a strongly-convex problem. However,
739 SSL faces the danger to turn the learning problem as non-convex as previously noted by
740 Sokolovska et al. (2008).
- 741 • The objective of a regulariser is to bias the risk towards optimum with smooth decision
742 functions whereas entropy-based SSL will lead to sharp decision functions.
- 743 • Regularisation usually does not depend on the data whereas H does in the SSL framework.

744 A entropy bias has been actually used by Pereyra et al. (2017) as a regulariser but as entropy
745 *maximisation* which should has an effect that is the opposite of the SSL method introduced by
746 Grandvalet & Bengio (2004), the entropy minimisation.

747 **D Proof that $\hat{\mathcal{R}}_{DeSSL}(\theta)$ is unbiased under MCAR**

748 **Theorem D.1.** *Under the MCAR assumption, $\hat{\mathcal{R}}_{DeSSL}(\theta)$ is an unbiased estimator of $\mathcal{R}(\theta)$.*

749 As a consequence of the theorem, under the MCAR assumption, $\hat{\mathcal{R}}_{CC}(\theta)$ is also unbiased as a special
750 case of $\hat{\mathcal{R}}_{DeSSL}(\theta)$ for $\lambda = 0$

751 **Proof:** We first recall that the DeSSL risk estimator $\hat{\mathcal{R}}_{DeSSL}(\theta)$ is defined for any λ by

$$\begin{aligned}\hat{\mathcal{R}}_{DeSSL}(\theta) &= \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i) \\ &= \sum_{i=1}^n \left(\frac{r_i}{n_l} L(\theta; x_i, y_i) + \lambda \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right) H(\theta; x_i) \right).\end{aligned}\tag{19}$$

752 By the law of total expectation:

$$\mathbb{E}[\hat{\mathcal{R}}_{DeSSL}(\theta)] = \mathbb{E}_r \left[\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)|r] \right].$$

753 As far as we are under the MCAR assumption, the data (x, y) and the missingness variable r are
754 independent thus, $\mathbb{E}_r \left[\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)|r] \right] = \mathbb{E}_r \left[\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)] \right]$.

755 We focus on $\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)]$. First, we replace $\hat{\mathcal{R}}_{DeSSL}(\theta)$ by its definition and then use the
756 linearity of the expectation. Then,

$$\begin{aligned}\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)] &= \mathbb{E} \left[\frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i) \right] \quad \text{by definition} \\ &= \frac{1}{n_l} \sum_{i=1}^{n_l} \mathbb{E}[L(\theta; x_i, y_i)] + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} \mathbb{E}[H(\theta; x_i)] - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} \mathbb{E}[H(\theta; x_i)] \quad \text{by linearity}\end{aligned}$$

757 The couples (x_i, y_i) are i.i.d. samples following the same distribution. Then, we have

$$\begin{aligned}\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)] &= \frac{1}{n_l} \sum_{i=1}^{n_l} \mathbb{E}[L(\theta; x, y)] + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} \mathbb{E}[H(\theta; x)] - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} \mathbb{E}[H(\theta; x)] \quad \text{i.i.d samples} \\ &= \mathbb{E}[L(\theta; x, y)] \\ &= \mathcal{R}(\theta).\end{aligned}$$

758 Finally, we have the results that , $\hat{\mathcal{R}}_{DeSSL}(\theta)$ is unbiased as $\mathcal{R}(\theta)$ is a constant,

$$\mathbb{E}[\hat{\mathcal{R}}_{DeSSL}(\theta)] = \mathbb{E} \left[\mathbb{E}_{x,y}[\hat{\mathcal{R}}_{DeSSL}(\theta)|r] \right] = \mathbb{E}_r [\mathcal{R}(\theta)] = \mathcal{R}(\theta).\tag{20}$$

759 **E Proof and comments about Theorem 3.1**

760 **Theorem 3.1** *The function $\lambda \mapsto \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r)$ reaches its minimum for:*

$$\lambda_{opt} = \frac{n_u \text{Cov}(L(\theta; x, y), H(\theta; x))}{n \mathbb{V}(H(\theta; x))} \quad (21)$$

761 *and*

$$\begin{aligned} \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r)|_{\lambda_{opt}} &= \left(1 - \frac{n_u}{n} \rho_{L,H}^2\right) \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)) \\ &\leq \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)), \end{aligned} \quad (22)$$

762 *where $\rho_{L,H} = \text{Corr}(L(\theta; x, y), H(\theta; x))$.*

763 **Proof:** For any $\lambda \in \mathbb{R}$, we want to compute the variance:

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r).$$

764 Under the MCAR assumption, x and y are both jointly independent of r . Also, the couples (x_i, y_i, r_i)
765 are independent. Therefore, we have

$$\begin{aligned} \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) &= \sum_{i=1}^n \mathbb{V}_{(x_i, y_i) \sim p(x, y|r)} \left(\frac{r_i}{n_l} L(\theta, x_i, y_i) + \lambda \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right) H(\theta, x_i) \right) \quad \text{i.i.d samples} \\ &= \sum_{i=1}^n \mathbb{V}_{(x_i, y_i) \sim p(x, y)} \left(\frac{r_i}{n_l} L(\theta, x_i, y_i) + \lambda \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right) H(\theta, x_i) \right) \quad (x, y) \text{ and } r \text{ independent} \end{aligned}$$

766 Using the fact that the couples (x_i, y_i) are i.i.d. samples following the same distribution, we have

$$\begin{aligned} \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) &= \sum_{i=1}^n \mathbb{V}_{(x, y) \sim p(x, y)} \left(\frac{r_i}{n_l} L(\theta, x, y) + \lambda \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right) H(\theta, x) \right) \\ &= \sum_{i=1}^n \frac{r_i^2}{n_l^2} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right)^2 \mathbb{V}(H(\theta, x)) \quad \text{using covariance} \\ &\quad + 2\lambda \frac{r_i}{n_l} \left(\frac{1-r_i}{n_u} - \frac{r_i}{n_l} \right) \text{Cov}(L(\theta, x, y), H(\theta, x)) \end{aligned}$$

767 Now, we remark that the variable r is binary and therefore $r^2 = r$, $(1-r)^2 = 1-r$ and $r(1-r) = 0$.
768 Using that and simplifying, we have

$$\begin{aligned} \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) &= \sum_{i=1}^n \frac{r_i}{n_l^2} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \frac{(1-r_i)n_l^2 + r_i n_u^2}{n_l^2 n_u^2} \mathbb{V}(H(\theta, x)) \\ &\quad - 2\lambda \frac{r_i}{n_l^2} \text{Cov}(L(\theta, x, y), H(\theta, x)) \end{aligned}$$

769 Finally, by summing and simplifying the expression (note that $n_l + n_u = n$), we compute the
770 expression variance,

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) = \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \frac{n}{n_l n_u} \mathbb{V}(H(\theta, x)) - \frac{2\lambda}{n_l} \text{Cov}(L(\theta, x, y), H(\theta, x))$$

771 So $\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r)$ is a quadratic function in λ and reaches its minimum for λ_{opt} such that:

$$\lambda_{opt} = \frac{n_u \text{Cov}(L(\theta, x, y), H(\theta, x))}{n \mathbb{V}(H(\theta, x))}$$

772 And, at λ_{opt} , the variance of $\hat{\mathcal{R}}_{DeSSL}(\theta)|r$ becomes

$$\begin{aligned}\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) &= \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) \left(1 - \frac{n_u}{n} \frac{\text{Cov}(L(\theta, x, y), H(\theta, x))^2}{\mathbb{V}(H(\theta, x))\mathbb{V}(L(\theta, x, y))} \right) \\ &= \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) \left(1 - \frac{n_u}{n} \text{Corr}(L(\theta, x, y), H(\theta, x))^2 \right) \\ &= \left(1 - \frac{n_u}{n} \rho_{L,H}^2 \right) \frac{1}{n_l} \mathbb{V}(L(\theta, x, y))\end{aligned}$$

773 *Remark E.1.* If H is perfectly correlated with L ($\rho_{L,H} = 1$), then the variance of the DeSSL estimator
774 is equal to the variance of the estimator with no missing labels.

775 *Remark E.2. Is it possible to estimate λ_{opt} in practice ?* The data distribution $p(x, y)$ being
776 unknown, the computation of λ_{opt} is not possible directly. Therefore, we need to use an estimator of
777 the covariance $\text{Cov}(L(\theta; x, y), H(\theta; x))$ and the variance $\mathbb{V}(H(\theta; x))$ (See Equation 23). Also, we
778 have to be careful not to introduce a new bias with the computation of λ_{opt} , indeed, if we compute
779 it using the training set, λ_{opt} becomes dependent of x and y and therefore $\hat{\mathcal{R}}_{DeSSL}(\theta)|r$ becomes
780 biased. A solution would be to use a validation dataset for its computation. Another approach is to
781 compute it using the splitting method (Avramidis & Wilson, 1993). Moreover, the computation of
782 λ_{opt} is tiresome and time-consuming in practice as it has to be updated for every different value of θ ,
783 so at each gradient step.

$$\hat{\lambda}_{opt} = \frac{\frac{1}{n_l} \sum_{i=1}^{n_l} (L(\theta; x_i, y_i) - \bar{L}(\theta))(H(\theta; x_i) - \bar{H}(\theta))}{\frac{1}{n} \sum_{i=1}^n (H(\theta; x_i) - \bar{H}(\theta))^2} \quad (23)$$

784 where $\bar{H}(\theta) = \frac{1}{n} \sum_{i=1}^n H(\theta; x_i)$ and $\bar{L}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i)$

785 *Remark E.3. About the sign of λ* As explained in the article, the theorem still has a *quantitative*
786 merit when it comes to choosing λ , by telling that the sign of λ is positive when H and L are
787 positively correlated which will generally be the case with the examples mentioned in the article. For
788 instance, concerning the entropy minimisation technique, the following proposition proves that the
789 log-likelihood is negatively correlated with its entropy and therefore it justifies the choice of $\lambda > 0$ in
790 the entropy minimisation.

791 **Proposition E.4.** *The log-likelihood of the true distribution $\log p(y|x)$ is negatively correlated with*
792 *its entropy $\mathbb{H}_{\tilde{y}}(p(\tilde{y}|x)) = -\mathbb{E}_{\tilde{y} \sim p(\cdot|x)}[\log p(\tilde{y}|x)]$.*

$$\text{Cov}(\log p(y|x), \mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))) < 0 \quad (24)$$

Proof.

$$\text{Cov}(\log p(y|x), \mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))) = \mathbb{E}_{x,y}[\log p(y|x)\mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))] - \mathbb{E}_{x,y}[\log p(y|x)]\mathbb{E}_x[\mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))] \quad (25)$$

$$= -\mathbb{E}_{x,y}[\log p(y|x)\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] + \mathbb{E}_{x,y}[\log p(y|x)]\mathbb{E}_x[\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] \quad (26)$$

$$(27)$$

793 By the law of total expectation, we have that $\mathbb{E}_x[\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] = \mathbb{E}_{x,\tilde{y}}[\log p(\tilde{y}|x)]$, then

$$\text{Cov}(\log p(y|x), \mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))) = -\mathbb{E}_{x,y}[\log p(y|x)\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] + \mathbb{E}_{x,y}[\log p(y|x)]^2 \quad (28)$$

$$= \mathbb{E}_{x,y}[\log p(y|x)]^2 - \mathbb{E}_{x,y}[\log p(y|x)\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] \quad (29)$$

$$(30)$$

794 On the other hand, also with the law of total expectation, $\mathbb{E}_{x,y}[\log p(y|x)\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] =$
795 $\mathbb{E}_x[\mathbb{E}_{y|x}[\log p(y|x)]\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]]$, so

$$\begin{aligned}\mathbb{E}_{x,y}[\log p(y|x)\mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] &= \mathbb{E}_x[\mathbb{E}_{y|x}[\log p(y|x)]^2] \\ &\geq \mathbb{E}_x[\mathbb{E}_{y|x}[\log p(y|x)]]^2 && \text{Jensen's inequality} \\ &\geq \mathbb{E}_{x,y}[\log p(y|x)]^2 && \text{total expectation law}\end{aligned}$$

796 Finally, we have the results,

$$\begin{aligned} \text{Cov}(\log p(y|x), \mathbb{H}_{\tilde{y}}(p(\tilde{y}|x))) &\leq \mathbb{E}_{x,y}[\log p(y|x)]^2 - \mathbb{E}_{x,y}[\log p(y|x)]^2 \\ &\leq 0 \end{aligned}$$

797

□

798 *Remark E.5.* We can also see the Pseudo-label as a form of entropy. Indeed, modulo the confidence
799 selection on the predicted probability, the Pseudo-label objective is the inverse of the Rényi min-
800 entropy:

$$\mathbb{H}_{\infty}(x) = -\max_y \log p(y|x)$$

801 **F Why debiasing with the labelled dataset?**

802 We remark that the debiasing can be performed with any subset of the training data, labelled and
 803 unlabelled. The choice of debiasing only with the labelled data can be explained both intuitively and
 804 computationally in regard to the Theorem 3.1. Intuitively, the debiasing term penalises the confidence
 805 on the labelled datapoints and then prevents the overfitting on the train dataset. As remarked in section
 806 3.1, Pereyra et al. (2017) showed that penalising low entropy models acts as a strong regulariser
 807 in supervised settings. This comforts the idea of penalising low entropy on the labelled dataset,
 808 i.e. debiasing the entropy minimisation with the labelled dataset. Considering Pseudo-Label-based
 809 methods, the objective for the labelled data is to predict the correct labels with moderate confidence.
 810 This is also similar to the concept of plausibility inference described by Barndorff-Nielsen (1976).

811 In regard to Theorem 3.1, we show that the optimum choice of subset for debiasing is either only the
 812 labelled data or the whole dataset and both are equivalent.

We consider a subset \mathcal{A} of the training set. We defined a as follow:

$$a_i = \begin{cases} 1/|\mathcal{A}| & \text{if } x_i \in \mathcal{A} \\ 0 & \text{otherwise} \end{cases} .$$

813 The unbiased estimator is then:

$$\hat{\mathcal{R}}_{DeSSL, \mathcal{A}}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \lambda \sum_{i=1}^n a_i H(\theta; x_i). \quad (31)$$

814

815 We compute the variance of this quantity as in the proof of Theorem 3.1 and show that:

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL, \mathcal{A}}(\theta) | r) = \sum_{i=1}^n \frac{r_i}{n_l^2} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \left(\frac{1-r_i}{n_u} - a_i \right)^2 \mathbb{V}(H(\theta, x)) - 2\lambda \frac{r_i a_i}{n_l} \text{Cov}(L(\theta, x, y), H(\theta, x)) \quad (32)$$

816 Suppose that no labelled datapoints are in \mathcal{A} . Then, the last term of the variance is null. Hence,
 817 having no labelled datapoints in \mathcal{A} leads to a variance increase. We also remark that debiasing with
 818 the entire dataset is equivalent that debiasing with the labelled datapoints. Indeed

$$\begin{aligned} \hat{\mathcal{R}}_{DeSSL}(\theta) &= \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{1}{n} \sum_{i=1}^n H(\theta; x_i) \\ &= \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda}{n} \sum_{i=1}^{n_l} H(\theta; x_i) - \frac{\lambda}{n} \sum_{i=1}^{n_u} H(\theta; x_i) \\ &= \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda n_l}{n n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda n_l}{n n_l} \sum_{i=1}^{n_l} H(\theta; x_i), \end{aligned}$$

819 which is equivalent to debiasing with only the labelled dataset by replacing λ by $\lambda \frac{n_l}{n}$.

820 At this point we can still sample a random subset composed of l labelled datapoints and u unlabelled
 821 datapoints. Therefore $a_i = 1/(l+u) \mathbf{1}\{x_i \in \mathcal{A}\}$, we show in the following that the optimum choice
 822 of the couple (l, u) are $(n_l, 0)$ and (n_l, n_U) , so only the labelled or the whole dataset.

We sample l labelled and u unlabelled datapoints to debiased the estimator, by simplifying the term
 in the sum of Equation 32 as follow:

$$\left(\frac{1-r_i}{n_u} - a_i \right)^2 = \begin{cases} \left(\frac{1}{n_u} - \frac{1}{l+u} \right)^2 & \text{if } x_i \in \mathcal{A} \text{ and } r_i = 0 \\ \frac{1}{(l+u)^2} & \text{if } x_i \in \mathcal{A} \text{ and } r_i = 1 \\ \frac{1}{n_u^2} & \text{if } x_i \notin \mathcal{A} \text{ and } r_i = 0 \\ 0 & \text{if } x_i \notin \mathcal{A} \text{ and } r_i = 1 \end{cases} .$$

823 Then, by summing the term and simplifying, we get:

$$\begin{aligned}\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) &= \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \left[u \left(\frac{1}{n_u} - \frac{1}{l+u} \right) + \frac{l}{(l+u)^2} + \frac{n_u - u}{n_u^2} \right] \mathbb{V}(H(\theta, x)) \\ &\quad - 2\lambda \frac{l}{n_l(l+u)} \text{Cov}(L(\theta, x, y), H(\theta, x)) \\ &= \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) + \lambda^2 \frac{n_l}{n_u} \frac{n_u - u + l}{l+u} \mathbb{V}(H(\theta, x)) - 2\lambda \frac{l}{n_l(l+u)} \text{Cov}(L(\theta, x, y), H(\theta, x))\end{aligned}$$

We want to minimise $\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r)$ with respect to (λ, l, u) . $\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r)$ reaches its minimum in λ at

$$\lambda_{opt} = \frac{n_u}{n_l} \frac{l}{n_u - u + l} \frac{\text{Cov}(L(\theta, x, y), H(\theta, x))}{\mathbb{V}(H(\theta, x))}.$$

824 Then,

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) = \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) - \frac{n_u}{n_l} \frac{l^2}{(n_u - u + l)(l+u)} \frac{\text{Cov}(L(\theta, x, y), H(\theta, x))^2}{\mathbb{V}(H(\theta, x))}.$$

825 We now want to minimise with respect to $0 \leq u \leq n_u$ and $1 \leq l \leq n_l$. We can easily show that the
826 $(n_u - u + l)(l+u)$ reaches its minimum for $u = 0$ or $u = n_u$ and for both value:

$$\mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta)|r) = \frac{1}{n_l} \mathbb{V}(L(\theta, x, y)) - \frac{n_u}{n_l} \frac{l}{n_u + l} \frac{\text{Cov}(L(\theta, x, y), H(\theta, x))^2}{\mathbb{V}(H(\theta, x))}.$$

827 Then $l/(n_u + l)$ is an increasing function, then reaches its maximum at $l = n_l$. So finally, the optimal
828 choices for the couple $(n_l, 0)$ and (n_l, n_u) . We showed that these couples are equivalent.

829 **G Proof of Theorem 3.2**

830 **Theorem 3.2** *If $\mathcal{S}(p_\theta, (x, y)) = -L(\theta; x, y)$ is a proper scoring rule, then*

$$\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x)\right) \quad (33)$$

831 *is also a proper scoring rule.*

Proof. The scoring rule considered in our SSL framework is:

$$\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x)\right).$$

The proper scoring rule of the fully supervised problem is

$$\mathcal{S}(p_\theta, (x, y, r)) = -L(\theta; x, y).$$

832 Let p be the true distribution of the data (x, y, r) . Under MCAR, r is independent of x and y , then
 833 $p(x, y, r) = p(r)p(x, y)$.

$$\begin{aligned} \mathcal{S}'(p_\theta, p) &= \int p(x, y, r)\mathcal{S}'(p_\theta, (x, y, r)) dx dy dr \\ &= \int p(x, y)p(r)\mathcal{S}'(p_\theta, (x, y, r)) dx dy dr && \text{by independence} \\ &= -\int p(x, y)p(r)\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x) dx dy dr \\ &= -\int_{x,y} p(x, y) \underbrace{\left(\int_r p(r)\frac{rn}{n_l} dr\right)}_{=1} L(\theta; x, y) dx dy \\ &\quad - \lambda n \int_{x,y} p(x, y) \underbrace{\left(\int_r p(r)\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right) dr\right)}_{=0} H(\theta; x) dx dy \\ &= -\int_{x,y} p(x, y)L(\theta; x, y) dx dy \\ &= \mathcal{S}(p_\theta, p) \end{aligned}$$

834 Therefore, if $\mathcal{S}(p_\theta, (x, y)) = -L(\theta; x, y)$ is a proper scoring rule, then
 835 *mathcal* $\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x)\right)$ is also a proper scoring rule.

836 □

837 **H Proof of Theorem 3.5**

838 Assumption 3.3: the minimum θ^* of \mathcal{R} is well-separated.

$$\inf_{\theta: d(\theta^*, \theta) \geq \epsilon} \mathcal{R}(\theta) > \mathcal{R}(\theta^*) \quad (34)$$

839 Assumption 3.4: uniform weak law of large numbers holds for a function L if:

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n L(\theta, x_i, y_i) - \mathbb{E}[L(\theta, x, y)] \right| \xrightarrow{p} 0 \quad (35)$$

840 **Theorem 3.5.** Under assumption A and assumption B for both L and H , $\hat{\theta} = \arg \min \hat{\mathcal{R}}_{DeSSL}$ is
841 asymptotically consistent with respect to n .

842 This result is a direct application of Theorem 5.7 from van der Vaart (2000, Chapter 5) that states
843 that under assumption A and B for L , $\hat{\theta} = \arg \min \hat{\mathcal{R}}$ is asymptotically consistent with respect to n .
844 Assumption A remains unchanged as we have M-estimators of the same \mathcal{R} . We now aim to prove that
845 under assumption B for both L and H , we have the assumption B on $\theta \rightarrow \frac{rn}{n_l} L(\theta; x, y) + \lambda(1 -$
846 $\frac{rn}{n_l}) H(\theta; x)$.

847 **Lemma H.1.** If the uniform law of large number holds for both L and H , then it holds for $\theta \rightarrow$
848 $\frac{rn}{n_l} L(\theta; x, y) + \lambda(1 - \frac{rn}{n_l}) H(\theta; x)$.

849 *Proof.* Suppose assumption B for L , then the same result holds if we replace n with n_l as n and n_l
850 are coupled by the law of r . Indeed, when n grows to infinity, n_l too and inversely. Therefore,

$$\sup_{\theta \in \Theta} \left| \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) - \mathbb{E}[L(\theta; x, y)] \right| \xrightarrow{p} 0$$

851 Now, suppose we have assumption B for H , then we can make the same remark than for L . Now, we
852 have to show that:

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \frac{rn}{n_l} L(\theta; x, y) + \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) - \mathbb{E}[L(\theta; x, y)] \right| \xrightarrow{p} 0$$

853 We first split the absolute value and the sup operator as

$$\begin{aligned} & \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \frac{rn}{n_l} L(\theta; x, y) + \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) - \mathbb{E}[L(\theta; x, y)] \right| \\ & \leq \sup_{\theta \in \Theta} \left| \frac{1}{n_l} \sum_{i=1}^n \frac{rn}{n_l} L(\theta; x, y) - \mathbb{E}[L(\theta; x, y)] \right| + \left| \frac{1}{n} \sum_{i=1}^n \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) \right| \\ & \leq \underbrace{\sup_{\theta \in \Theta} \left| \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x, y) - \mathbb{E}[L(\theta; x, y)] \right|}_{\xrightarrow{p} 0} + \sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) \right|. \end{aligned}$$

854 So we now have to prove that the second term is also converging to 0 in probability. Again by splitting
855 the absolute value and the sup, we have

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) \right| = \sup_{\theta \in \Theta} \left| \frac{\lambda}{n} \sum_{i=1}^n \frac{(1-r)n}{n_u} H(\theta; x) - \frac{\lambda}{n} \sum_{i=1}^n \frac{rn}{n_l} H(\theta; x) \right|$$

856 Then we have that,

$$\begin{aligned}
& \sup_{\theta \in \Theta} \left| \frac{\lambda}{n_u} \sum_{i=1}^n (1-r)H(\theta; x) - \frac{\lambda}{n_l} \sum_{i=1}^n rH(\theta; x) \right| \\
&= \sup_{\theta \in \Theta} \left| \frac{\lambda}{n_u} \sum_{i=1}^n (1-r)H(\theta; x) - \mathbb{E}[H(\theta; x, y)] - \left(\frac{\lambda}{n_l} \sum_{i=1}^n rH(\theta; x) - \mathbb{E}[H(\theta; x, y)] \right) \right| \\
&= \sup_{\theta \in \Theta} \left| \frac{\lambda}{n_u} \sum_{i=n_l+1}^{n_l+n_u} H(\theta; x) - \mathbb{E}[H(\theta; x, y)] - \left(\frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x) - \mathbb{E}[H(\theta; x, y)] \right) \right| \\
&\leq \underbrace{\sup_{\theta \in \Theta} \left| \frac{\lambda}{n_u} \sum_{i=n_l+1}^{n_l+n_u} H(\theta; x) - \mathbb{E}[H(\theta; x, y)] \right|}_{\xrightarrow{\frac{p}{n}} 0} + \underbrace{\sup_{\theta \in \Theta} \left| \left(\frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x) - \mathbb{E}[H(\theta; x, y)] \right) \right|}_{\xrightarrow{\frac{p}{n}} 0}.
\end{aligned}$$

857 Thus,

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n \frac{rn}{n_l} L(\theta; x, y) + \lambda n \left(\frac{1-r}{n_u} 1 - \frac{r}{n_l} \right) H(\theta; x) - \mathbb{E}[L(\theta; x, y)] \right| \xrightarrow{\frac{p}{n}} 0$$

858 And we now just have to apply the results of van der Vaart (2000, Theorem 5.7) to have the asymptotic
859 consistent of $\hat{\theta} = \arg \min \hat{\mathcal{R}}_{DeSSL}$.

860

□

861 *Remark H.2.* A sufficient condition on the function H to verify assumption B, the uniform weak
862 law of large numbers, is to be bounded (Newey & McFadden, 1994, Lemma 2.4). For instance,
863 the entropy $H = -\sum_y p_\theta(y|x) \log(p_\theta(y|x))$ is bounded and therefore, the entropy minimisation is
864 asymptotically consistent.

865 **I Asymptotic normality of DeSSL**

866 In the following, we study a modified version of the objective to simplify the proof. Let us consider
867 the following DeSSL objective $L'(\theta; x, y, r) = \frac{r}{\pi}L(\theta; x, y) + \lambda \left(\frac{1-r}{1-\pi} - \frac{r}{\pi} \right) H(\theta; x)$ which has the
868 same properties than the original one (unbiasedness, variance reduction property, consistency and
869 benefit from generalisation error bounds). The idea is to replace n_l with πn to simplify the expression.
870 The value n_l converges to πn then the following Theorem should hold with the true DeSSL objective.
871 We define the cross-covariance matrix between random vectors $\nabla L(\theta; x, y)$ and $\nabla H(\theta; x)$ as
872 $K_\theta(i, j) = \text{Cov}(\nabla L(\theta; x, y)_i, \nabla H(\theta; x)_j)$.

873 **Theorem I.1.** *Suppose L and H are smooth functions in $\mathcal{C}^2(\Theta, \mathbb{R})$. Assume $\mathcal{R}(\theta)$ admit a second-
874 order Taylor expansion at θ^* with a non-singular second order derivative V_{θ^*} . Under the MCAR
875 assumption, we have that $\hat{\theta}_{DeSSL}$ is asymptotically normal with covariance:*

$$\begin{aligned} \Sigma_{DeSSL} &= \frac{1}{\pi} V_{\theta^*}^{-1} \mathbb{E} [\nabla L(\theta^*; x, y) \nabla L(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\ &\quad + \frac{\lambda^2}{\pi(1-\pi)} V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x, y) \nabla H(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\ &\quad - \frac{\lambda}{\pi} V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1}. \end{aligned}$$

876 As a consequence, we can minimise the trace of the covariance. Indeed, $\text{Tr}(\Sigma_{DeSSL})$ reaches its
877 minimum at

$$\lambda_{opt} = (1 - \pi) \frac{\text{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})}{\text{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1})}, \quad (36)$$

878 and at λ_{opt} :

$$\text{Tr}(\Sigma_{DeSSL}) - \text{Tr}(\Sigma_{CC}) = -\frac{1 - \pi}{\pi} \frac{\text{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})^2}{\text{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1})} \leq 0. \quad (37)$$

879 The complete case is the special case of DeSSL with $\lambda = 0$. Then, the Theorem holds for the
880 complete case.

881 *Proof.* We define $L'(\theta; x, y, r) = \frac{r}{\pi}L(\theta; x, y) + \lambda \left(\frac{1-r}{1-\pi} - \frac{r}{\pi} \right) H(\theta; x)$ The assumptions of the
882 theorem are sufficient assumptions to apply Theorem 5.23 of Van der Vaart 1998 to the couple
883 $(\hat{\theta}_{DeSSL}, L')$. Hence, we obtain the following representation for representation $\hat{\theta}_{DeSSL}$:

$$\sqrt{n}(\hat{\theta}_{DeSSL} - \theta^*) = \frac{1}{\sqrt{n}} V_{\theta^*}^{-1} \sum_{i=1}^n \frac{r_i}{\pi} \nabla L(\theta^*; x_i, y_i) + \lambda \left(\frac{1-r_i}{1-\pi} - \frac{r_i}{\pi} \right) \nabla H(\theta^*; x_i) + o_p(1). \quad (38)$$

884

$$\sqrt{n}(\hat{\theta}_{DeSSL} - \theta^*) \xrightarrow{\mathcal{L}} \mathcal{N}(0, \Sigma_{DeSSL}),$$

885 The asymptotic normality follows with variance:

$$\Sigma_{DeSSL} = V_{\theta^*}^{-1} \mathbb{E} [\nabla L'(\theta^*; x, y) \nabla L'(\theta^*; x, y)^T] V_{\theta^*}^{-1}.$$

886 Using the MCAR assumption, we simplify the expression of Σ_{DeSSL} :

$$\begin{aligned}
\Sigma_{DeSSL} &= V_{\theta^*}^{-1} \mathbb{E} [\nabla L'(\theta^*; x, y) \nabla L'(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\
&= \frac{1}{\pi^2} V_{\theta^*}^{-1} \mathbb{E} [r \nabla L(\theta^*; x, y) \nabla L(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\
&\quad + \lambda^2 V_{\theta^*}^{-1} \mathbb{E} \left[\left(\frac{1-r}{(1-\pi)^2} + \frac{r}{\pi^2} \right) \nabla H(\theta^*; x, y) \nabla H(\theta^*; x, y)^T \right] V_{\theta^*}^{-1} \\
&\quad - \frac{\lambda}{\pi^2} V_{\theta^*}^{-1} \mathbb{E} [r \nabla L(\theta^*; x, y) \nabla H(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\
&= \frac{1}{\pi} V_{\theta^*}^{-1} \mathbf{Cov}(L(\theta^*; x, y)) V_{\theta^*}^{-1} + \frac{\lambda^2}{\pi(1-\pi)} V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x, y) \nabla H(\theta^*; x, y)^T] V_{\theta^*}^{-1} - \frac{\lambda}{\pi} V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1}.
\end{aligned}$$

887 We remark that the complete case is the particular case of DeSSL with $\lambda = 0$. Then,

$$\begin{aligned}
\Sigma_{DeSSL} &= \Sigma_{CC} + \frac{\lambda^2}{\pi(1-\pi)} V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x, y) \nabla H(\theta^*; x, y)^T] V_{\theta^*}^{-1} \\
&\quad - \frac{\lambda}{\pi} V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1}.
\end{aligned}$$

888 The asymptotic relative efficiency of consequence, the asymptotic relative efficiency $\hat{\theta}_{DeSSL}$ com-
889 pared to $\hat{\theta}_{CC}$ is defined as the quotient $\frac{\mathbf{Tr}(\Sigma_{DeSSL})}{\mathbf{Tr}(\Sigma_{CC})}$. This quotient can be minimised with respect to λ :
890

$$\lambda_{opt} = (1-\pi) \frac{\mathbf{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})}{\mathbf{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1})}, \quad (39)$$

891 and at λ_{opt} :

$$\frac{\mathbf{Tr}(\Sigma_{DeSSL})}{\mathbf{Tr}(\Sigma_{CC})} = 1 - \frac{1-\pi}{\pi} \frac{\mathbf{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})^2}{\mathbf{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1}) \mathbf{Tr}(\Sigma_{CC})} \leq 1. \quad (40)$$

892 □

893 **Remark I.2. On the sign of λ .** It is easy to show that a sufficient condition to have $\lambda_{opt} > 0$ is to
894 have K_{θ^*} positive semi-definite. Indeed, using that V_{θ^*} is positive definite and Proposition 6.1 of
895 Serre (2010), we show that $\mathbf{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1}) > 0$ and then $\lambda_{opt} > 0$.

896 **Remark I.3. Why minimising the trace of Σ_{DeSSL} ?** Minimising the trace of Σ_{DeSSL} leads to an
897 estimator with a smaller asymptotic MSE, see Chen et al. (2020).

Remark I.4. Fully supervised setting. We also remark that our theorem matches the theorem for the supervised setting. Indeed, observing all the labelled corresponds to the case $\pi = 1$ and we obtain:

$$\Sigma_{DeSSL} = \Sigma_{CC} = \Sigma_{\text{Fully supervised}}.$$

898 **J Proof of Theorem 3.6**

899 Our proof will be based on the following result from Shalev-Shwartz & Ben-David (2014, Theorem
900 26.5).

901 **Theorem J.1.** *Let \mathcal{H} be a set of parameters, $z \sim \mathcal{D}$ a random variable living in a space \mathcal{Z} , $c > 0$,
902 and $\ell : \mathcal{H} \times \mathcal{Z} \rightarrow [-c, c]$. We denote*

$$L_{\mathcal{D}}(h) = \mathbb{E}_z[\ell(h, z)], \text{ and } L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h, z_i), \quad (41)$$

903 where z_1, \dots, z_m are i.i.d. samples from \mathcal{D} . For any $\delta > 0$, with probability at least $1 - \delta$, we have

$$L_{\mathcal{D}}(h) \leq L_S(h) + 2\mathbb{E}_{(\varepsilon_i)_{i \leq m}} \left[\sup_{h \in \mathcal{H}} \left(\frac{1}{m} \sum_{i=1}^m \varepsilon_i \ell(h, z_i) \right) \right] + 4c \sqrt{\frac{2 \log(4/\delta)}{m}}, \quad (42)$$

904 where $\varepsilon_1, \dots, \varepsilon_m$ are i.i.d. Rademacher variables independent from z_1, \dots, z_m .

905 We can now restate and prove our generalisation bound.

906 **Theorem 3.6.** *We assume that both L and H are bounded and that the labels are MCAR. Then,
907 there exists a constant $\kappa > 0$, that depends on λ, L, H , and the ratio of observed labels, such that,
908 with probability at least $1 - \delta$, for all $\theta \in \Theta$,*

$$\mathcal{R}(\theta) \leq \hat{\mathcal{R}}_{DeSSL}(\theta) + 2R_n + \kappa \sqrt{\frac{\log(4/\delta)}{n}}, \quad (43)$$

909 where R_n is the Rademacher complexity

$$R_n = \mathbb{E}_{(\varepsilon_i)_{i \leq n}} \left[\sup_{\theta \in \Theta} \left(\frac{1}{n_l} \sum_{i=1}^{n_l} \varepsilon_i L(\theta; x_i, y_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} \varepsilon_i H(\theta; x_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} \varepsilon_i H(\theta; x_i) \right) \right], \quad (44)$$

910 with $\varepsilon_1, \dots, \varepsilon_m$ i.i.d. Rademacher variables independent from the data.

911 *Proof.* We use Theorem J.1 with $z = (x, y, r)$, $\mathcal{H} = \Theta$, $m = n$, and

$$\ell(h, z) = \frac{nr_i}{n_l} L(\theta; x_i, y_i) + \lambda \left(\frac{n(1-r_i)}{n_u} - \frac{nr_i}{n_l} \right) H(\theta; x_i). \quad (45)$$

912 The unbiasedness of our estimate under the MCAR assumption, proven in Appendix D, ensures that
913 the condition of Equation (41) is satisfied with $L_{\mathcal{D}}(h) = \mathcal{R}(\theta)$ and $L_S(h) = \hat{\mathcal{R}}_{DeSSL}(\theta)$. Now,
914 since L and H are bounded, there exists $M > 0$ such that $|L| < M$ and $|H| < M$. We can then
915 bound ℓ :

$$|\ell(h, z)| \leq \frac{n}{n_l} M + \lambda \max \left\{ \frac{n}{n_u}, \frac{n}{n_l} \right\} M = c. \quad (46)$$

916 Now that we have chosen a c that bounds ℓ , we can use Theorem J.1 and finally get Equation (43)
917 with $\kappa = 4c\sqrt{2}$. \square

918 **K DeSSL with H applied on all available data**

919 For consistency-based SSL methods it is common to use all the available data for the consistency
920 term:

$$\hat{\mathcal{R}}_{SSL}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n} \sum_{i=1}^n H(\theta; x_i). \quad (47)$$

921 With the same idea, we debias the risk estimate with the labelled data:

$$\begin{aligned} \hat{\mathcal{R}}_{DeSSL}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n} \sum_{i=1}^n H(\theta; x_i) \\ - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i). \end{aligned} \quad (48)$$

922 Under MCAR, this risk estimate is unbiased and the main theorem of the article hold with minor
923 modifications. In Theorem 3.1, λ_{opt} is slightly different and the expression of the variance at λ_{opt}
924 remains the same. The scoring rule in Theorem 3.2 is different but the theorem remains the same.
925 Both Theorem 3.5 and 3.6 remain the same with very similar proofs.

926 **Theorem K.1.** *The function $\lambda \mapsto \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta))$ reaches its minimum for:*

$$\lambda_{opt} = \frac{\text{Cov}(L(\theta; x, y), H(\theta; x))}{\mathbb{V}(H(\theta; x))} \quad (49)$$

927 *and*

$$\begin{aligned} \mathbb{V}(\hat{\mathcal{R}}_{DeSSL}(\theta))|_{\lambda_{opt}} &= \left(1 - \frac{n_u}{n} \rho_{L,H}^2\right) \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)) \\ &\leq \mathbb{V}(\hat{\mathcal{R}}_{CC}(\theta)) \end{aligned} \quad (50)$$

928 *where $\rho_{L,H} = \text{Corr}(L(\theta; x, y), H(\theta; x))$.*

929 When H is applied on all labelled and unlabelled data, the scoring rule used in the learning process
930 is then $\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{r n}{n_l} L(\theta; x, y) + \lambda \left(1 - \frac{r n}{n_l}\right) H(\theta; x)\right)$ and we have \mathcal{S}' is a proper scoring
931 rule.

932 **L MNIST and MedMNIST**

933 **L.1 MNIST**

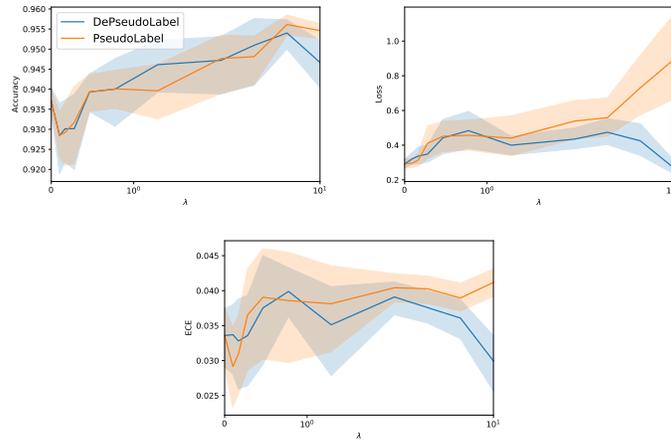


Figure 6: The influence of λ on Pseudo-label and DePseudo-label for a Lenet trained on MNIST with $n_l = 1000$: (Left) Test accuracy; (Middle) Mean test cross-entropy; (Right) Mean test ECE, with 95% CI

934 **L.2 MNIST label noise**

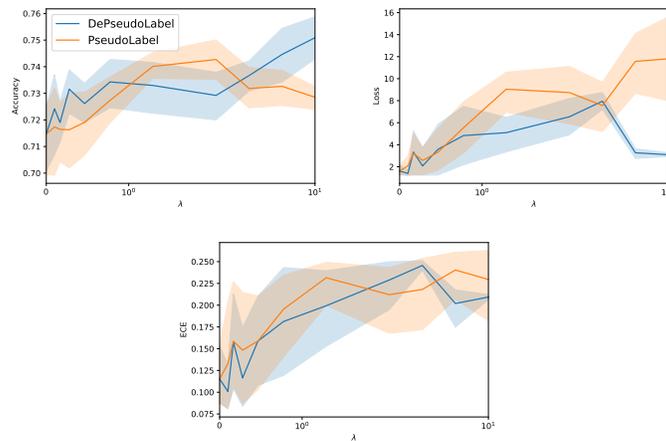


Figure 7: The influence of λ on Pseudo-label and DePseudo-label for a Lenet trained on MNIST with label noise with $n_l = 1000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

Table 3: Test AUC of Complete Case , PseudoLabel and DePseudoLabel on five datasets of MedMNIST.

DATASET	COMPLETE CASE	PSEUDOLABEL	DEPSEUDOLABEL
DERMA	84.26 ± 0.50	82.64 ± 1.19	83.82 ± 0.95
PNEUMONIA	94.28 ± 0.46	94.34 ± 0.91	94.15 ± 0.33
RETINA	70.70 ± 0.74	70.12 ± 1.01	69.97 ± 1.44
BREAST	74.67 ± 3.68	74.86 ± 3.18	75.33 ± 3.05
BLOOD	97.83 ± 0.23	97.83 ± 0.23	97.72 ± 0.15

936 **M PseudoLabel and DePseudoLabel on CIFAR: p-values**

937 **M.1 CIFAR-10**

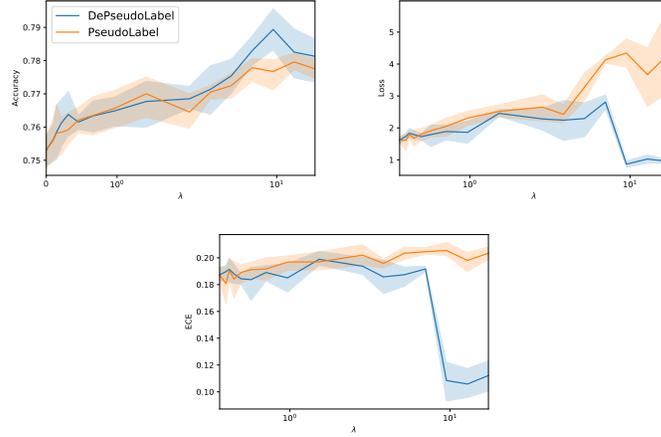


Figure 8: The influence of λ on Pseudo-label and DePseudo-label on CIFAR-10 with $n_l=4000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

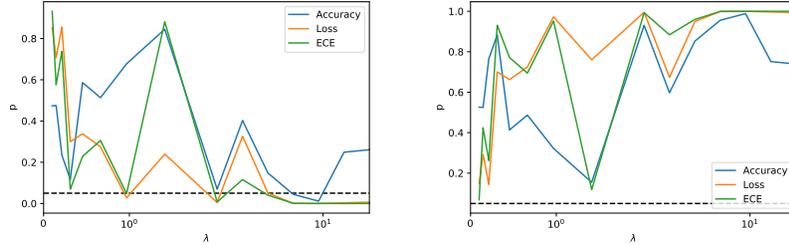


Figure 9: p-values of a paired student test between PseudoLabel and DePseudoLabel (Right) DePseudoLabel is better than PseudoLabel; (Left) DePseudoLabel is worse than PseudoLabel.

938 **M.2 Computation of λ_{opt} on the test set.**

939 As explained in the main text, the estimation of $Cov(L(\theta; x, y), H(\theta; x))$ with few labels led to
 940 extremely unstable unsatisfactory results. However, we test the formula on CIFAR-10 and different
 941 methods to provide intuition on the order of λ_{opt} and the range of the variance reduction regime
 942 (between 0 and $2\lambda_{opt}$). To do so, we estimate λ_{opt} on the test set for CIFAR-10 by training a CNN13
 943 using only 4,000 labelled data on 200 epochs. The value of λ_{opt} is 1.67, 31.16 and 0.66 for entropy
 944 minimisation, pseudo label and Fixmatch. Therefore, the reduced variance regime covers the intuitive
 945 choices of λ in the SSL literature. Unfortunately, computing λ_{opt} on the test set is not applicable in
 946 practice.

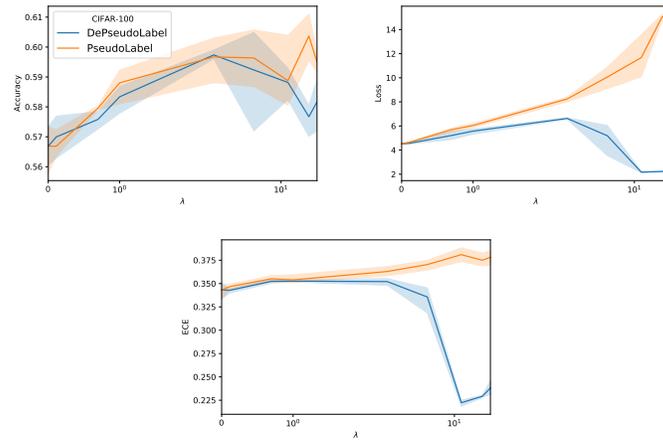


Figure 10: The influence of λ on Pseudo-label and DePseudo-label on CIFAR-100 with $n_l=4000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

948 **N Fixmatch (Sohn et al., 2020)**

949 **N.1 Per class accuracy**

950 In recent work, Zhu et al. (2022) exposed the disparate effect of SSL on different classes. Indeed,
 951 classes with a high complete case accuracy benefit more from SSL than classes with a low baseline
 952 accuracy. They introduced a metric called the benefit ratio (\mathcal{BR}) that quantifies the impact of SSL on
 953 a class C :

$$\mathcal{BR}(C) = \frac{acc_{SSL}(C) - acc_{CC}(C)}{acc_S(C) - acc_{CC}(C)}, \quad (51)$$

954 where $acc_{SSL}(C)$, $acc_{CC}(C)$ and $acc_S(C)$ are respectively the accuracy of the class with an SSL
 955 trained model, a complete-case model and a fully supervised model (a model that has access to all
 956 labels). Inspired by this work, we report the per class accuracy and the benefit ratio in Table N.1. We
 957 see that the “poor” classes such as bird, cat and dog tend to benefit from DeFixmatch much more
 958 than from Fixmatch. We compute $acc_S(C)$ using a pre-trained model with the same architecture¹.
 959 Zhu et al. (2022) also promote the idea that a fair SSL algorithm should benefit different sub-classes
 960 equally, then having $\mathcal{BR}(C) = \mathcal{BR}(C')$ for all C, C' . While perfect equality seems unachievable in
 961 practice, we propose to look at the standard deviation of the \mathcal{BR} through the different classes. While
 962 the standard deviation of Fixmatch is 0.12, the one of DeFixmatch is 0.06. Therefore, DeFixmatch
 963 improves the sub-populations’ accuracies more equally.

Table 4: Mean accuracy per class and mean benefit ratio (\mathcal{BR}) on 5 folds for Fixmatch, DeFixmatch and the Complete Case. Bold: “poor” complete case accuracy classes.

	COMPLETE CASE	FIXMATCH		DEFIXMATCH	
	ACCURACY	ACCURACY	\mathcal{BR}	ACCURACY	\mathcal{BR}
AIRPLANE	86.94	95.94	0.88	96.62	0.94
AUTOMOBILE	95.26	97.54	0.68	98.22	0.89
BIRD	80.46	90.80	0.68	92.64	0.80
CAT	70.08	82.50	0.56	87.16	0.78
DEER	88.88	95.86	0.78	97.26	0.94
DOG	79.66	87.16	0.53	90.98	0.81
FROG	93.12	97.84	0.80	98.62	0.94
HORSE	90.96	96.94	0.83	97.64	0.92
SHIP	94.12	97.26	0.67	98.06	0.84
TRUCK	93.18	96.82	0.84	97.20	0.93

964 **N.2 Fixmatch details**

As first detailed in Appendix B, Fixmatch is a pseudo-label based method with data augmentation. Indeed, Fixmatch uses weak augmentations of x (flip-and-shift) for the pseudo-labels selection and then minimises the likelihood with the prediction of the model on a strongly augmented version of x . Weak augmentations are also used for the supervised part of the loss. In this context,

$$L(\theta; x, y) = \mathbb{E}_{x_1 \sim weak(x)} [-\log(p_\theta(y|x_1))]$$

and

$$H(\theta; x) = \mathbb{E}_{x_1 \sim weak(x)} \left[\mathbb{1}[\max_y p_{\hat{\theta}}(y|x_1) > \tau] \mathbb{E}_{x_2 \sim strong(x)} [-\log(p_\theta(\arg \max_y p_{\hat{\theta}}(y|x_1)|x_2))] \right]$$

965 where x_1 is a weak augmentation of x and x_2 is a strong augmentation. We tried to debias an
 966 implementation of Fixmatch ¹ however training was very unstable and led to model that were much
 967 worst than the complete case. We believed that this behaviour is because the supervised part of

¹<https://github.com/LeeDoYup/FixMatch-pytorch>

968 the loss does not include strong augmentation. Indeed, our theoretical results encourage to have a
969 strong correlation between L and H , therefore including strong augmentations in the supervised term.
970 Moreover, a solid baseline for CIFAR-10 using only labelled data integrated strong augmentations
971 (Cubuk et al., 2020). We modify the implementation, see Code in supplementary materials. Therefore,
972 the supervised loss term can be written as:

$$L(\theta; x, y) = \frac{1}{2} (\mathbb{E}_{x_1 \sim \text{weak}(x)}[-\log(p_\theta(y|x_1))] + \mathbb{E}_{x_2 \sim \text{strong}(x)}[-\log(p_\theta(y|x_2))]), \quad (52)$$

973 where x_1 is a weak augmentation of x and x_2 is a strong augmentation. This modification encourages
974 us to choose $\lambda = \frac{1}{2}$ as the original Fixmatch implementation used $\lambda = 1$. We also remark that this
975 modification degrades the performance of Fixmatch (less than 2%) reported in the work of Sohn
976 et al. (2020). However, including strong augmentations in the supervised part greatly improves the
977 performance of the Complete Case.

978 **O CIFAR and SVHN: Oliver et al. (2018) implementation of**
 979 **consistency-based model.**

980 In this section, we present the results on CIFAR and SVHN by debiasing the implementation of
 981 (Oliver et al., 2018) of II-Model, Mean-Teacher and VAT ². We mimic the experiments of Oliver et al.
 982 (2018, figure-4) with the same configuration and the exact same hyperparameters (Oliver et al., 2018,
 983 Appendix B and C). We perform an early stopping independently on both cross-entropy and accuracy.
 984 As reported below, we reach almost the same results as the biased methods.

985 **O.1 CIFAR-10**

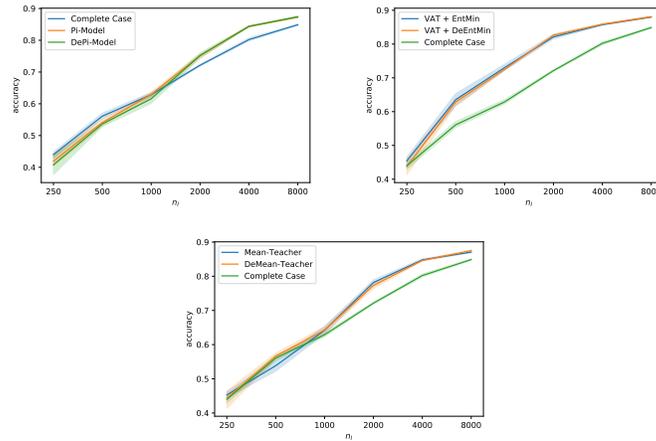


Figure 11: Test accuracy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) II-model and DeII-model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

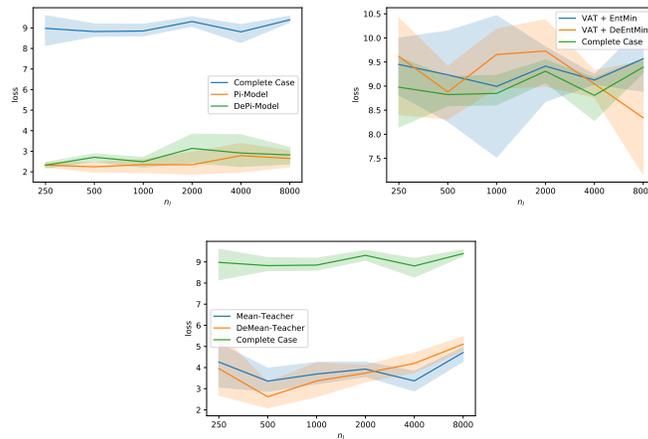


Figure 12: Test cross-entropy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) II-model and DeII-model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

²<https://github.com/brain-research/realistic-ssl-evaluation>

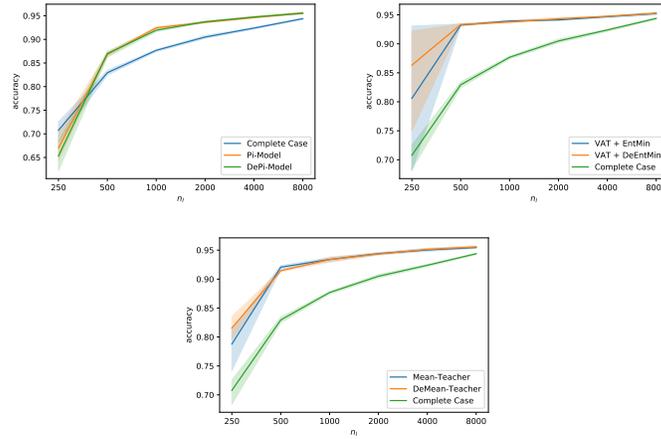


Figure 13: Test accuracy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) Π -model and De Π -model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

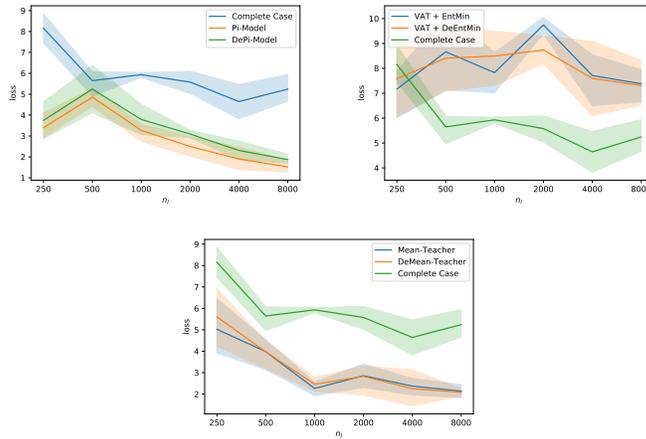


Figure 14: Test cross-entropy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) Π -model and De Π -model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

988 In this section, we tested these methods against the benchmarks of Chapelle et al., 2006, Chapter
 989 21 and UCI datasets already used in an SSL context in (Guo et al., 2010). We trained a logistic
 990 regression for the case of 100 labelled datapoints and finetune λ with a very small validation set, 20
 991 datapoints. We evaluated the performance in accuracy and cross-entropy of PseudoLabel, EntMin,
 992 DePseudoLabel and DeEntMin

993 **P.1 SSL Benchmark**

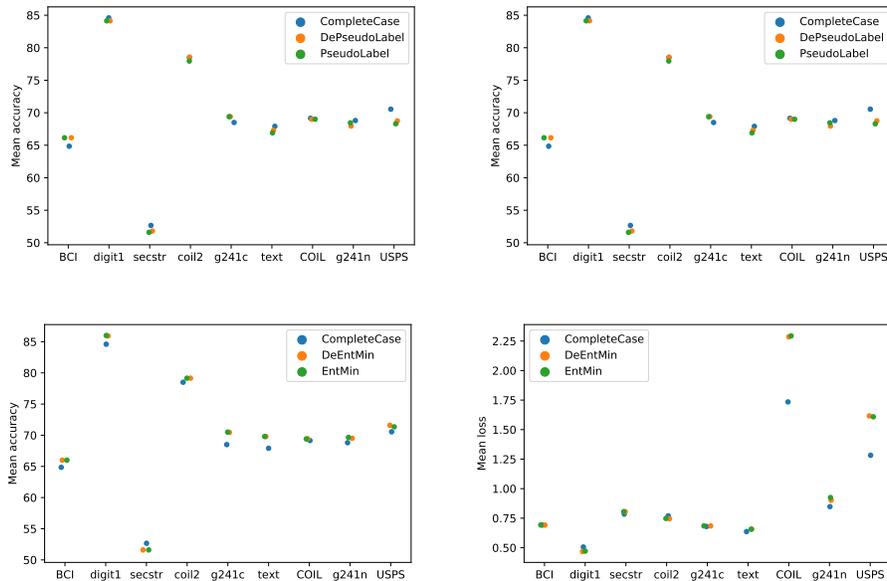


Figure 15: Mean accuracy and cross-entropy for each SSL datasets (Chapelle et al., 2006) on a logistic regression. (Top-Left) PseudoLabel and DePseudoLabel accuracy (Top-Right) PseudoLabel and DePseudoLabel cross-entropy (Bottom-Left) EntMin and DeEntMin accuracy (Bottom-Right) EntMin and DeEntMin cross-entropy.

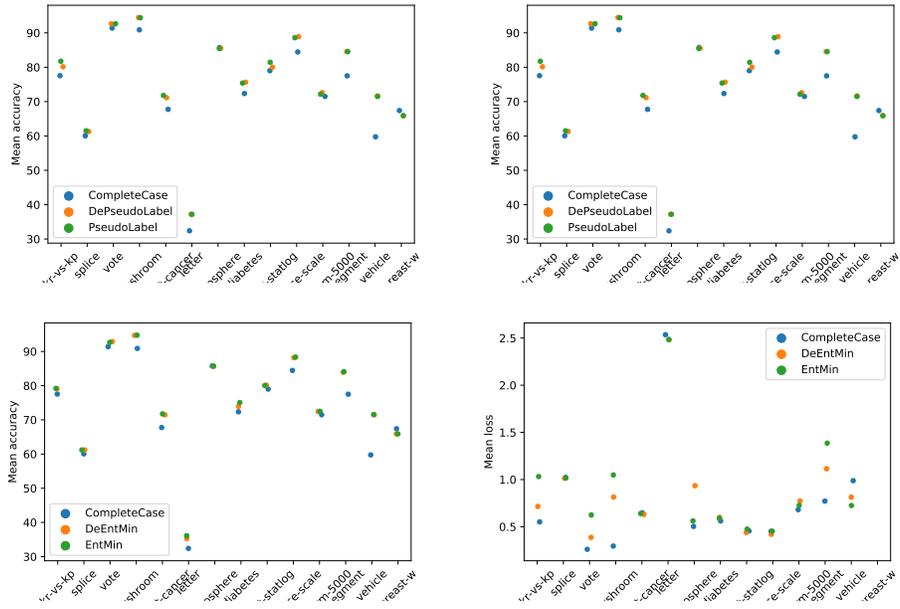


Figure 16: Mean accuracy and cross-entropy for each UCI datasets (Guo et al., 2010) on a logistic regression. (Top-Left) PseudoLabel and DePseudoLabel accuracy (Top-Right) PseudoLabel and DePseudoLabel cross-entropy (Bottom-Left) EntMin and DeEntMin accuracy (Bottom-Right) EntMin and DeEntMin cross-entropy.

995 **Q Computation details**

996 **Q.1 Computation resources**

997 Deep Learning experiments of this work required approximately 9,200 hours of GPU computation.
998 In particular, Fixmatch was trained using 4 GPUs. Here are the details:

- 999 • MNIST : 300 hours
- 1000 • medMNIST: 3 hours
- 1001 • CIFAR-10: 525 hours
- 1002 • CIFAR-100: 1500 hours
- 1003 • Fixmatch : 960 hours
- 1004 • Realistic SSL evaluation on both CIFAR and SVHN: 5880 hours

1005 **Q.2 Computation libraries and tools**

- 1006 • Python (Van Rossum & Drake Jr, 1995)
- 1007 • PyTorch (Paszke et al., 2019)
- 1008 • TensorFlow (Abadi et al., 2015)
- 1009 • Scikit-learn (Pedregosa et al., 2011)
- 1010 • Seaborn (Waskom et al., 2017)
- 1011 • Python imaging library (Lundh et al., 2012)
- 1012 • Numpy (Harris et al., 2020)
- 1013 • Pandas (McKinney et al., 2010)
- 1014 • RandAugment (Cubuk et al., 2020)
- 1015 • Fixmatch-Pytorch ³
- 1016 • Realistic-SSL-evaluation (Oliver et al., 2018)

³<https://github.com/LeeDoYup/FixMatch-pytorch>