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

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

Most semi-supervised methods rely on these distributional assumptions to ensure performance in entropy minimisation, pseudo-labelling and consistency-based methods. However, no proof is given that guarantees the effectiveness of state-of-the-art methods (Tarvainen & Valpola, 2017; Miyato et al., 2018; Sohn et al., 2020; Pham et al., 2021). To illustrate that SSL requires specific assumptions, we show in a toy example that pseudo-labelling fails at learning. To do so, we draw samples from two uniform distributions with a small overlap. Both supervised and semi-supervised neural networks are trained using the same labelled dataset. While the supervised algorithm learns perfectly the true distribution of $p(1|x)$, the semi-supervised learning methods (both entropy minimisation and pseudo-label) underestimate $p(1|x)$ for $x \in [1, 3]$ (see Figure 1). We also test our proposed method (DeSSL) on this dataset and show that the unbiased version of each SSL technique learns the true distribution accurately. See Appendix A for the results with Entropy Minimisation.

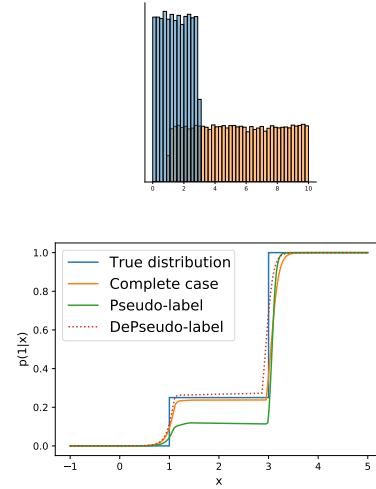


Figure 1: (Left) Data histogram. (Right) Posterior probabilities $p(1|x)$ of the same model trained following either complete case (only labelled data), Pseudo-label or our DePseudo-label. $n_l = 25,000$, $n_u = 25,000$.

1.1 Contributions

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

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

2 Semi-supervised learning

2.1 Learning with labelled data

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

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

2.2 Learning with both labelled and unlabelled data

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

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

2.2.1 Complete case: throwing the unlabelled data away

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

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

2.2.2 Incorporating unlabelled data

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

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

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

2.2.3 Some examples of surrogates

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

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

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

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

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

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

147 where the **Div** is a non-negative function that measures the divergence between two distributions, $\hat{\theta}$ is
 148 a fixed copy of the current parameter θ (the gradient is not propagated through $\hat{\theta}$) and $0 \leq C_1 \leq C_2$.

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

152 2.3 Theoretical guarantees

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

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

2.4 Related works

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

3 DeSSL: Unbiased semi-supervised learning

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

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

Under the MCAR assumption, this estimator is unbiased for any value of the parameter λ . For a proof of this result see Appendix D.

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

3.1 Does the DeSSL risk estimator make sense?

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

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

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

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

The debiasing term of our objective will penalise the confidence of the model on the labelled data. Pereyra et al. (2017) actually show that penalising the entropy in a supervised context acts as a strong regulator for supervised model and improves on the state-of-the-art on common benchmarks. This comforts us in the idea of debiasing using labelled data in the case of entropy-minimisation. Similarly, the debiasing term in pseudo-label turns the problem into plausability inference as described by Barndorff-Nielsen (1976).

Our objective also resembles doubly-robust risk estimates used for SSL in the context of kernel machines by Liu & Goldberg (2020) and for deep learning in a recent preprint (Hu et al., 2022). In both cases, their focus is quite different, as they consider weaker conditions than MCAR, but very specific choices of H .

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

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

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

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

and at λ_{opt} :

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

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

A proof of this theorem is available as Appendix E. This theorem provides a formal justification to the heuristic idea that H should be a surrogate of L . Indeed, DeSSL is a more accurate risk estimate when H is strongly positively correlated with L , which is likely to be the case when H is equal or equivalent to an expectation of L . Then, choosing λ positive is a coherent choice. We also demonstrate in Appendix E that L and H are positively correlated when L is the negative likelihood and H is the entropy. Other SSL methods have variance reduction guarantees and already has shown great promises in SSL, see Fox-Roberts & Rosten (2014) and Sakai et al. (2017). In a purely supervised context, Chen et al. (2020) show that the effectiveness of data augmentation techniques lays partially on the variance reduction of the risk estimate. A natural application of this theorem would be to tune λ automatically by estimating λ_{opt} . In our case however, the estimation of $\text{Cov}(L(\theta; x, y), H(\theta; x))$ with few labels led to extremely unstable unsatisfactory results.

3.3 Calibration

The calibration of a model is its capacity of predicting probability estimates that are representative of the true distribution. This property is determinant in real-world application when we need reliable predictions. A scoring rule \mathcal{S} is a function assigning a score to the predictive distribution

262 $p_\theta(y|x)$ 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
 263 e.g. Gneiting & Raftery, 2007). A scoring rule measures both the accuracy and the quality of
 264 predictive uncertainty, meaning that better calibration is rewarded. The expected scoring rule is
 265 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
 266 that $\mathcal{S}(p_\theta, p) \leq \mathcal{S}(p, p)$ (Gneiting & Raftery, 2007). The motivation behind having proper scoring
 267 rules comes from the following: suppose that the true data distribution p is accessible by our
 268 set of models. Then, the scoring rule encourages to predict $p_\theta = p$. The opposite of a proper
 269 scoring rule can then be used to train a model to encourage the calibration of predictive uncertainty:
 270 $L(\theta; x, y) = -\mathcal{S}(p_\theta, (x, y))$. Most common losses used to train models are proper scorings rule such
 271 as log-likelihood.

272 **Theorem 3.2.** *If $\mathcal{S}(p_\theta, (x, y)) = -L(\theta; x, y)$ is a proper scoring rule, then $\mathcal{S}'(p_\theta, (x, y, r)) =$
 273 $-(\frac{rn}{n_l} L(\theta; x, y) + \lambda n(\frac{1-r}{n_u} - \frac{r}{n_l}) H(\theta; x))$ is also a proper scoring rule.*

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

277 3.4 Consistency

278 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
 279 properties of $\hat{\theta}$ depend on the behaviours of the functions L and H . We will thus require the following
 280 standard assumptions.

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

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

283 **Theorem 3.5.** *Under the MCAR assumption, Assumption 3.3 and Assumption 3.4, $\hat{\theta} =$
 284 $\arg \min \hat{\mathcal{R}}_{DeSSL}$ is consistent.*

285 For a proof of this theorem see Appendix F. This theorem is a simple application of van der Vaart's
 286 (2000) Theorem 5.7 proving the consistency of a M-estimator. Also, this results holds for the complete
 287 case, with $\lambda = 0$ which prove that the complete case is a solid baseline under the MCAR assumption.

288 **Coupling of n_l and n_u under the MCAR assumption** Under the MCAR assumption, n_l and n_u
 289 are random variables. We have that $r \sim \mathcal{B}(\pi)$ (i.e. any x has the probability π of being labelled).
 290 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
 291 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
 292 go to infinity too and even if $n_u \gg n_l$.

293 3.5 Rademacher complexity and generalisation bounds

294 In this section, we prove an upper bound for the generalisation error of DeSSL. The unbiasedness of
 295 $\hat{\mathcal{R}}_{DeSSL}$ can directly be used to derive generalisation bounds based on the Rademacher complexity
 296 (Bartlett & Mendelson, 2002), defined in our case as

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

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

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

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

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

304 4 Experiments

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

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

323 4.1 MNIST

324 MNIST is an advantageous dataset for SSL since classes are
 325 well-separated. We compare PseudoLabel and DePseudoLabel
 326 for a LeNet-like architecture using $n_l = 1000$ labelled data
 327 on 10 different splits of the training dataset into a labelled and
 328 unlabelled set. Models are then evaluated using the standard
 329 10,000 test samples. We used 10% of n_l as the validation set.
 330 We test the influence of the hyperparameter λ and report the
 331 accuracy, the cross-entropy and the expected calibration error
 332 (ECE, Guo et al., 2017) at the epoch of best validation accuracy,
 333 see Figure 2 and Appendix J. In this example SSL and DeSSL
 334 have the almost the same accuracy for all λ , however, DeSSL
 335 seems to be always better calibrated. In order to break the cluster
 336 assumption, we reproduced the same experiment on a modified
 337 MNIST. Indeed, we had label noise by replacing the true label for
 338 20% of the dataset by a randomly sampled label, see Appendix
 339 J. In this setting, DeSSL performs better for large λ in term of
 340 accuracy and also provides a better calibration.

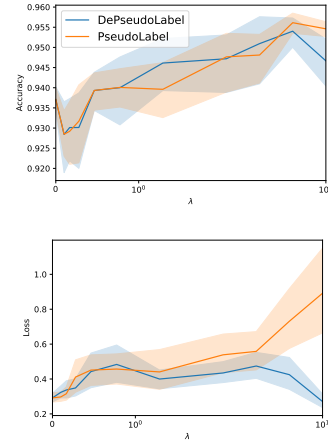


Figure 2: The influence of λ on Pseudo-label and DePseudo-label for a LeNet trained on MNIST with $n_l = 1000$: (Left) Mean test accuracy; (Right) Mean test cross-entropy, with 95% CI.

341 4.2 MedMNIST

342 We compare PseudoLabel and DePseudoLabel on different datasets of MedMNIST, a large-scale
 343 MNIST-like collection of biomedical images. We selected the five smallest 2D datasets of the
 344 collection, for these dataset it is likely that the cluster assumption no longer holds. We trained a
 345 5-layer CNN with a fixed $\lambda = 1$ and n_l at 10% of the training data. We report in Table 1 the mean
 346 accuracy and cross-entropy on 5 different splits of the labelled and unlabelled data and the number
 347 of labelled data used. We report the AUC in Appendix J. DePseudoLabel compete with PseudoLabel
 348 in terms of accuracy and even success when PseudoLabel’s accuracy is less than the complete
 349 case. Moreover, DePseudoLabel is always better in term of cross-entropy, so calibration, whereas
 350 PseudoLabel is always worse than the complete case.

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

4.3 CIFAR

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

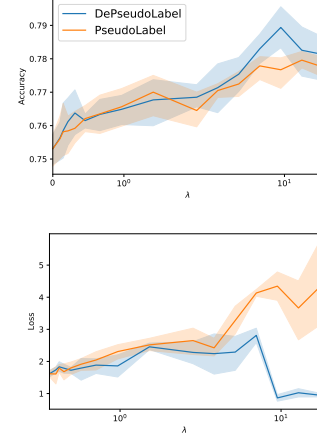


Figure 3: Influence of λ on Pseudo-label and DePseudo-label for a CNN trained on CIFAR with $n_l = 4000$: (Left) Mean test accuracy; (Right) Mean test cross-entropy, with 95% CI.

4.4 Fixmatch (Sohn et al., 2020)

We debiased a version of Fixmatch, see Appendix L for further details. For this experiment, we use $n_l = 4000$ on 5 different folds. First, we report that a strong baseline using data augmentation reach 87.27% accuracy. Then, we observe that on the debiasing method improve both accuracy and cross-entropy of this modified version of Fixmatch. Inspired by Zhu et al. (2022), we show that our method improved performance on “poor” classes more equally than the biased version. Indeed, DeFixmatch improves Fixmatch by 1.57% overall but by 4.91% on the worst class. We report in Appendix L 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

5 Conclusion

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

References

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

439 Grandvalet, Y. and Bengio, Y. Semi-supervised learning by entropy minimization. *Advances in*
440 *Neural Information Processing Systems*, 2004.

441 Guo, C., Pleiss, G., Sun, Y., and Weinberger, K. Q. On calibration of modern neural networks.
442 *International Conference on Machine Learning*, 2017.

443 Guo, L.-Z., Zhang, Z.-Y., Jiang, Y., Li, Y.-F., and Zhou, Z.-H. Safe deep semi-supervised learning for
444 unseen-class unlabeled data. *International Conference on Machine Learning*, 2020.

445 Guo, Y., Niu, X., and Zhang, H. An extensive empirical study on semi-supervised learning. *IEEE*
446 *International Conference on Data Mining*, 2010.

447 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser,
448 E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M. H., Brett, M.,
449 Haldane, A., Fernández del Río, J., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K.,
450 Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., and Oliphant, T. E. Array programming with
451 NumPy. *Nature*, 585:357–362, 2020. doi: 10.1038/s41586-020-2649-2.

452 Hu, X., Niu, Y., Miao, C., Hua, X.-S., and Zhang, H. On non-random missing labels in semi-
453 supervised learning. In *International Conference on Learning Representations*, 2022.

454 Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance
455 reduction. *Advances in Neural Information Processing Systems*, 2013.

456 Kawakita, M. and Takeuchi, J. Safe semi-supervised learning based on weighted likelihood. *Neural*
457 *Networks*, 53:146–164, 2014.

458 Kingma, D. P., Mohamed, S., Rezende, D. J., and Welling, M. Semi-supervised learning with deep
459 generative models. In *Advances in Neural Information Processing Systems*, pp. 3581–3589, 2014.

460 Krause, A., Perona, P., and Gomes, R. Discriminative clustering by regularized information maxi-
461 mization. *Advances in neural information processing systems*, 23, 2010.

462 Krizhevsky, A. Learning multiple layers of features from tiny images. Technical report, MIT, NYU,
463 2009.

464 Laine, S. and Aila, T. Temporal ensembling for semi-supervised learning. *International Conference*
465 *on Learning Representations*, 2017.

466 LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.

468 Lee, D.-H. Pseudo-Label : The simple and efficient semi-supervised learning method for deep
469 neural networks. *Workshop on challenges in representation learning, International conference on*
470 *machine learning*, 2013.

471 Li, Y.-F. and Zhou, Z.-H. Towards making unlabeled data never hurt. *IEEE transactions on pattern*
472 *analysis and machine intelligence*, 37:175–188, 2014.

473 Li, Y.-F., Kwok, J. T., and Zhou, Z.-H. Towards safe semi-supervised learning for multivariate
474 performance measures. *AAAI Conference on Artificial Intelligence*, 2016.

475 Little, R. J. and Rubin, D. B. *Statistical Analysis with Missing Data*. John Wiley & Sons, 2019.

476 Liu, T. and Goldberg, Y. Kernel machines with missing responses. *Electronic Journal of Statistics*,
477 14:3766–3820, 2020.

478 Lundh, F., Ellis, M., et al. Python imaging library (pil), 2012.

479 McKinney, W. et al. Data structures for statistical computing in python. In *Proceedings of the 9th*
480 *Python in Science Conference*, volume 445, pp. 51–56. Austin, TX, 2010.

481 McLachlan, G. J. Estimating the linear discriminant function from initial samples containing a small
482 number of unclassified observations. *Journal of the American statistical association*, 72:403–406,
483 1977.

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

530 Sohn, K., Berthelot, D., Li, C.-L., Zhang, Z., Carlini, N., Cubuk, E. D., Kurakin, A., Zhang, H.,
531 and Raffel, C. FixMatch: Simplifying semi-supervised learning with consistency and confidence.
532 *Avances in Neural Information Processing Systems*, 2020.

533 Sokolovska, N., Cappé, O., and Yvon, F. The asymptotics of semi-supervised learning in discrimina-
534 tive probabilistic models. In *International Conference on Machine Learning*, 2008.

535 Tarvainen, A. and Valpola, H. Mean teachers are better role models: Weight-averaged consistency
536 targets improve semi-supervised deep learning results. *Advancer in Neural Information Processing*
537 *Systems*, 2017.

538 Trapp, M., Madl, T., Peharz, R., Pernkopf, F., and Trappl, R. Safe semi-supervised learning of
539 sum-product networks. *Conference on Uncertainty in Artificial Intelligence*, 2017.

540 Tsiatis, A. A. *Semiparametric theory and missing data*. Springer, 2006.

541 Tsuchiya, T., Charoenphakdee, N., Sato, I., and Sugiyama, M. Semisupervised ordinal regression
542 based on empirical risk minimization. *Neural Computation*, 33:3361–3412, 2021.

543 van der Vaart, A. W. *Asymptotic statistics*. Cambridge university press, 2000.

544 van Engelen, J. E. and Hoos, H. H. A survey on semi-supervised learning. *Machine Learning*, 109:
545 373–440, 2020.

546 Van Rossum, G. and Drake Jr, F. L. *Python reference manual*. Centrum voor Wiskunde en Informatica
547 Amsterdam, 1995.

548 Verma, V., Kawaguchi, K., Lamb, A., Kannala, J., Bengio, Y., and Lopez-Paz, D. Interpolation
549 consistency training for semi-supervised learning. *International Joint Conference on Artificial*
550 *Intelligence*, 2019.

551 Waskom, M., Botvinnik, O., O’Kane, D., Hobson, P., Lukauskas, S., Gemperline, D. C., Augspurger,
552 T., Halchenko, Y., Cole, J. B., Warmenhoven, J., de Ruiter, J., Pye, C., Hoyer, S., Vanderplas,
553 J., Villalba, S., Kunter, G., Quintero, E., Bachant, P., Martin, M., Meyer, K., Miles, A., Ram,
554 Y., Yarkoni, T., Williams, M. L., Evans, C., Fitzgerald, C., Brian, Fonnesbeck, C., Lee, A.,
555 and Qalieh, A. mwaskom/seaborn: v0.8.1 (september 2017), September 2017. URL <https://doi.org/10.5281/zenodo.883859>.
556

557 Wei, C., Shen, K., Chen, Y., and Ma, T. Theoretical analysis of self-training with deep networks on
558 unlabeled data. In *International Conference on Learning Representations*, 2021.

559 Xie, Q., Dai, Z., Hovy, E., Luong, M.-T., and Le, Q. V. Unsupervised data augmentation for
560 consistency training. *Advances in Neural Information Processing Systems*, 2019.

561 Yang, J., Shi, R., and Ni, B. MedMNIST classification decathlon: A lightweight AutoML benchmark
562 for medical image analysis. In *IEEE 18th International Symposium on Biomedical Imaging (ISBI)*,
563 pp. 191–195, 2021a.

564 Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., and Ni, B. MedMNIST v2: A
565 large-scale lightweight benchmark for 2D and 3D biomedical image classification. *arXiv preprint*
566 *arXiv:2110.14795*, 2021b.

567 Zhang, B., Wang, Y., Hou, W., Wu, H., Wang, J., Okumura, M., and Shinozaki, T. FlexMatch: Boost-
568 ing semi-supervised learning with curriculum pseudo labeling. *Advances in Neural Information*
569 *Processing Systems*, 2021a.

570 Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. mixup: Beyond empirical risk minimization.
571 *Internation Conference on Learning Representations*, 2017.

572 Zhang, S., Wang, M., Liu, S., Chen, P.-Y., and Xiong, J. How unlabeled data improve generalization
573 in self-training? A one-hidden-layer theoretical analysis. In *International Conference on Learning*
574 *Representations*, 2021b.

575 Zhu, X., Ghahramani, Z., and Lafferty, J. D. Semi-supervised learning using Gaussian fields and
576 harmonic functions. *International conference on machine learning*, 2003.

577 Zhu, Z., Luo, T., and Liu, Y. The rich get richer: Disparate impact of semi-supervised learning. In
578 *International Conference on Learning Representations*, 2022.

Checklist

1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes] See section 5
- (c) Did you discuss any potential negative societal impacts of your work? [N/A] Theoretical work
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]

2. If you are including theoretical results...

- (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices

3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] code and instructions to run Fixmatch.
- (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendices
- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Confidence intervals on all figures.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] An estimation in Appendix

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- (a) If your work uses existing assets, did you cite the creators? [Yes] See Appendix
- (b) Did you mention the license of the assets? [Yes]
- (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] Code
- (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
- (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]

5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

617 A Toy example

618 We trained a 4 layer neural network (1/20/100/20/1) with ReLU activation function using 25,000
 619 labelled and 25,000 unlabelled points draw from two 1D uniform laws with an overlap. We used
 620 $\lambda = 1$ and a confidence threshold for Pseudo-label $\tau = 0.70$. We optimised the model's weights
 621 using a stochastic gradient descent (SGD) optimiser with a learning rate of 0.1.

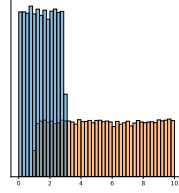


Figure 4: Data histogram

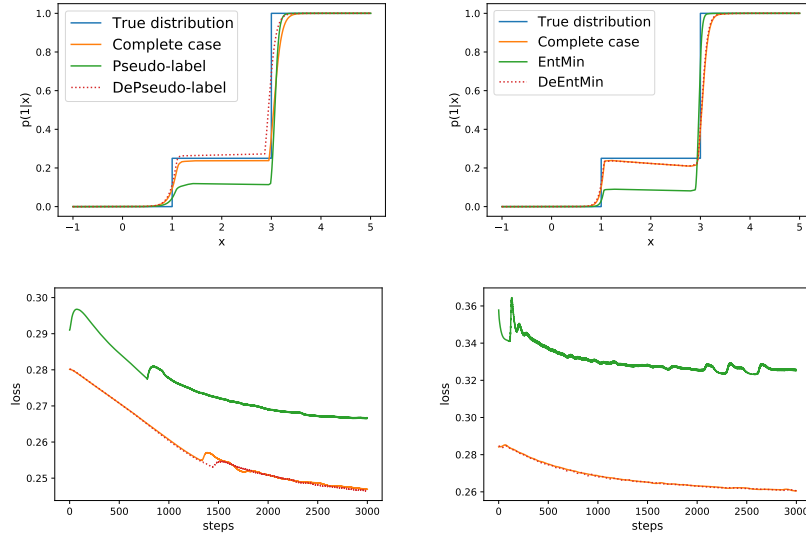


Figure 5: 4 layer neural net trained using SSL methods on a 1D dataset drawn from two uniform laws. (Top-left) Posterior probabilities $p(1|x)$ of the same model trained following either complete case (only labelled data), Pseudo-label or our DePseudo-label. (Top-right) Same for EntMin and DeEntMin (Bottom-left)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

B Details on surrogates and more examples

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

B.1 Entropy-based

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

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

Pseudo-label: As presented in the core article, the unsupervised objective of pseudo-label can be 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) encouraged the pseudo-labels method for deep semi-supervised learning. Then, Rizve et al. (2021) recently improved the pseudo-label selection by introducing an uncertainty-aware mechanism on the confidence of the model concerning the predicted probabilities. Pham et al. (2021) reaches state-of-the-art on the Imagenet challenge using pseudo-labels on a large dataset of additional images.

B.2 Pseudo-label and data augmentation

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

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

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

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

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

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

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

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

664 B.3 Consistency-based

665 Consistency-based method aim to smooth the decision function of the models or have more stable
 666 predictions. These objectives H are not directly a form of expectation of L but are equivalent to an
 667 expectation of L . For all the following methods we are able to write the unsupervised objective H
 668 such that:

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

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

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

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

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

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

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

681 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
 682 $\mathbb{E}_{\pi_x(\tilde{x}+r, \tilde{y})} [L(\theta; \tilde{x}, \tilde{y})]$. VAT is also a method between consistency-based and entropy-based method
 683 as long as we use the KL-divergence or the cross-entropy as the measure of divergence.

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

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

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

691 **II-Model** The II-Models are intrinsically stochastic models (for example a model with dropout)
 692 encouraged to make consistent prediction through several passes of the same x in the model. The
 693 SSL loss is using the stochastic behavior of the model where the model f_{θ} and penalises different
 694 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
 695 through the model f_{θ} . A generic form of the unsupervised part of the risk estimator is then

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

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

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

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

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

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

C On the semi-supervised bias

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

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

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

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

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

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

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

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

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

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

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

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

744 By the law of total expectation:

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

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

747 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
748 linearity of the expectation. Then,

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

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

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

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

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

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

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

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

753 *and*

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

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

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

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

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

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

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

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

759 Now, we remark that the variable r is binary and therefore $r^2 = r$, $(1-r)^2 = 1-r$ and $r(1-r) = 0$.

760 Using that and simplifying, we have

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

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

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

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

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

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

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

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

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

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

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

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

783 **Proposition E.4.** *The log-likelihood of the true distribution $\log p(y|x)$ is negatively correlated with*
784 *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.

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

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

$$(27)$$

785 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

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

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

$$(30)$$

786 On the other hand, also with the law of total expectation, $\mathbb{E}_{x,y}[\log p(y|x) \mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]] =$
787 $\mathbb{E}_x[\mathbb{E}_{y|x}[\log p(y|x)] \mathbb{E}_{\tilde{y}|x}[\log p(\tilde{y}|x)]]$, so

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

788 Finally, we have the results,

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

789

□

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

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

793 **F Proof of Theorem 3.2**

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

$$\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x)\right) \quad (31)$$

795 *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

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

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

797 $p(x, y, r) = p(r)p(x, y)$.

$$\mathcal{S}'(p_\theta, p) = \int p(x, y, r) \mathcal{S}'(p_\theta, (x, y, r)) dx dy dr \quad (32)$$

$$= \int p(x, y) p(r) \mathcal{S}'(p_\theta, (x, y, r)) dx dy dr \quad \text{by independence} \quad (33)$$

$$= - \int p(x, y) p(r) \frac{rn}{n_l} L(\theta; x, y) + \lambda n \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) H(\theta; x) dx dy dr \quad (34)$$

$$= - \int_{x,y} p(x, y) \underbrace{\left(\int_r p(r) \frac{rn}{n_l} dr \right)}_{=1} L(\theta; x, y) dx dy \quad (35)$$

$$- \lambda n \int_{x,y} p(x, y) \underbrace{\left(\int_r p(r) \left(\frac{1-r}{n_u} - \frac{r}{n_l} \right) dr \right)}_{=0} H(\theta; x) dx dy \quad (36)$$

$$= - \int_{x,y} p(x, y) L(\theta; x, y) dx dy \quad (37)$$

$$= \mathcal{S}(p_\theta, p) \quad (38)$$

798 Therefore, if $\mathcal{S}(p_\theta, (x, y)) = -L(\theta; x, y)$ is a proper scoring rule, then

799 $\mathcal{S}'(p_\theta, (x, y, r)) = -\left(\frac{rn}{n_l}L(\theta; x, y) + \lambda n\left(\frac{1-r}{n_u} - \frac{r}{n_l}\right)H(\theta; x)\right)$ is also a proper scoring rule.

800 □

801 **G Proof of Theorem 3.5**

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

$$\inf_{\theta: d(\theta^*, \theta) \geq \epsilon} \mathcal{R}(\theta) > \mathcal{R}(\theta^*) \quad (39)$$

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

$$\sup_{\theta \in \Theta} \left| \frac{1}{n} \sum_{i=1}^n L(\theta, x_i, y_i) - \mathbb{E}[L(\theta, x, y)] \right| \xrightarrow{p} 0 \quad (40)$$

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

806 This result is a direct application of Theorem 5.7 from van der Vaart (2000, Chapter 5) that states
807 that under assumption A and B for L , $\hat{\theta} = \arg \min \hat{\mathcal{R}}$ is asymptotically consistent with respect to n .
808 Assumption A remains unchanged as we have M-estimators of the same \mathcal{R} . We now aim to prove that
809 under assumption B for both L and H , we have the assumption B on $\theta \rightarrow \frac{rn}{n_l} L(\theta; x, y) + \lambda(1 -$
810 $\frac{rn}{n_l}) H(\theta; x)$.

811 **Lemma G.1.** If the uniform law of large number holds for both L and H , then it holds for $\theta \rightarrow$
812 $\frac{rn}{n_l} L(\theta; x, y) + \lambda(1 - \frac{rn}{n_l}) H(\theta; x)$.

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

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

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

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

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

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

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

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

820 Then we have that,

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

821 Thus,

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

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

824

□

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

H Proof of Theorem 3.6

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

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

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

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

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

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

We can now restate and prove our generalisation bound.

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

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

where R_n is the Rademacher complexity

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

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

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

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

The unbiasedness of our estimate under the MCAR assumption, proven in Appendix D, ensures that 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, since L and H are bounded, there exists $M > 0$ such that $|L| < M$ and $|H| < M$. We can then bound ℓ :

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

Now that we have chosen a c that bounds ℓ , we can use Theorem H.1 and finally get Equation (43)

with $\kappa = 4c\sqrt{2}$. \square

849 I DeSSL with H applied on all available data

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

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

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

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

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

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

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

858 and

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

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

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

863 J MNIST and MedMNIST

864 J.1 MNIST

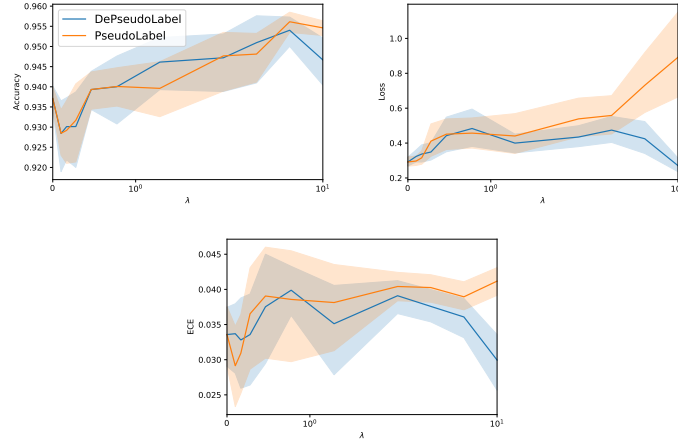


Figure 6: The influence of λ on Pseudo-label and DePseudo-label for a Lenet trained on MNIST with $n_l = 1000$: (Left) Test accuracy; (Middle) Mean test cross-entropy; (Right) Mean test ECE, with 95% CI

865 J.2 MNIST label noise

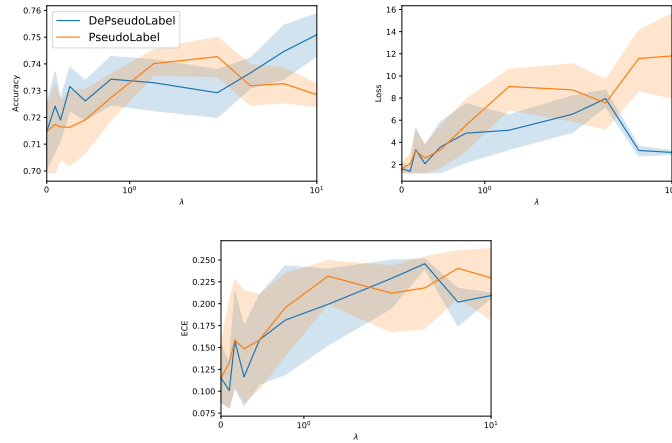


Figure 7: The influence of λ on Pseudo-label and DePseudo-label for a Lenet trained on MNIST with label noise with $n_l = 1000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

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

867 K PseudoLabel and DePseudoLabel on CIFAR: p-values

868 K.1 CIFAR-10

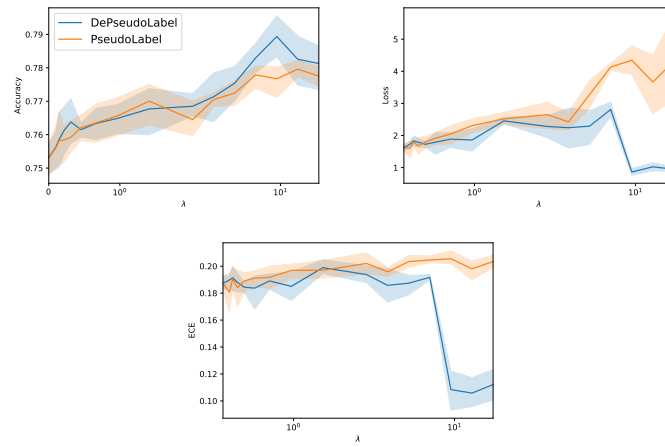


Figure 8: The influence of λ on Pseudo-label and DePseudo-label on CIFAR-10 with $n_l = 4000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

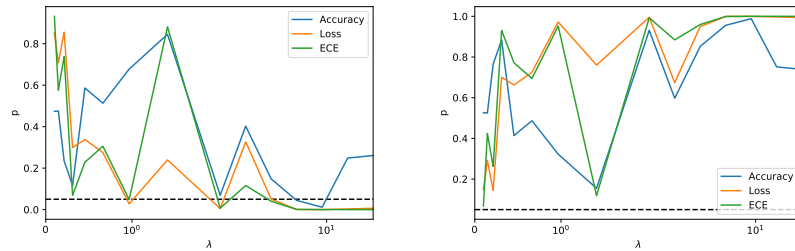


Figure 9: p-values of a paired student test between PseudoLabel and DePseudoLabel (Right) DePseudoLabel is better than PseudoLabel; (Left) DePseudoLabel is worst than PseudoLabel.

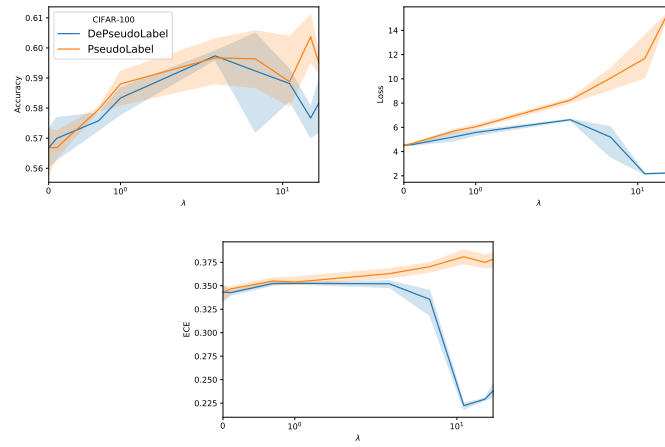


Figure 10: The influence of λ on Pseudo-label and DePseudo-label on CIFAR-100 with $n_l = 4000$: (Left) Mean test accuracy; (Middle) Mean test cross-entropy; (Right) Test ECE, with 95% CI.

L Fixmatch (Sohn et al., 2020)

L.1 Per class accuracy

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

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

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

L.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

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

where x_1 is a weak augmentation of x and x_2 is a strong augmentation. We tried to debias an implementation of Fixmatch¹ however training was very unstable and led to a model that was much worse than the complete case. We believed that this behaviour is because the supervised part of

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

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

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

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

M CIFAR and SVHN: Oliver et al. (2018) implementation of consistency-based model.

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

M.1 CIFAR-10

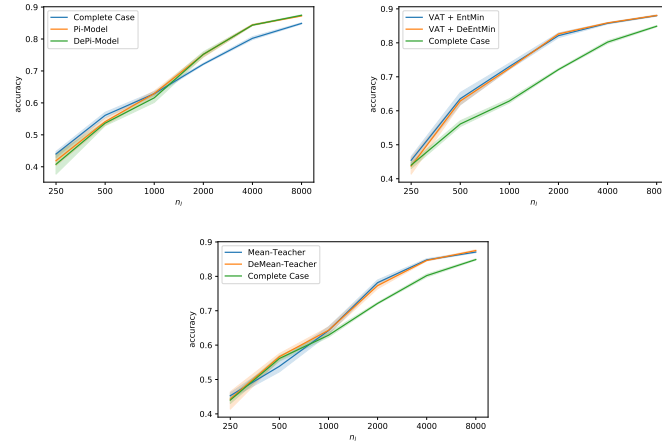


Figure 11: Test accuracy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) II-model and DeII-model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

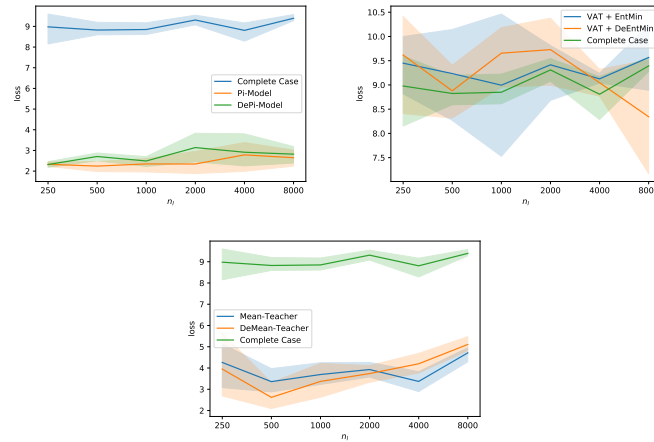


Figure 12: Test cross-entropy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) II-model and DeII-model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

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

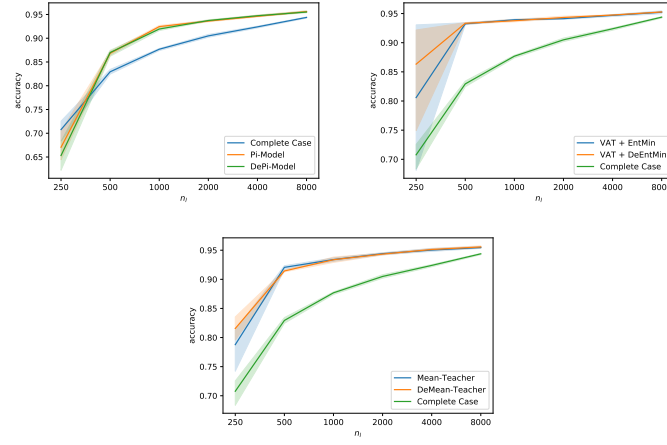


Figure 13: Test accuracy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) Π -model and De Π -model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

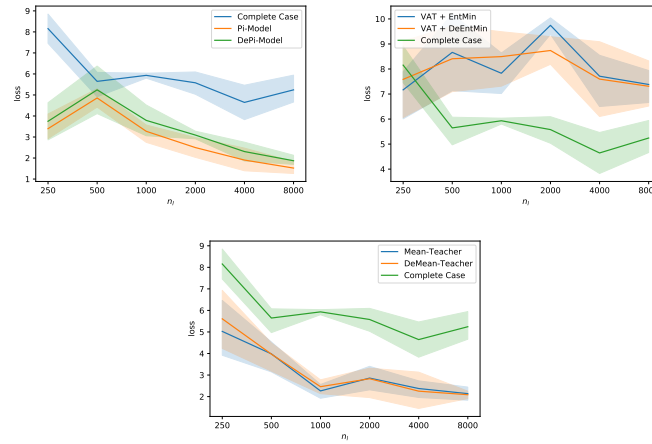


Figure 14: Test cross-entropy for each SSL approaches on CIFAR-10 with various amounts of labelled data n_l . (Left) Π -model and De Π -model. (Right) VAT+EntMin and VAT+DeEntMin. (Bottom) Mean-teacher and DeMean-teacher. Shadows represent 95% CI.

910 In this section, we tested these methods against the benchmarks of Chapelle et al., 2006, Chapter
 911 21 and UCI datasets already used in a SSL context in (Guo et al., 2010). We trained a logistic
 912 regression for the case of 100 labelled datapoints and finetune λ with a very small validation set, 20
 913 datapoints. We evaluated the performance in accuracy and cross-entropy of PseudoLabel, EntMin,
 914 DePseudoLabel and DeEntMin

915 N.1 SSL Benchmark

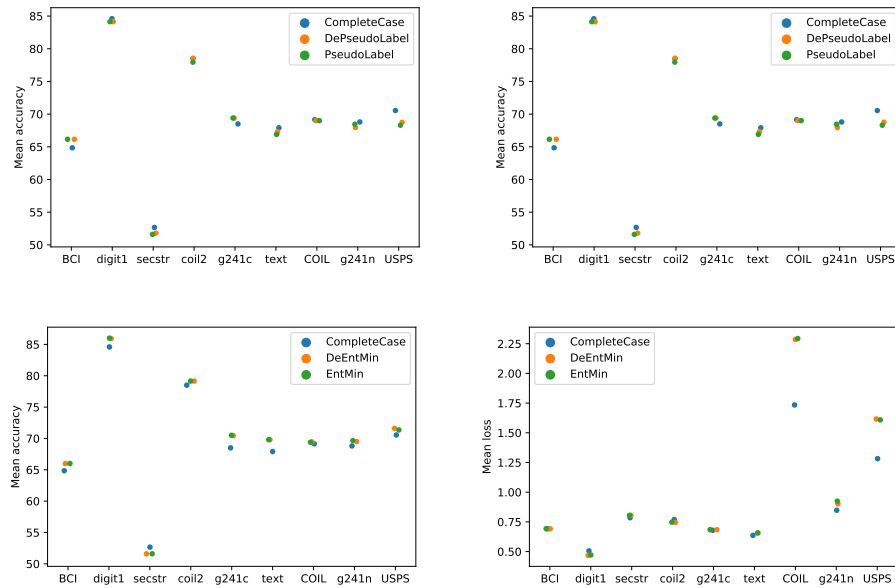


Figure 15: Mean accuracy and cross-entropy for each SSL datasets (Chapelle et al., 2006) on a logistic regression. (Top-Left) PseudoLabel and DePseudoLabel accuracy (Top-Right) PseudoLabel and DePseudoLabel cross-entropy (Bottom-Left) EntMin and DeEntMin accuracy (Bottom-Right) EntMin and DeEntMin cross-entropy.

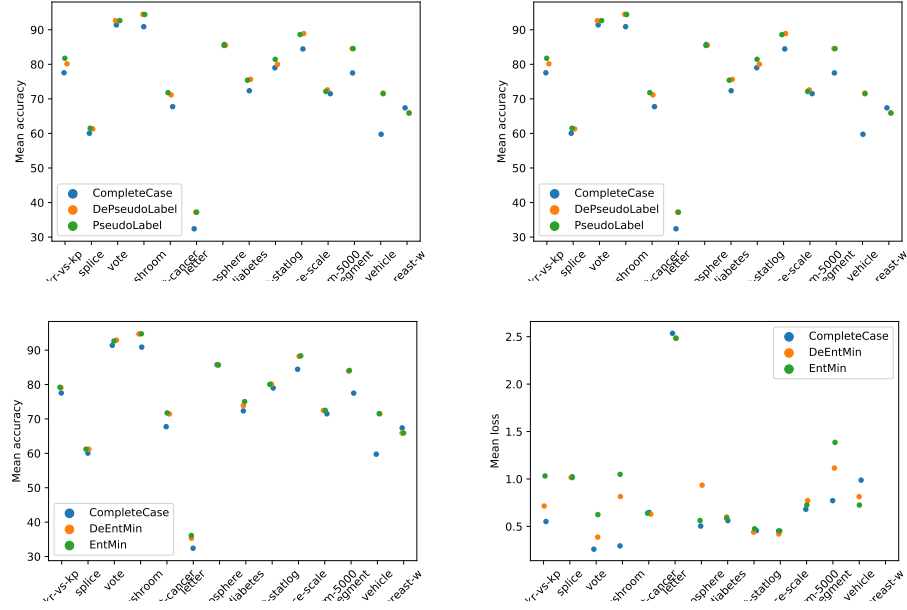


Figure 16: Mean accuracy and cross-entropy for each UCI datasets (Guo et al., 2010) on a logistic regression. (Top-Left) PseudoLabel and DePseudoLabel accuracy (Top-Right) PseudoLabel and DePseudoLabel cross-entropy (Bottom-Left) EntMin and DeEntMin accuracy (Bottom-Right) EntMin and DeEntMin cross-entropy.

917 **O Computation details**

918 **O.1 Computation resources**

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

- 921 • MNIST : 300 hours
- 922 • medMNIST: 3 hours
- 923 • CIFAR-10: 525 hours
- 924 • CIFAR-100: 1500 hours
- 925 • Fixmatch : 960 hours
- 926 • Realistic SSL evaluation on both CIFAR and SVHN: 5880 hours

927 **O.2 Computation libraries and tools**

- 928 • Python (Van Rossum & Drake Jr, 1995)
- 929 • PyTorch (Paszke et al., 2019)
- 930 • TensorFlow (Abadi et al., 2015)
- 931 • Scikit-learn (Pedregosa et al., 2011)
- 932 • Seaborn (Waskom et al., 2017)
- 933 • Python imaging library (Lundh et al., 2012)
- 934 • Numpy (Harris et al., 2020)
- 935 • Pandas (McKinney et al., 2010)
- 936 • RandAugment (Cubuk et al., 2020)
- 937 • Fixmatch-Pytorch ³
- 938 • Realistic-SSL-evaluation (Oliver et al., 2018)

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