



**HAL**  
open science

# Are labels informative in semi-supervised learning? Estimating and leveraging the missing-data mechanism

Aude Sportisse, Hugo Schmutz, Olivier Humbert, Charles Bouveyron,  
Pierre-Alexandre Mattei

## ► To cite this version:

Aude Sportisse, Hugo Schmutz, Olivier Humbert, Charles Bouveyron, Pierre-Alexandre Mattei. Are labels informative in semi-supervised learning? Estimating and leveraging the missing-data mechanism. International Conference on Machine Learning (ICML), 2023, Hawaii, United States. hal-03982898

**HAL Id: hal-03982898**

**<https://hal.science/hal-03982898v1>**

Submitted on 14 Feb 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

---

# Are labels informative in semi-supervised learning? Estimating and leveraging the missing-data mechanism

---

Aude, Sportisse<sup>1,2</sup> Hugo Schmutz<sup>1,2,3,4</sup> Olivier Humbert<sup>3,5</sup> Charles Bouveyron<sup>1,2</sup> Pierre-Alexandre Mattei<sup>1,2</sup>

## Abstract

Semi-supervised learning is a powerful technique for leveraging unlabeled data to improve machine learning models, but it can be affected by the presence of “informative” labels, which occur when some classes are more likely to be labeled than others. In the missing data literature, such labels are called missing not at random. In this paper, we propose a novel approach to address this issue by estimating the missing-data mechanism and using inverse propensity weighting to debias any SSL algorithm, including those using data augmentation. We also propose a likelihood ratio test to assess whether or not labels are indeed informative. Finally, we demonstrate the performance of the proposed methods on different datasets, in particular on two medical datasets for which we design pseudo-realistic missing data scenarios.

## 1. Introduction

Technological advancements have enabled the collection and storage of vast amounts of data, offering real hope for better prediction of phenomena. Unfortunately, this also leads to dirty data, and more specifically, missing data. In this paper, we focus on the scenario where a large amount of data is available, but labeling the data is costly, time-consuming, or even risky (for instance, medical data collection requires invasive tests for patients). Semi-supervised learning (SSL) (Chapelle et al., 2009; Van Engelen & Hoos, 2020) has emerged as a crucial problem to leverage both labeled and unlabeled data in predictive models. The unlabeled data are treated as observations with missing labels, as previously done in various studies (Grandvalet & Bengio, 2004; Ahfock & McLachlan, 2019; Hu et al., 2022; Schmutz et al., 2022). Recently, SSL algorithms have been extended to deep learning techniques demonstrating remarkable empirical successes, particularly through the systematic use of data augmentation (Berthelot et al., 2019a; Xie et al., 2020; Sohn et al., 2020; Rizve et al., 2021).

One of the challenges in semi-supervised learning is that the distribution of labels in the unlabeled dataset is unknown.

For instance, it is uncertain whether a class that is well-represented in the labeled images is also well-represented in the unlabeled images. The traditional approach is to assume that the label distributions are identical in the labeled and unlabeled datasets. This assumption implies that people label classes in equal proportions, regardless of the class nature or the quality of the images. However, it disregards the potential unbalance of popularity among classes. For example, in a medical context, doctors may prioritize labeling the class of sick patients or leave unlabeled the data with an ambiguous diagnosis. When the label distribution differs in the labeled and unlabeled datasets, the missing labels are said to be informative or Missing Not At Random (MNAR). The missingness of a label must be taken into account to obtain results from the available data that can be generalized to the entire population (Rubin, 1976). This is usually modeled by the missing-data mechanism, i.e. the probability of a sample to be observed (depending on the values of the label itself). Recently, it has been shown that classical SSL algorithms indeed fail to provide accurate results for the less observed classes in presence of informative labels. As there is a selection bias in the sample, MNAR data also raise the issue that some models can lead to non-identifiable parameters (Baker & Laird, 1988; Miao et al., 2016). A major challenge is that testing whether the data is indeed MNAR is difficult (d’Haultfoeuille, 2010), but it is necessary to provide a guideline for choosing which algorithm to apply. The main objective of this paper is to address these issues by estimating the missing-data mechanism.

Beyond the scope of deep learning, in the missing-data literature, significant works have considered the case of MNAR responses (Tang & Ju, 2018). Ibrahim & Lipsitz (1996) and Ibrahim et al. (2001) estimate the parameters of both the model and the missing-data mechanism using the Expectation-Maximization (EM) algorithm in binomial regression and generalized linear models. They also propose a likelihood ratio test statistic for selecting the variables related to the missingness but leave the identifiability of the parameters in perspective. It is in the semi-parametric setting that a lot of work has been done to obtain identification results, most often using a *shadow* variable (Miao et al., 2019; Miao & Tchetgen Tchetgen, 2016) that adds auxiliary information (Molenberghs et al., 2008). Some works

(Shao & Wang, 2016; Morikawa et al., 2017) also propose to debias classical estimators by using inverse probability weighting (IPW) techniques, weighting each sample by the inverse of its probability of being observed as determined by the missing-data mechanism. However, only the recent work of Hu et al. (2022) proposes an extension to deep learning, debiasing the risk estimator with a propensity score, but they do not directly model the missing-data mechanism, which is the main focus of our study (see Section 3.2 for a comprehensive comparison).

Our key contributions are summarized as follows:

- We consider a general self-masked MNAR model and prove its identifiability, showing in the process of identifiability of the model of (Hu et al., 2022).
- We propose two estimates of the missing-data mechanism and show their consistency.
- Based on these estimators, we propose an algorithm using IPW techniques able to debias any SSL algorithm in presence of informative labels.
- We provide a heuristic procedure to test whether the labels are indeed MNAR.
- We first demonstrate the efficiency of our methods on classical datasets. Furthermore, we propose two pseudo-realistic MNAR scenarios using medMNIST datasets (Yang et al., 2021). These contrast with the toy missing-data scenarios often used in existing works, even when the method is designed to handle informative labels.

## 2. Informative labels

### 2.1. Missing labels typology

In this paper, we study a dataset of  $n$  samples, denoted as  $D = (x_i, y_i)_{i=1}^n$ , where  $x_i$  represents the features and  $y_i$  represents the labels, which are drawn from the distribution  $p(x, y) = p(x)p(y|x)$ . Some of the labels are supposed to be missing, and thus the dataset is split into two subsets: a labeled dataset  $D_\ell = (x_i, y_i)_{i=1}^{n_\ell}$  of size  $n_\ell$  and an unlabeled dataset  $D_u = (x_i)_{i=n_\ell+1}^n$  of size  $n_u = n - n_\ell$ . The distribution of the labeled dataset is denoted as  $p^\ell(\cdot)$  (resp.  $p^u(\cdot)$  for the unlabeled dataset). In the following, we consider a discrete set of labels denoted as  $\mathcal{C} = \{1, \dots, K\}$ , with  $K$  the number of classes.

Most of the semi-supervised learning methods make the following assumption:

**A1.** The marginal distributions of the features and of the labels are identical in the labeled and unlabeled dataset, i.e.  $p^\ell(x) = p^u(x)$ ,  $\forall x$  and  $p^\ell(y) = p^u(y)$ ,  $\forall y$ .

Assumption **A1.** means that people label classes in equal proportions, regardless of their nature (label) or the quality of the images (features). Modeling two separate distribu-

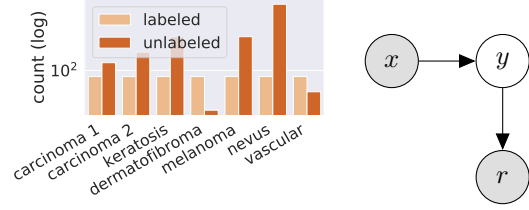


Figure 1. Left panel: MNAR labels for dermaMNIST (log count of labeled and unlabeled images per class). Right panel: Structural causal graph of the self-masked mechanism (Assumption **A2.**). The nodes in grey represent fully observed variables and the edges from  $x$  to  $y$  means that  $x$  causes  $y$ .

tions,  $p^\ell(\cdot)$  and  $p^u(\cdot)$ , is not always convenient, so we instead use a notation commonly used in missing-data studies. We introduce an additional random variable called missing-data indicator,  $r \in \{0, 1\}$ , where  $r = 1$  if  $y$  is observed and  $r = 0$  if  $y$  is missing. For example, this notation implies  $p^\ell(x) = p(x|r = 1)$  and  $p^u(x) = p(x|r = 0)$ . According to Rubin’s (1976) typology, labels can be: (i) Missing Completely At Random (MCAR) if the cause of the missingness is completely independent from the data values, i.e.  $r \perp\!\!\!\perp x, y$  (equivalent to Assumption **A1.**), (ii) Missing At Random (MAR) if the cause of the missingness can be explained by the features,  $r \perp\!\!\!\perp y|x$  and (iii) Missing Not At Random (MNAR) in all other cases. For example, labels will be MAR if medical doctors are less likely to label analyses that are of poor quality but they will be MNAR if they prefer to label the class of sick patients first. This last situation creates an “unbalanced class popularity” for which Assumption **A1.** is no longer valid.

In this paper, we focus on the MNAR case, specifically, the labels are assumed to be “self-masked MNAR”, meaning their unavailability only depends on their own values. This assumption allows us to model the unbalanced class popularity situation (see Figure 1) and is widely used in the missing-data literature (Mohan, 2018; Sportisse et al., 2020). Our assumption is formalized as follows:

**A2.** The labels are self-masked MNAR, i.e.  $r \perp\!\!\!\perp x|y$ .

Assumption **A2.** is weaker than Assumption **A1.**, since the equality of marginal distributions, either for features or labels, in both the labeled and unlabeled datasets, is relaxed. It only requires that the conditional distribution of the features given the class is the same in the labeled and unlabeled datasets (i.e.  $p^\ell(x|y) = p^u(x|y)$ ,  $\forall x, y$ ). For example, it does not cover the case where the radiography of sick patients does not have the same resolution whether it is labeled or not.

*Remark 2.1* (More general assumptions). The notation introduced here does not allow to consider different sets for the labels present in the labeled dataset and in the unlabeled

dataset. Label distribution mismatch has already been considered, such as when new classes appear in the unlabeled dataset (Guo et al., 2020; Cao et al., 2021) or when none of the classes present in the labeled dataset are present in the unlabeled dataset (Chen et al., 2020; Huang et al., 2021). However, these works are beyond the context of our work, as they do not allow to directly account for the informative nature of the labels, which is the main focus of this paper.

## 2.2. Non-ignorable missing-data mechanism

This typology of missing-data mechanism is important in determining the appropriate method to use: statistical inference can be performed on  $p(x, y)$  for MCAR or MAR labels, but it should be performed on  $p(x, y, r)$  for MNAR labels (Little & Rubin, 2019). We denote the parameter of interest,  $\theta \in \Theta$ , of  $p(y|x; \theta)$ . In most cases, this parameter corresponds to the weights of a neural network, or in simpler cases, a logistic regression. The parameter  $\phi$  of the missing data mechanism  $p(r|x, y; \phi)$  lives in  $\Phi = [0, 1]^K$ . In the following, we assume that the parameters are distinct, meaning that the joint parameter space is equal to  $\Theta \times \Phi$ . Following the common notation introduced by Le Morvan et al. (2020), the observed label vector is  $(y \odot r)$ , where  $\odot$  represents the term-by-term product, such that  $(y \odot r)_i = y_i$  if  $r_i = 1$  and  $(y \odot r)_i = \text{NA}$  if  $r_i = 0$ .

The traditional method for estimating  $\theta$  is to minimize the negative observed log-likelihood:

$$\begin{aligned} \ell(\theta, \phi) &= - \sum_{i=1}^n \log p(x_i, y_i \odot r_i, r_i; \theta, \phi) \\ &= - \sum_{i=1}^n \log p(r_i|x_i, y_i \odot r_i; \phi) p(y_i \odot r_i|x_i; \theta) p(x_i) \\ &= - \sum_{i=1}^n \begin{cases} \log p(r_i|x_i, y_i; \phi) p(y_i|x_i; \theta) & \text{if } r_i = 1 \\ \log \sum_{\tilde{y} \in \mathcal{C}} p(r_i|x_i, \tilde{y}; \phi) p(\tilde{y}|x_i; \theta) & \text{if } r_i = 0 \end{cases} + C \\ &= - \sum_{i=1}^n r_i \log p(y_i|x_i; \theta) + C \text{ under M(C)AR assumption,} \end{aligned} \quad (1)$$

where  $C$  is a constant independent of  $\theta$ . For M(C)AR labels, we use in the last step that  $p(r_i|x_i, y; \phi)$  does not depend on  $y$  and that  $\sum_{\tilde{y} \in \mathcal{C}} p(\tilde{y}|x_i; \theta) = 1$ ; the result follows if  $\phi$  is considered as a nuisance term. This simple calculation is a common technique in the missing data literature (Little & Rubin, 2019). It implies that for MCAR or MAR labels, it is not necessary to estimate the missing-data mechanism and minimizing the complete likelihood (on labeled data only) is sufficient. However, for MNAR labels, the missing-data mechanism cannot be ignored and must be taken into account.

## 2.3. Identification of the joint distribution

To fix ideas, Figure 1 (right panel) shows the causal relationships between the variables  $x, y$  and the missing-data indicator  $r$  through a structural causal graph (Neuberg, 2003). Assumption A2. allows us to get the nonparametric identification of the joint distribution  $p(y, x, r)$ , i.e. it can be expressed with quantities involving only observed data. Specifically, in the self-masked setting, the features act as shadow variables, providing enough auxiliary observed information to achieve the identifiability of the parameters.

**Proposition 2.2** (Identification of the joint distribution). *Under Assumptions A2. (self-masked MNAR), the joint distribution  $p(y, x, r)$  is identified.*

This result is a corollary of Theorem 1 in (Miao et al., 2019) and is proved in Appendix A. It is worth noting that this result also demonstrates the identification in (Hu et al., 2022).

## 3. Debiasing classical SSL algorithms

### 3.1. Complete-case: learning with labeled data

In classical supervised learning, the aim is to learn a predictive model  $p(y|x; \theta)$ , parametrized by  $\theta \in \Theta$ , by minimizing the theoretical risk:

$$\theta^* = \operatorname{argmin}_{\theta \in \Theta} \mathcal{R}(\theta) := \mathbb{E}_{(x,y) \sim p(x,y)} [\ell_\ell(\theta; x, y)],$$

where  $\ell_\ell$  is typically the negative log-likelihood function  $\ell_\ell(\theta; x_i, y_i) = -\log p(y_i|x_i; \theta)$  but can be any loss function. The theoretical risk is never observed, as it requires the knowledge of the true distribution  $p(x, y)$ . A typical learning procedure is then the minimization of the empirical risk, which is an unbiased estimate of the theoretical risk:

$$\hat{\theta} = \operatorname{argmin}_{\theta \in \Theta} \hat{\mathcal{R}}(\theta) := \frac{1}{n} \sum_{i=1}^n \ell_\ell(\theta; x_i, y_i).$$

This quantity is known in the supervised learning setting, but is still unobserved in the presence of missing labels. For MCAR labels, the natural estimator for  $\mathcal{R}(\theta)$  is the complete-case empirical risk, computed for the labeled data only:  $\hat{\mathcal{R}}^{CC}(\theta) := \frac{1}{n_\ell} \sum_{i=1}^n r_i \ell_\ell(\theta; x_i, y_i)$ . It is unbiased for MCAR labels but not in the other cases. For MAR labels, Liu & Goldberg (2020) propose to use the IPW estimator defined as  $\hat{\mathcal{R}}^{\text{IPW, MAR}}(\theta) := \frac{1}{n} \sum_{i=1}^n \frac{r_i \ell_\ell(\theta; x_i, y_i)}{\pi^{\text{MAR}}(x_i)}$ , where  $\pi^{\text{MAR}}(x) = \mathbb{P}(r = 1|x)$ . Similarly, for self-masked MNAR labels, we propose the following IPW estimator:

$$\hat{\mathcal{R}}_\phi(\theta) := \frac{1}{n} \sum_{i=1}^n \frac{r_i \ell_\ell(\theta; x_i, y_i)}{\phi_{y_i}}, \quad (2)$$

where  $\phi = (\phi_0, \dots, \phi_K) \in \Phi = [0, 1]^K$ , and  $\forall k \in \mathcal{C}$ ,  $\phi_k := \mathbb{P}(r = 1|y = k)$ .

The idea behind the IPW technique is that one labeled sample  $(x_i, y_i, r_i = 1)$  should not be counted only once but should take into account that there are unlabeled samples  $(x_j, y_j, r_j = 0), j \neq i$  that belong to the same class  $(y_j = y_i)$ . As a result, it is then counted  $1/\phi_{y_i}$  times. For example, if the probability of being observed in a class is one third, an observed sample from that class will be counted three times.

**Proposition 3.1** (Unbiasedness of the IPW estimator). *The IPW estimator proposed in (2) is unbiased, if the mechanism is well specified, i.e.  $\mathbb{E}[r|y] = \phi_y$ .*

*Proof.*

$$\begin{aligned} \mathbb{E} \left[ \frac{r}{\phi_y} \ell_\ell(\theta; x, y) \right] &= \mathbb{E} \left[ \mathbb{E} \left[ \frac{r}{\phi_y} \ell_\ell(\theta; x, y) | y \right] \right] \\ &= \mathbb{E} \left[ \mathbb{E} \left[ \frac{r}{\phi_y} | y \right] \mathbb{E} [\ell_\ell(\theta; x, y) | y] \right] \\ &= \mathbb{E} \left[ \frac{\mathbb{E}[r|y]}{\phi_y} \mathbb{E} [\ell_\ell(\theta; x, y) | y] \right] = \mathcal{R}(\theta), \end{aligned}$$

using  $r \perp\!\!\!\perp x|y$  in the second equality.  $\square$

### 3.2. Incorporating the unlabeled data

A major drawback of the classical IPW estimator in (2) is that it only uses labeled data and not all available data. To address this, traditional SSL algorithms for MCAR labels add a regularization term to the classical supervised objective:

$$\hat{\mathcal{R}}^{\text{SSL}}(\theta) := \frac{1}{n} \sum_{i=1}^n r_i \frac{\ell_\ell(\theta; x_i, y_i)}{n_\ell/n} + \frac{\lambda}{n} \sum_{i=1}^n (1 - r_i) \frac{\ell_u(\theta; x_i)}{n_u/n}, \quad (3)$$

where  $\lambda > 0$  is a regularization term. The function  $\ell_u$  is a loss function which does not depend on the labels; [Schmutz et al. \(2022\)](#) note that  $\ell_u$  can be viewed in many cases as a surrogate of  $\ell_\ell$ .

For example, [Grandvalet & Bengio \(2004\)](#) use the Shannon entropy. Another popular approach is to use “pseudo-labels” ([Rizve et al., 2021](#)) for unlabeled data by selecting the class with the highest posterior probability. Only the pseudo-labels that have a predicted probability higher than a predefined threshold  $\tau$  are used as targets. Both methods encourage the model to have a high level of confidence when imputing unlabeled data, but pseudo-label methods only use data points that have already been predicted with high confidence. Recently, state-of-the-art methods such as ([Sohn et al., 2020](#); [Berthelot et al., 2019a](#)) have also been developed to make the model more robust to data augmentation of the features.

For MNAR labels, the standard estimator in (3) is biased, and we propose the following estimator, which is unbiased

if the mechanism is correctly specified, using the same argument as Proposition 3.1:

$$\hat{\mathcal{R}}_\phi^{\text{SSL}}(\theta) := \frac{1}{n} \sum_{i=1}^n \frac{r_i \ell_\ell(\theta; x_i, y_i)}{\phi_{y_i}} - \frac{\lambda}{n} \sum_{i=1}^n \frac{r_i - \phi_{y_i}}{\phi_{y_i}} \ell_u(\theta; x_i). \quad (4)$$

This estimator has the significant advantage of being able to debias any SSL algorithm, including methods using data augmentation, with the knowledge of the weights  $\phi_{y_i}$ . This is the MNAR counterpart of the estimator proposed by [Schmutz et al. \(2022\)](#) for MCAR data and of the one suggested by [Liu & Goldberg \(2020\)](#) for MAR data with  $\lambda = 1$ . The only difference is the form of the mechanism: for MCAR data,  $\phi_{y_i} = n_\ell/n, \forall i$  and for MAR data, they use  $\pi^{\text{MAR}}(x)$  defined in Section 3.1 instead.

**Comparison with the work of [Hu et al. \(2022\)](#)** Our estimator also shares similarities with the “doubly-robust” estimate suggested by [Hu et al. \(2022\)](#), for debiasing the classical SSL estimators of the risk for self-masked MNAR labels (Assumption [A2.](#)). They build their estimator upon a very interesting strategy also used in the missing-data literature (see the recent review of [Rabe-Hesketh & Skrondal \(2022\)](#)) and leverage from it to indirectly account for the informative nature of the missing labels. The form of the risk estimator is as (4), but they target a composite likelihood (which does not encompass the cross entropy) by starting from  $\hat{\theta} \in \text{argmin}_\theta - \log p(x|y; \theta)$  instead of  $\hat{\theta} \in \text{argmin}_\theta - \log p(y|x; \theta)$ . The biggest advantage of their strategy is that it allows not to directly model the missing-data mechanism, which can be tedious in some missing data scenarios. Besides, the fundamental difference between our work and theirs is their weight does not involve the missing-data mechanism, but only the class proportions  $p(y)$ . Their method thus encourages the model to be accurate for the least frequent classes (when  $p(y)$  is small) but will not detect or favour the least labeled classes (when  $\phi_y$  is small). On the contrary, our method will benefit from the estimation of the missing-data mechanism, which can be obtained at no extra computational cost (see Section 4.5).

## 4. Estimating the missing-data mechanism

In Section 3, we proposed two unbiased estimates of the risk when the labels are informative. However, both require the knowledge of the missing-data mechanism. In practice, we will use the estimators given in (2) and (4) by plugging-in an estimation  $\hat{\phi}$  of the mechanism, resulting in  $\hat{\mathcal{R}}_{\hat{\phi}}$  and  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta)$ . In this section, we provide two estimators of the missing-data mechanism by using either the method of moments or the method of maximum likelihood.

#### 4.1. Moment estimator

A possible estimator of the missing-data mechanism is obtained by the method of moments applied to  $p(r = 1, y = y) = \mathbb{E}[\mathbb{1}_{\{r=1, y\}}] = \phi_y p(y)$ . It implies

$$\phi_y^M = \frac{\sum_{i=1}^n \mathbb{1}_{\{r=1, y_i=y\}}}{n} \frac{1}{p(y)}, \forall y \in \mathcal{C}. \quad (5)$$

This estimator allows us to leverage the information we have on the labeled data, because  $p(r = 1, y)$  is estimated by counting the number of labeled data in each class ( $\mathbb{1}_{\{r=1, y_i=y\}}$ ). The challenge now is to estimate the class distribution  $p(y)$ . This allows for two simple cases where the mechanism can be calculated directly: (i) when we know that the entire dataset is balanced (use  $p(y) = 1/K$ ) and (ii) when we have prior information on the class proportions (use  $p(y) = p_{\text{prior}}(y)$ ). This last case can happen when we have data from the general population (e.g. we know the prevalence rate of a disease).

When the class proportions are unknown, we propose to estimate  $p(y)$  as follows:

$$\hat{p}(y; \theta) = \frac{1}{n} \sum_{i=1}^n p(y_i | x_i; \theta), \forall \theta \in \Theta, \quad (6)$$

which is a consistent estimator of  $p(y)$  noting that  $p(y) = \int p(y|x; \theta)p(x)dx$ . The estimator of the missing-data mechanism is thus:

$$(\hat{\phi}_y^M)_\theta = \frac{\sum_{i=1}^n \mathbb{1}_{\{r=1, y_i=y\}}}{n} \frac{1}{\hat{p}(y, \theta)}, \forall y \in \mathcal{C}. \quad (7)$$

*Remark 4.1* (Computation of  $p_\theta(y)$ ). The estimator of the class proportions defined in (6) cannot be incorporated as such in a SGD algorithm, typically used to estimate  $\theta$ . We propose two ways to compute it within a mini-batch:

- use a moving averaging strategy inspired by (Hu et al., 2022), by using a buffer  $\hat{p}_{\text{buffer}}(y)$  updated at each iteration with  $\hat{p}(y, \theta = \theta_b)$  (6), where  $\theta_b$  is the parameter of the current mini-batch  $b^1$ :

$$\hat{p}_{\text{buffer}}(y) = \mu \hat{p}_{\text{buffer}}(y) + (1 - \mu) \hat{p}(y, \theta_b). \quad (8)$$

- propagate the gradients through  $(\hat{\phi}_y^M)_\theta$ .

#### 4.2. Maximum likelihood estimator

The second estimator of the missing-data mechanism relies on the method of maximum likelihood, already carried out by (Ibrahim & Lipsitz, 1996; Ibrahim et al., 2001) outside the scope of deep learning. It is obtained by minimizing the negative observed log-likelihood (1):

$$\hat{\theta}^L, \hat{\phi}^L = \operatorname{argmin}_{\theta \in \Theta, \phi \in \Phi} \ell(\theta, \phi; x, y \odot r, r). \quad (9)$$

<sup>1</sup>In this method, while  $(\hat{\phi}_y^M)_\theta$  depends on  $\theta$ , we do not propagate the gradients through  $\theta$ .

We highlight the following points:

- (Two-steps algorithm for SSL). Even if (9) gives an estimator  $\hat{\theta}^L$  of  $\theta$ , the latter can be really improved by incorporating the unlabeled data as in  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta)$  in a second step (see Algorithm 1 in Section 4.5).
- (MCAR setting). For not informative labels, the unlabeled data are not used for the estimation, as noted in Section 2.2. Besides, as expected, the minimum of the function is attained for a mechanism equal to the proportion of the labeled data. Indeed, we have  $\frac{\partial \ell(\theta, \phi_0)}{\partial \phi_0} = -\frac{n_l}{\phi_0} + \frac{n-n_l}{1-\phi_0}$  and  $\frac{\partial \ell(\theta, \phi_0)}{\partial \phi_0} = 0 \Leftrightarrow \phi_0 = \frac{n_l}{n}$ .
- (Convexity). The negative observed log-likelihood (1) is convex in  $\phi$ , for a fixed  $\theta \in \Theta$  (Appendix B.2).

*Remark 4.2* (Solving (9) in practice). To our knowledge, there is no closed form for the minimization problem. In practice, we propose to calculate the gradients by using the automatic differentiation package in PyTorch (Paszke et al., 2017). To comply with the constraint  $\phi \in \Phi$ , we consider  $\sigma(\phi_k) = \frac{1}{1+\exp(-\phi_k)}$ ,  $\forall k \in \mathcal{C}$  instead of  $\phi_k$ . In addition, we suggest solving (9) subject to the constraint of  $\sum_y \frac{\sum_{i=1}^n \mathbb{1}_{\{r=1, y_i=y\}}}{n} \frac{1}{\hat{\phi}_y^L} = 1$ , to comply with  $\sum_y p(y) = 1$  (see (5)), by using the `mdmm` package.

#### 4.3. Theoretical results

In this section, we provide theoretical results that validate the relevance of the chosen estimators. The proofs are detailed in Appendix B. We first demonstrate the consistency of the moment estimator for a fixed  $\theta \in \Theta$  by applying general results such as the law of large numbers and Slutsky's theorem.

**Proposition 4.3** (Consistency of  $\hat{\phi}^M$ ). *The moment estimator defined by (7) is consistent for a fixed  $\theta \in \Theta$ .*

Additionally, the consistency and asymptotic normality of the maximum likelihood estimator are obtained by applying Theorem 5.7 and Theorem 5.23 of (Van der Vaart, 2000) (stated in the more general case of M-estimators). We consider the negative observed log-likelihood for a fixed  $\theta$ , denoted as  $\ell_\theta : \phi \mapsto \ell(\theta, \phi; x, y \odot r, r)$ . In Appendix, we prove that the associated statistical model  $\mathcal{P}_\phi = \{p(r|y \odot r; \phi) : \phi \in \Phi\}$  is identifiable and that under interchangeability of differentiation w.r.t.  $\phi$  and integration over  $(x, y, r)$  (Assumption A3.), the Fisher information evaluated at the oracle estimate is invertible. Besides, we assume that  $\phi$  is in the interior of the set  $\Phi = [0, 1]^K$ , i.e. it cannot be on its boundary (Assumption A4.).

**Proposition 4.4** (Consistency, asymptotic normality of  $\hat{\phi}^L$ ). *Under Assumptions A3. and A4., the estimator  $\hat{\phi}^L$  is consistent and asymptotically normal.*

Finally, the consistency of the estimator of the missing-data

mechanism directly implies the consistency of the risks that we minimize in our SSL algorithms.

**Proposition 4.5** (Consistency of the risk). *If  $\hat{\phi}$  is a consistent estimator of  $\phi$  and if the mechanism is well specified, the risks  $\hat{\mathcal{R}}_{\hat{\phi}}$  and  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta)$  are consistent estimators of the theoretical risk  $\mathcal{R}(\theta) = \mathbb{E}[\ell_{\ell}(\theta; x, y)]$ .*

*Remark 4.6* (Consistency of  $\hat{\mathcal{R}}_{\hat{\phi}}$  using unlabeled data). As a consequence of the ignorability of the MCAR mechanism (see Section 2.2), the estimator of the theoretical risk using only labeled data is consistent in presence of MCAR labels. Proposition 4.5 shows that the IPW estimator  $\hat{\mathcal{R}}_{\hat{\phi}}$  is consistent for MNAR labels. It is worth noting that its expression refers only to labeled data but involves an estimator of the missing-data mechanism, computed on both labeled and unlabeled data (see Equations (7) and (9)). This underlines the relevance of unlabeled data in SSL with MNAR labels.

*Remark 4.7* (Double-robustness of the SSL risk  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}$ ). An interesting property is double-robustness, meaning that the estimator is consistent even if either the missing-data mechanism estimation or imputation is inaccurate. Hu et al. (2022) prove that double-robustness of their debiased risk (see Section 3.2) holds, if they assume that under inaccurate propensity estimation, the imputation is perfect (in the sense that the model always predicts the right class). This is a strong assumption. In our work, double-robustness (in that sense) of  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}$  is directly implied by Proposition 4.5 and by the unbiasedness of the risk. The strong assumption above could also be relaxed, applying Theorem 2 of (Miao et al., 2019) to our case: this is left as a perspective of our work.

#### 4.4. Testing the assumption on the mechanism

We present here a heuristics for estimating the missing-data mechanism to test if the labels are MCAR or not in the case of semi-supervised learning. The aim of such a test is to encourage the use of a specific method if the labels are not MCAR, or to support the selection of a traditional method if they are.

We want to test

$$H_0 : \phi \in \Phi^{\text{MCAR}} \text{ against } H_1 : \phi \notin \Phi^{\text{MCAR}},$$

where  $\Phi^{\text{MCAR}} = \{\phi \in [0, 1]^K, \forall k, k', \phi_k = \phi_{k'}\}$ . For the maximum likelihood estimator given in (9), we consider the following test statistic:

$$-2 \left( \inf_{\theta, \phi} \ell(\theta, \phi) - \inf_{\theta} \ell(\theta, \phi^{\text{MCAR}}) \right). \quad (10)$$

Under the same assumptions of Proposition 4.4, we know that the test stastic  $2(\ell_{\theta}(\phi) - \ell_{\theta}(\phi^{\text{MCAR}}))$ , for a fixed  $\theta \in \Theta$ , converges in distribution to a chi-squared random variable  $\chi_d^2$  (Theorem 16.7 of Van der Vaart (2000)), where  $d$  is the difference in degrees of freedom between the null hypothesis ( $H_0$ ) and the alternative hypothesis ( $H_1$ ), i.e.  $d = K - 1$

for  $K$  classes. We conjecture that an extension of Proposition 4.4 developed for a fixed  $\theta \in \Theta$  can be obtained by considering the profile likelihood  $\phi \mapsto \inf_{\theta \in \Theta} \ell(\theta, \phi)$ , and by applying the asymptotic results on it (Murphy & Van der Vaart, 2000). This implies that the test statistic (10) also converges in distribution to  $\chi_d^2$ . Based on this asymptotic distribution, it is possible to calculate a p-value from (10) and to easily test if the MCAR assumption is rejected.

#### 4.5. Algorithms

To ensure clarity, we explain how the proposed estimators can be applied to any SSL algorithm. As previously mentioned, the goal is to plug-in the estimation of the mechanism into the estimators of the theoretical risk.

The moment estimator presented in (7) is continuously updated throughout the SSL algorithm (Algorithm 2), and thus the estimation of the mechanism does not add any additional computational cost when using the moment estimator, as the estimation of the mechanism and the model are performed in a single step. On the other hand, when using the maximum likelihood method, the estimation process is divided into two steps: (i) the estimation of the mechanism by optimizing (1) (Algorithm 1) and (ii) the estimation of the model  $\theta$  using the SSL algorithm (Algorithm 2 using the estimator (i) as input for  $\hat{\phi}$ ). In both algorithms, the hyperparameters are classical: the sizes of the mini-batch ( $N_{\mathcal{B}}, N'_{\mathcal{B}}$ ), the learning rates ( $\gamma_{\phi}, \gamma_{\theta}, \gamma'_{\theta}$ ) and the number of epochs ( $N, N'$ ).

*Remark 4.8* (Adaptive threshold). SSL algorithms (Rizve et al., 2021; Sohn et al., 2020) that utilize pseudo-label techniques often employ a fixed threshold to select relevant imputations from unlabeled data. However, several recent studies (Hu et al., 2022; Wei et al., 2021; Berthelot et al., 2019a) have noted that an adaptive threshold can improve the performance of the classifier on the rarest classes, particularly in unbalanced semi-supervised learning. For instance, Hu et al. (2022) propose to use an adaptive threshold based on class proportions, setting higher requirements for popular classes and lower requirements for rare classes. Another possible adaptive threshold, suggested by our estimation of the missing-data mechanism, would depend on the missingness proportion of a class and set the highest requirement for the most observed class:

$$\forall k \in \mathcal{C}, \tau(y_k) = \tau_0 \left( \frac{\mathbb{P}(r = 1 | y = k)}{\max_y \mathbb{P}(r = 1 | y)} \right)^{\beta} = \tau_0 \left( \frac{\phi_y}{\max_y \phi_y} \right)^{\beta},$$

with  $\tau_0$  the classical threshold and  $\beta$  the hyper-parameter that determines how adaptive the cutoff is.

### 5. Numerical experiments

In this study, we evaluate the effectiveness of our proposed estimates of the missing-data mechanism using the benchmark dataset MNIST (LeCun & Cortes, 2010). Additionally, our debiased approach (Algorithm 2) of the classical SSL method using pseudo-labels (Rizve et al., 2021) is compared

**Algorithm 1** Maximum likelihood estimator for  $\phi$ 

**Input:** labeled data  $D_\ell$ , unlabeled data  $D_u$   
 Initialize  $\theta_0$  (at random),  $\phi_0$  (MCAR case:  $n_\ell/n$ ).  
**for**  $k = 0$  **to**  $N$  **iteratively do**  
   Sample a Mini-Batch  $\mathcal{B}$  of size  $N_B$  from  $D_\ell$  and from  $D_u$ .  
    $\phi_{k+1} = \phi_k - \gamma_\phi \partial_{\phi} \frac{1}{N_B} \sum_{i \in \mathcal{B}} \ell(\theta_k, \phi_k)$   
   Sample a Mini-Batch  $\mathcal{B}$  of size  $N_B$  from  $D_\ell$  and from  $D_u$ .  
    $\theta_{k+1} = \theta_k - \gamma_\theta \partial_{\theta} \frac{1}{N_B} \sum_{i \in \mathcal{B}} \ell(\theta_k, \phi_k)$   
**end for**  
**Output:**  $\phi_N = \hat{\phi}^L, \theta_N$

**Algorithm 2** Debaised SSL algorithm for informative labels

**Input:** labeled data  $D_\ell$ , unlabeled data  $D_u$ ,  $\hat{\phi}$  (if available)  
 Initialize  $\theta_0$  (at random)  
**for**  $k = 1$  **to**  $N'$  **do**  
   Sample a Mini-Batch  $\mathcal{B}$  of size  $N'_B$  from  $D_\ell$  and from  $D_u$ .  
   **if**  $\hat{\phi}$  is not provided **then**  
     Compute  $\hat{\phi}_y, \forall y \in \mathcal{C}$  by the method of moments (7).  
   **end if**  
    $\theta_{k+1} = \theta_k - \gamma'_\theta \partial_{\theta} \frac{1}{N'_B} \sum_{i \in \mathcal{B}} \hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta_k)$   
**end for**  
**Output:**  $\theta_{N'}$

in both its original implementation (PI) and its debaised version for MCAR labels (**DePI**) (Schmutz et al., 2022), using both the MNIST dataset and two datasets of MedMNIST (Yang et al., 2021). Furthermore, we compare our debaised version of Fixmatch (Sohn et al., 2020), designed to handle informative labels, with its original counterpart (**Fix**) and its debaised version for MCAR labels (**DeFix**) (Schmutz et al., 2022) on the CIFAR-10 dataset (Krizhevsky et al., 2009).

To evaluate the accuracy of our proposed estimates of the missing-data mechanism, we calculate the normalized Mean Squared Error (MSE) using the formula  $\|\hat{\phi} - \phi^*\|_2 / \|\phi^*\|_2^2$ . This provides a measure of how well our estimate of the missing-data mechanism ( $\hat{\phi}$ ) approximates the true mechanism ( $\phi^*$ ). We consider four different estimators of the missing-data mechanism.

- **MLE:** the maximum likelihood estimator derived from Algorithm 1. As highlighted in Section 4.2, we use the estimation of  $\theta$  given by Algorithm 2 when assessing the model’s performance.
- **ME:** the moment estimator derived from Algorithm 2 by using a moving averaging strategy (8) for the class distribution.
- **MEg:** the moment estimator derived from Algorithm 2 by propagating the gradients through  $\theta$ .
- **CADR:** the estimator derived from Hu et al. (2022). Although the authors did not propose an estimation of the missing-data mechanism, we are able to derive it directly from their estimation of the class proportions (see (7)).

**5.1. MNIST and CIFAR-10 for toy mechanisms**

The MNIST dataset is an advantageous choice for SSL as the classes are well-separated, allowing us to verify the effectiveness of our method in simple cases. In order to randomly select the labeled and unlabeled data per class according to a specific distribution, we follow the method proposed by Hu et al. (2022). The number of labeled data (or unlabeled data) in each class  $k$  is determined by  $n_k = n_1 \gamma^{-\frac{k-1}{K-1}}, \forall k \in \mathcal{C}$ , where  $\gamma$  controls the degree of imbalance among the classes, with  $\gamma = 1$  resulting in a balanced distribution of labeled data among classes. Additionally,  $n_1$  represents the maximum (or minimum) number of labeled data among all the classes. In particular, we consider two cases (see Figure ??, Appendix C):

- S1.** when the dataset is balanced, we randomly select labeled data in each class with  $n_1 = 400$  and  $\gamma = 10$ , and the remaining data is considered as unlabeled.
- S2.** when the dataset is unbalanced, we randomly select labeled data (resp. unlabeled data) with  $n_1 = 400$  and  $\gamma = 10$  for (resp.  $\gamma = 0.1$ ).

**S1.** (resp. **S2.**) leads to a percentage of observed labels of 3% (resp. 9%). We trained a 3-layer CNN for both Algorithm 1 and 2. In terms of estimation of the missing-data mechanism, all methods have comparable and low MSE values in the balanced setting **S1.**, as reported in Appendix C (Table 4). In the unbalanced case **S2.**, the estimation of the missing-data mechanism **CADR** underestimates the observed proportions in the four rarest classes (i.e. classes 0 to 3), as seen in Figure 2, which leads to a highest MSE (see Table 1). For model estimation, while in the balanced case **S1.** the methods have comparable results, there is an improvement in both test accuracy and test loss with our methods that include mechanism estimation (**MLE**, **MEg** and **ME**), especially for the less observed classes (i.e. classes 5 to 9). Note also that in both cases the method of Hu et al. (2022) has the highest test loss, which can be explained by the fact that the objective function that they minimize is quite different as explained in Section 3.2.

 Table 1. Test accuracy and test loss on MNIST, Setting **S2.**

METHOD	LOSS	ACCURACY	MSE $\phi$
PL	0.141 ± 0.018	92.95 ± 0.55	0.594
DEPL	0.138 ± 0.015	93.18 ± 0.71	0.594
CADR	0.160 ± 0.029	89.15 ± 0.99	0.106 ± 0.012
MLE (OURS)	0.116 ± 0.021	94.29 ± 0.11	0.027 ± 0.012
MEG (OURS)	0.103 ± 0.009	94.83 ± 0.38	0.022 ± 0.004
ME (OURS)	0.111 ± 0.005	94.59 ± 0.28	0.037 ± 0.002

In the CIFAR-10 dataset and considering Setting **S2.**, we compare the original version of Fixmatch (Sohn et al., 2020)



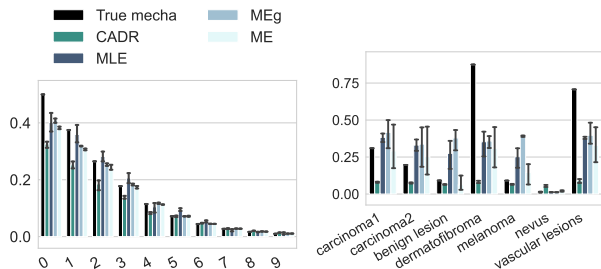


Figure 2. Estimation of the mechanism (coordinates of  $\phi$  for each class) on MNIST (Setting **S2.**) and on dermaMNIST.

with its debiased versions<sup>2</sup>. In Table 2, our method **ME** demonstrates improved performance with higher overall accuracy. Again, we observe in Figure 4 (Appendix C) that the classes with less observations (classes 5 to 10, particularly "dog" and "ship") are more accurately predicted when the missing-data mechanism is taken into account.

Table 2. Test accuracy and test loss for CIFAR10, Setting **S2.**

METHOD	LOSS	ACCURACY
PL	0.426 $\pm$ 0.017	90.91 $\pm$ 0.12
DEPL	0.536 $\pm$ 0.020	89.71 $\pm$ 0.16
CADR	0.452 $\pm$ 0.006	91.14 $\pm$ 0.30
ME (OURS)	0.321 $\pm$ 0.016	91.88 $\pm$ 0.24

## 5.2. medMNIST with pseudo-realistic mechanisms

We first consider the dermaMNIST dataset (Codella et al., 2019), which consists of 10,015 dermatoscopic images categorized into 7 different skin conditions. This dataset is unbalanced, its the most frequent class (71% of the images) being benign naevi (aka moles). To simulate a more realistic MNAR scenario, we assume that a medical doctor would like to classify the conditions equally and select 70 images per class for labeling, resulting in 7% of observed labels (see Figure 1, left panel). Note that despite this selection, the dataset remains unbalanced due to the original distribution of the classes. Our three estimators (**MLE**, **ME**, **MEg**) of the missing-data mechanism detect that the class of naevi is very little observed compared to other classes (see Figure 2), whereas **CADR** gives a mechanism where all classes are equally observed. In Table 3, only our methods determine if a lesion is a nevus or not with a high accuracy, which

<sup>2</sup>Only the results for the moment estimator using averaging strategy (**ME**) are reported, as Algorithm 1 (**MLE**) has not yet been implemented with data augmentation and **MEg** proved difficult to calibrate. It can be challenging to find the right balance between too much initialization using  $\hat{\phi} = n_{\ell}/n$  (leading to deviation from the optimal solution) and too little (leading to numerical problems). Therefore, we recommend considering the **ME** estimator when using data augmentation in practice.

can be used as a pre-processing step before the images are reviewed by a medical doctor. Note that, even if the class proportions are known, **CADR** fails to give accurate results for the nevus class (see Table 5 in Appendix C), which shows the relevance of taking into account the missing-data mechanism in this case. Finally, we use **MLE** together with the test presented in Section 4.4 to assess whether or not labels are informative. If we generate MCAR labels, the likelihood ratio test is rightfully unable to reject the MCAR hypothesis ( $p$ -value of  $0.68 \pm 0.2$  over 10 runs). But if we give it images with MNAR labels, the test rejects the MCAR hypothesis with very high confidence ( $p$ -value  $< 10^{-4}$  for all 10 runs).

Table 3. Test accuracies on dermaMNIST and noduleMNIST3D.

METHOD	DERMAMNIST	NODULEMNIST
PL	57.72 $\pm$ 1.95	84.91
CADR	49.36 $\pm$ 1.91	80.32
MLE (OURS)	66.4 $\pm$ 0.81	85.8
MEg (OURS)	66.65 $\pm$ 1.76	82.26
ME (OURS)	65.8 $\pm$ 0.78	85.16

We now consider the noduleMNIST3D dataset (Armato III et al., 2011) on images from thoracic CT scans, which is of particular interest to simulate the MNAR labels. We have access to the subtlety score  $s$ , which describes from 1 (extremely subtle) to 5 (obvious) the difficulty of nodule detection. According to these scores, we simulate the missing-data mechanism using  $p(r|y) = \sum_{s=1}^5 p(r|s)p(s|y)$ , with  $p(s|y)$  computed on the data. The only quantities to choose are the probability of being observed given the subtlety score, we fix a low probability when the detection was difficult ( $p(r|s \in \{1, 2, 3\}) = 0.1$ ) and a high probability when the detection was easy ( $p(r|s \in \{4, 5\}) = 0.9$ ). At the end, the class of benign nodules has a missing-data proportion (43%) higher than the class of malignant nodules (8%). On the contrary to the missing-data setting chosen for dermaMNIST, the more observed class is also the less frequent (Figure 6). **MLE** performs better in terms of accuracy (Table 3) and all the methods designed for informative labels has a highest specificity than the classical one **PI** for the class of malignant nodules (Table 6 in Appendix C).

## 6. Conclusion

For future works, we would be eager to (i) provide a more realistic theoretical grounding without freezing  $\theta$  and (ii) propose another statistical test for the moment estimator, for example using bootstrap strategies. Note that the latter perspective is quite challenging, because of the sample bias in the informative data.

## References

- Ahfock, D. and McLachlan, G. J. On missing label patterns in semi-supervised learning. *arXiv preprint arXiv:1904.02883*, 2019.
- Armato III, S. G., McLennan, G., Bidaut, L., McNitt-Gray, M. F., Meyer, C. R., Reeves, A. P., Zhao, B., Aberle, D. R., Henschke, C. I., Hoffman, E. A., et al. The lung image database consortium (lidc) and image database resource initiative (idri): a completed reference database of lung nodules on ct scans. *Medical physics*, 38(2):915–931, 2011.
- Baker, S. G. and Laird, N. M. Regression analysis for categorical variables with outcome subject to nonignorable nonresponse. *Journal of the American Statistical Association*, 83(401):62–69, 1988.
- Berthelot, D., Carlini, N., Cubuk, E. D., Kurakin, A., Sohn, K., Zhang, H., and Raffel, C. Remixmatch: Semi-supervised learning with distribution alignment and augmentation anchoring. *International Conference on Learning Representations*, 2019a.
- Cao, K., Brbic, M., and Leskovec, J. Open-world semi-supervised learning. *arXiv preprint arXiv:2102.03526*, 2021.
- Chapelle, O., Scholkopf, B., and Zien, A. Semi-supervised learning. *IEEE Transactions on Neural Networks*, 20(3):542–542, 2009.
- Chen, Y., Zhu, X., Li, W., and Gong, S. Semi-supervised learning under class distribution mismatch. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3569–3576, 2020.
- Codella, N., Rotemberg, V., Tschandl, P., Celebi, M. E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368*, 2019.
- d’Haultfoeuille, X. A new instrumental method for dealing with endogenous selection. *Journal of Econometrics*, 154(1):1–15, 2010.
- Grandvalet, Y. and Bengio, Y. Semi-supervised learning by entropy minimization. *Advances in neural information processing systems*, 17, 2004.
- Guo, L.-Z., Zhang, Z.-Y., Jiang, Y., Li, Y.-F., and Zhou, Z.-H. Safe deep semi-supervised learning for unseen-class unlabeled data. In *International Conference on Machine Learning*, pp. 3897–3906. PMLR, 2020.
- Hu, X., Niu, Y., Miao, C., Hua, X.-S., and Zhang, H. On non-random missing labels in semi-supervised learning. In *International Conference on Learning Representations*, 2022.
- Huang, Z., Xue, C., Han, B., Yang, J., and Gong, C. Universal semi-supervised learning. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- Ibrahim, J. G. and Lipsitz, S. R. Parameter estimation from incomplete data in binomial regression when the missing data mechanism is nonignorable. *Biometrics*, pp. 1071–1078, 1996.
- Ibrahim, J. G., Chen, M.-H., and Lipsitz, S. R. Missing responses in generalised linear mixed models when the missing data mechanism is nonignorable. *Biometrika*, 88(2):551–564, 2001.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Le Morvan, M., Prost, N., Josse, J., Scornet, E., and Varoquaux, G. Linear predictor on linearly-generated data with missing values: non consistency and solutions. In *International Conference on Artificial Intelligence and Statistics*, pp. 3165–3174. PMLR, 2020.
- LeCun, Y. and Cortes, C. MNIST handwritten digit database. 2010. URL <http://yann.lecun.com/exdb/mnist/>.
- Little, R. J. and Rubin, D. B. *Statistical analysis with missing data*, volume 793. John Wiley & Sons, 2019.
- Liu, T. and Goldberg, Y. Kernel machines with missing responses. *Electronic Journal of Statistics*, 14(2):3766–3820, 2020.
- Miao, W. and Tchetgen Tchetgen, E. J. On varieties of doubly robust estimators under missingness not at random with a shadow variable. *Biometrika*, 103(2):475–482, 2016.
- Miao, W., Ding, P., and Geng, Z. Identifiability of normal and normal mixture models with nonignorable missing data. *Journal of the American Statistical Association*, 111(516):1673–1683, 2016.
- Miao, W., Liu, L., Tchetgen, E. T., and Geng, Z. Identification, doubly robust estimation, and semiparametric efficiency theory of nonignorable missing data with a shadow variable. *arXiv preprint arXiv:1509.02556*, 2019.
- Mohan, K. On handling self-masking and other hard missing data problems. In *AAAI Symposium 2018*, 2018.

- Molenberghs, G., Beunckens, C., Sotito, C., and Kenward, M. G. Every missingness not at random model has a missingness at random counterpart with equal fit. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(2):371–388, 2008.
- Morikawa, K., Kim, J. K., and Kano, Y. Semiparametric maximum likelihood estimation with data missing not at random. *Canadian Journal of Statistics*, 45(4):393–409, 2017.
- Murphy, S. A. and Van der Vaart, A. W. On profile likelihood. *Journal of the American Statistical Association*, 95(450):449–465, 2000.
- Neuberg, L. G. Causality: models, reasoning, and inference, by judea pearl, cambridge university press, 2000. *Econometric Theory*, 19(4):675–685, 2003.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. Automatic differentiation in pytorch. 2017.
- Rabe-Hesketh, S. and Skrondal, A. Ignoring non-ignorable missingness. *Psychometrika*, pp. 1–20, 2022.
- Rizve, M. N., Duarte, K., Rawat, Y. S., and Shah, M. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. *arXiv preprint arXiv:2101.06329*, 2021.
- Rubin, D. B. Inference and missing data. *Biometrika*, 63(3): 581–592, 1976.
- Schmutz, H., Humbert, O., and Mattei, P.-A. Don’t fear the unlabelled: safe deep semi-supervised learning via simple debiasing. *arXiv preprint arXiv:2203.07512*, 2022.
- Shao, J. and Wang, L. Semiparametric inverse propensity weighting for nonignorable missing data. *Biometrika*, 103(1):175–187, 2016.
- Sohn, K., Berthelot, D., Li, C.-L., Zhang, Z., Carlini, N., Cubuk, E. D., Kurakin, A., Zhang, H., and Raffel, C. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *Advances in Neural Information Processing Systems*, 33:596–608, 2020.
- Sportisse, A., Boyer, C., and Josse, J. Estimation and imputation in probabilistic principal component analysis with missing not at random data. *Advances in Neural Information Processing Systems*, 33:7067–7077, 2020.
- Tang, N. and Ju, Y. Statistical inference for nonignorable missing-data problems: a selective review. *Statistical Theory and Related Fields*, 2(2):105–133, 2018.
- Van der Vaart, A. W. *Asymptotic statistics*, volume 3. Cambridge university press, 2000.
- Van Engelen, J. E. and Hoos, H. H. A survey on semi-supervised learning. *Machine Learning*, 109(2):373–440, 2020.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Wei, C., Sohn, K., Mellina, C., Yuille, A., and Yang, F. Crest: A class-rebalancing self-training framework for imbalanced semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10857–10866, 2021.
- Xie, Q., Dai, Z., Hovy, E., Luong, T., and Le, Q. Un-supervised data augmentation for consistency training. *Advances in Neural Information Processing Systems*, 33: 6256–6268, 2020.
- Yang, J., Shi, R., Wei, D., Liu, Z., Zhao, L., Ke, B., Pfister, H., and Ni, B. Medmnist v2: A large-scale lightweight benchmark for 2d and 3d biomedical image classification. *arXiv preprint arXiv:2110.14795*, 2021.

## A. Proof of identification

**Proposition A.1** (Identification of the joint distribution). *Under Assumptions A2. (self-masked MNAR), the joint distribution  $p(y, x, r)$  is identified.*

*Proof.* This proof is a direct application of Theorem 1 of (Miao et al., 2019) (stated in a more general case). For clarity, we give it again here in our setting for clarity by proving the intermediate results.

Let us consider the odds ratio  $OR(y) = \frac{p(y|r=0)p(y=1|r=1)}{p(y|r=1)p(y=1|r=0)}$  where  $y = 1$  is used as a reference value. The goal is to determine  $p(y|r = 0, x)$  by only quantities involving observed data, and it will imply that the joint distribution  $p(y, x, r)$  is identified. Proposition 2 of (Miao et al., 2019) gives the two following equalities:

$$p(y|r = 0, x) = \frac{OR(y)p(y|r = 1, x)}{\mathbb{E}[OR(y)|r = 1, x]} \quad (11)$$

$$\mathbb{E}[\tilde{OR}(y)|r = 1, x] = \frac{p(x|r = 0)}{p(x|r = 1)} \text{ with } \tilde{OR}(y) = \frac{OR(y)}{\mathbb{E}[OR(y)|r=1]} \quad (12)$$

The first equation (11) indicates that the identification of the odds ratio function involves the identification of  $p(y|r = 0, x)$ . Note that in (12), as  $p(x|r = 0)$ ,  $p(x|r = 1)$  and  $p(y|r = 1, x)$  are obtained from observed data, so is  $\mathbb{E}[\tilde{OR}(y)|r = 1, x]$ . We just have to prove that (11) has a unique solution. Let us consider  $\tilde{OR}^*(y) \neq \tilde{OR}(y)$ , we have

$$\mathbb{E}[\tilde{OR}^*(y)|r = 1, x] = \frac{p(x|r = 0)}{p(x|r = 1)},$$

which implies

$$\mathbb{E}[\tilde{OR}^*(y) - \tilde{OR}(y)|r = 1, x] = 0 \Leftrightarrow \tilde{OR}^*(y) = \tilde{OR}(y),$$

This is obtained by using Condition 1 of (Miao et al., 2019). In our case, it amounts to assuming that for any image, there is a non-zero probability that it has any label.  $\tilde{OR}(y)$  is identified and so is  $OR(y)$ , noting that  $OR(y) = \frac{\tilde{OR}(y)}{\tilde{OR}(y=0)}$ . □

## B. Proofs of the theoretical results of Section 4

### B.1. Moment estimator

The moment estimator has the following form, for a fixed  $\theta \in \Theta$ :

$$\hat{\phi}_y^M = \frac{\sum_{i=1}^n \mathbb{1}_{\{r_i=1, y_i=y\}}}{n} \frac{1}{\hat{p}(y, \theta)}, \forall y \in \mathcal{C}$$

with

$$\hat{p}(y, \theta) := \hat{p}_\theta(y) = \frac{1}{n} \sum_{i=1}^n p(y_i|x_i; \theta)$$

In this section, we prove the consistency of this estimator.

**Proposition B.1** (Consistency of  $\hat{\phi}^M$ ). *The moment estimator defined by (7) is consistent for a fixed  $\theta \in \Theta$ .*

*Proof.* We have by the law of large numbers,

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{r_i=1, y_i=y\}} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mathbb{E}[\mathbb{1}_{\{r_i=1, y_i=y\}}] = p(r = 1, y) \quad (13)$$

$$\frac{1}{n} \sum_{i=1}^n p(y_i|x_i; \theta) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mathbb{E}[p(y|x)] = p(y) \quad (14)$$

We can apply the continuous mapping theorem to  $f(x) = 1/x, x \in ]0, 1]$  (assuming that the probability of each class is greater than 0) to get

$$\frac{1}{\frac{1}{n} \sum_{i=1}^n p(y_i|x_i)} \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \frac{1}{p(y)} \quad (15)$$

In (13) and (15), we have the convergence in probability to a constant, which implies by Slutsky's theorem that the product converges in probability to the product of the constants (Slutsky's theorem gives the property in distribution but convergence in probability and law is equivalent if the limit is a constant), which gives:

$$\forall y \in \mathcal{C}, \hat{\phi}_y^M \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \phi_y = \frac{p(r=1, y)}{p(y)}$$

□

## B.2. Maximum likelihood estimator

In this section, we study the negative observed log-likelihood  $\phi \mapsto \ell(\theta, \phi; x, y \odot r, r)$  and derive the consistency and the asymptotic normality of the maximum likelihood estimator  $\hat{\phi}^L$  given in (9). All the results are obtained for a fixed  $\theta \in \Theta$ . Let us recall:

$$\ell(\theta, \phi; x, y \odot r, r) = - \sum_{i=1}^n r_i \log p(y_i|x_i; \theta) \phi_{y_i} - \sum_{i=1}^n (1-r_i) \log \sum_{\tilde{y} \in \mathcal{C}} p(\tilde{y}|x_i; \theta) (1-\phi_{\tilde{y}}) \quad (16)$$

For simplicity, as  $\theta$  is fixed in this section, let us define the function  $\ell_\theta : \mathbb{R}^{\mathcal{C}} \mapsto \mathbb{R}$  such that  $\ell_\theta(\phi) = \ell(\theta, \phi, x, y \odot r, r)$ .

### B.2.1. CONVEXITY

**Proposition B.2** (Convexity of  $\ell_\theta(\phi; x, y \odot r, r)$ ). *For a fixed  $\theta \in \Theta$ , the negative observed log-likelihood  $\phi \mapsto \ell(\theta, \phi; x, y \odot r, r)$  is convex.*

*Proof.* Let us remark that:

$$\begin{aligned} \ell_\theta(\phi) &= - \sum_{i=1}^n r_i \log p(y_i|x_i) \phi_{y_i} - \sum_{i=1}^n (1-r_i) \log \sum_{\tilde{y} \in \mathcal{C}} p(\tilde{y}|x_i) (1-\phi_{\tilde{y}}) \\ &= - \sum_{i=1}^n r_i \log p(y_i|x_i) \phi_{y_i} - \sum_{i=1}^n (1-r_i) \log(-\phi_k B_{i,k} + C_{i,k}), \end{aligned}$$

where  $B_{i,k} = p(y_i = k|x_i; \theta)$  and  $C_{i,k} = B_{i,k} + \sum_{\tilde{y} \in \mathcal{C} \setminus \{k\}} p(\tilde{y}|x_i; \theta) (1-\phi_{\tilde{y}})$ . Let us compute the Hessian  $H$  of  $F$ , such that

$$H(\phi) := \left( \frac{\partial^2 \ell_\theta}{\partial \phi_\ell \partial \phi_k}(\phi) \right)_{k, \ell \in \mathcal{C}}$$

We have:

$$\forall k \in \mathcal{C}, \frac{\partial \ell_\theta}{\partial \phi_k}(\phi) = - \left( \sum_{i=1}^n r_i \frac{1}{\phi_k} \mathbb{1}_{y_i=k} - \sum_{i=1}^n (1-r_i) \frac{B_{i,k}}{-\phi_k B_{i,k} + C_{i,k}} \right).$$

Thus,

$$\begin{aligned} \forall k \in \mathcal{C}, \frac{\partial^2 \ell_\theta}{\partial \phi_k \partial \phi_k}(\phi) &= \sum_{i=1}^n r_i \frac{1}{\phi_k^2} \mathbb{1}_{y_i=k} + \sum_{i=1}^n (1-r_i) \frac{(B_{i,k})^2}{(-\phi_k B_{i,k} + C_{i,k})^2} \\ \forall \ell \in \mathcal{C}, \ell \neq k, \frac{\partial^2 \ell_\theta}{\partial \phi_\ell \partial \phi_k}(\phi) &= \sum_{i=1}^n (1-r_i) \frac{B_{i,k} B_{i,\ell}}{(-\phi_k B_{i,k} + C_{i,k})^2} \end{aligned}$$

In the following, let us denote  $A_i = -\phi_k B_{i,k} + C_{i,k} = \sum_{\tilde{y} \in \mathcal{C}} p(\tilde{y}|x_i; \theta) (1-\phi_{\tilde{y}})$ . We have now to prove that the Hessian is positive-definite. As it is a symmetric matrix, we can show that  $\forall v \in \mathbb{R}^{\mathcal{C}}, v \neq \vec{0}, v^T H v > 0$ .

$$v^T H v = \sum_{k=1}^K \frac{v_k^2}{\phi_k^2} \sum_{i=1}^n r_i \mathbb{1}_{y_i=k} + \underbrace{\sum_{k=1}^K v_k^2 \sum_{i=1}^n (1-r_i) \frac{(B_{i,k})^2}{A_i^2} + 2 \sum_{1 \leq k < \ell \leq K} v_k v_\ell \sum_{i=1}^n (1-r_i) \frac{B_{i,k} B_{i,\ell}}{A_i^2}}_{=T} > 0$$

The first term is trivially greater or equal to 0. Moreover, it is never equal to 0, if at least one sample is observed ( $n_\ell > 0$ ). For the last two terms, note that:

$$T = \sum_{i=1}^n \left( \sum_{k=1}^K v_k B_{i,k} \sqrt{\frac{(1-r_i)}{A_i^2}} \right)^2 \geq 0$$

□

*Remark B.3* (Domain of definition of  $\ell_\theta(\phi; x, y \odot r, r)$ ). We look at the natural domain of the negative observed log-likelihood  $\phi \mapsto \ell(\theta, \phi; x, y \odot r, r)$ , for a fixed  $\theta \in \Theta$ . The goal is to know if we can minimize the function without constraint on  $\phi$  (if its domain of definition is included in  $[0, 1]$ )

In (16), the first term implies  $\forall y_i, p(y_i|x_i; \theta) \phi_{y_i} > 0$  i.e.  $\phi_k > 0, \forall k \in \{1, \dots, K\}$ . The second term requires  $\forall i \in \{n_\ell + 1; n\}$ ,

$$\begin{aligned} & \sum_{\tilde{y} \in \mathcal{C}} p(\tilde{y}|x_i; \theta) (1 - \phi_{\tilde{y}}) > 0 \\ \Leftrightarrow & \forall k \in \{1, \dots, K\}, (1 - \phi_k) p(\tilde{y} = k|x_i; \theta) + \sum_{\tilde{y} \in \mathcal{C} \setminus \{k\}} p(\tilde{y}|x_i; \theta) (1 - \phi_{\tilde{y}}) > 0 \\ \Leftrightarrow & \forall k \in \{1, \dots, K\}, \phi_k < 1 + \frac{1}{p(\tilde{y} = k|x_i; \theta)} \sum_{\tilde{y} \in \mathcal{C} \setminus \{k\}} p(\tilde{y}|x_i; \theta) (1 - \phi_{\tilde{y}}) \end{aligned}$$

As  $\frac{1}{p(\tilde{y}=k|x_i; \theta)} \sum_{\tilde{y} \in \mathcal{C} \setminus \{k\}} p(\tilde{y}|x_i; \theta) (1 - \phi_{\tilde{y}}) > 0$ , this last inequality does not necessarily implies  $\phi_k \leq 1$ . Therefore, a reparametrization trick or constrained optimization is essential. The domain of definition for  $\phi_k, k \in \{1, \dots, K\}$  of the negative log-likelihood  $\ell$  is then

$$D_{\phi_k} = \left] 0; 1 + \min_{i \in \{1, \dots, n\}} \frac{1}{p(\tilde{y} = k|x_i; \theta)} \sum_{\tilde{y} \in \mathcal{C} \setminus \{k\}} p(\tilde{y}|x_i; \theta) (1 - \phi_{\tilde{y}}) \right[.$$

### B.2.2. CONSISTENCY AND ASYMPTOTIC NORMALITY

The consistency and asymptotic normality of the maximum likelihood estimator is obtained by applying Theorem 5.7 and Theorem 5.23 of (Van der Vaart, 2000). Let us assume the following:

**A3.** We can interchange differentiation with respect to  $\phi$  and integration over  $(x, y, r)$ .

**A4.**  $\forall k \in \mathcal{C}$ , there exists a compact interval  $U_k$  such that  $\phi_k \in U_k \subset ]0, 1[$ .

To get the results, we need first to show the identifiability of the statistical model  $\mathcal{P}_\phi = \{p(r|y \odot r; \phi) : \phi \in \Phi\}$  and the nonsingularity of the Fisher information  $I_{\phi^*}$ , with  $\phi^*$  the oracle point (i.e.  $\phi^* = \operatorname{argmin}_{\phi \in \Phi} \ell_\theta(\phi)$ ).

**Lemma B.4** (Identifiability of  $\mathcal{P}_\phi$ ). *The model  $\mathcal{P}_\phi$  is identifiable.*

*Proof.* Let us consider  $\phi, \phi' \in \Phi$  such that  $p(r|y \odot r; \phi) = p(r|y \odot r; \phi')$ .

$$\begin{aligned} p(r|y \odot r; \phi) &= p(r|y \odot r; \phi'), \forall y, r \\ \Leftrightarrow r \phi_y + \sum_{\tilde{y}} (1-r)(1 - \phi_{\tilde{y}}) &= r \phi'_y + \sum_{\tilde{y}} (1-r)(1 - \phi'_{\tilde{y}}), \forall y, r \\ \Leftrightarrow r(\phi_y - \phi'_y) + \sum_{\tilde{y}} (1-r)(\phi'_{\tilde{y}} - \phi_{\tilde{y}}) &= 0, \forall y, r \end{aligned}$$

The case  $r = 1$  leads to  $\phi_y = \phi'_y, \forall y$ . □

**Lemma B.5** (Nonsingularity of the Fisher information  $I_{\phi^*}$ ). *Under Assumption A3., the Fisher information at the oracle point  $\phi^*$  is nonsingular.*

*Proof.* Assumption A3. implies that  $\forall \ell, k \in \mathcal{C}, (I_{\phi^*})_{(\ell,k)} = -\mathbb{E}_{x,y,r} \left[ \frac{\partial^2 \log \ell_{\theta}(\phi^*)}{\partial \phi_{\ell} \partial \phi_k} \right] = -\frac{1}{n} \mathbb{E}_{x,y,r} [(H(\phi))_{(\ell,k)}]$ . We can simply use the strict convexity of the function  $\phi \mapsto \ell_{\theta}(\phi)$  proved in Proposition B.2.  $\square$

**Proposition B.6** (Consistency of  $\hat{\phi}^L$ ). *Under Assumption A4., the maximum likelihood estimator  $\hat{\phi}^L$  defined in (9) is consistent, for a fixed  $\theta \in \Theta$ .*

*Proof.* As said in Section 5.5 of Van der Vaart (2000) for the application to the maximum likelihood estimators, the method consists of applying Theorem 5.7 of Van der Vaart (2000) by noting that the MLE is a M-estimator:  $\hat{\phi}^L \in \operatorname{argmin}_{\phi \in \Phi} M_n(\phi)$ , with  $M_n(\phi) = \frac{1}{n} \sum_{i=1}^n \log \frac{p(x_i, y_i \odot r_i, r_i; \phi^*)}{p(x_i, y_i \odot r_i, r_i; \phi)}$  and  $M(\phi) = \mathbb{E} \left[ \log \frac{p(x, y \odot r, r; \phi^*)}{p(x, y \odot r, r; \phi)} \right]$ .

- We have the identifiability of  $\mathcal{P}_{\phi}$  by Lemma B.4. We also have the strong identifiability, which is equivalent to the identifiability, since the restricted set of  $\Phi, \otimes_{k \in \mathcal{C}} U_k$ , is compact (see Assumption A4.).
- We have to show that the uniform weak law of large numbers hold for the function  $\phi \mapsto \log \frac{p(x_i, y_i \odot r_i, r_i; \phi^*)}{p(x_i, y_i \odot r_i, r_i; \phi)}$ , i.e.

$$\sup_{\phi \in \Phi} |M_n(\phi) - M(\phi)| \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} 0.$$

Following Theorem 4.2 of Wainwright (2019), a sufficient condition is that there exists  $M > 0$  such that,  $\left| \log \frac{p(x, y \odot r, r; \phi^*)}{p(x, y \odot r, r; \phi)} \right| < M, \forall \phi, x, y, r$ , i.e. the uniform boundedness. We will prove that there exists  $M_1, M_2 > 0$  such that  $M_1 \leq \frac{p(x, y \odot r, r; \phi^*)}{p(x, y \odot r, r; \phi)} \leq M_2$ .

We have

$$\frac{p(x, y \odot r, r; \phi^*)}{p(x, y \odot r, r; \phi)} = \frac{r\phi_y^* + (1-r) \sum_{\bar{y}} (1 - \phi_{\bar{y}}^*)}{r\phi_y + (1-r) \sum_{\bar{y}} (1 - \phi_{\bar{y}})}$$

Using Assumption A4., we get  $M_1 = \min \left( \frac{\min_k a_k}{\max_k b_k}, \frac{1 - \max_k b_k}{1 - \min_k a_k} \right)$  and  $M_2 = \max \left( \frac{1 - \min_k a_k}{1 - \max_k b_k}, \frac{\max_k b_k}{\min_k a_k} \right)$

- By identifiability of the statistical model  $\mathcal{P}_{\phi}$ , we have:  $M(\phi) = M(\phi^*)$  implies  $\phi = \phi^*$ . Therefore  $\forall \phi \neq \phi^*, M(\phi) > M(\phi^*)$  and the function  $\phi \mapsto M(\phi)$  admits a strict minimum in  $\phi_0$ .  $\square$

**Proposition B.7** (Asymptotic normality of  $\hat{\phi}^L$ ). *Under Assumption A3. and A4., the maximum likelihood estimator  $\hat{\phi}^L$  defined in (9) is asymptotically normal, for a fixed  $\theta \in \Theta$ .*

*Proof.* The proof directly follows from Theorem 5.39 of Van der Vaart (2000) (Corollary of Theorem 5.23 for the maximum likelihood estimators). To apply the theorem, we check the following conditions:

- The statistical model  $\mathcal{P}_{\phi} = \{p(r|y \odot r; \phi) : \phi \in \Phi\}$  is differentiable in quadratic mean at  $\phi^*$ , because  $p(r|y \odot r; \phi) = r\phi_y + \sum_{\bar{y}} (1-r)(1 - \phi_{\bar{y}})$  is trivially twice differentiable.
- The score function  $S(\phi; r, y \odot r)$  is uniformly bounded in  $y, r$  and  $\phi$ , ranging over a compact and continuous in  $\phi$  for all  $y, r$ , i.e. we want to show that there exists a real number  $M$ ,

$$\|S(\phi; r, y)\|_1 = \sum_{k=1}^K \left| \frac{\partial \log p(r|y \odot r; \phi)}{\partial \phi_k} \right| \leq M, \forall \phi, y, r$$

The score function for its coordinate  $k$  is:  $S(\phi_k; r, y \odot r) = \frac{\partial \log p(r|y \odot r; \phi)}{\partial \phi_k} = \frac{r\mathbb{1}_{y=k} - (1-r)}{r\phi_y + (1-r) \sum_{\bar{y}} (1 - \phi_{\bar{y}})}$ . If  $r = 1$ , this amounts to bound  $1/\phi_k$  and if  $r = 0$ , this amounts to bound  $\frac{1}{\sum_{\bar{y}} (1 - \phi_{\bar{y}})} \leq \frac{1}{K(1 - \max_k \phi_k)}$ . Using Assumption A4. is sufficient to get the bound. Let us denote  $U_k = [a_k, b_k], \forall k \in \mathcal{C}$ , one has  $a_k \leq \phi_k \leq b_k \forall k \in \mathcal{C}$  and we can choose for the bound  $M = K \max \left( \frac{1}{\min_k a_k}, \frac{1}{K(1 - \max_k b_k)} \right)$ .

- By Lemma B.5, the Fisher information  $I_{\phi^*}$  is nonsingular.

- By Proposition B.6, the estimator  $\hat{\phi}^L$  is consistent.

□

### B.3. Consistency of the SSL risk

**Proposition B.8** (Consistency of the risk). *If  $\hat{\phi}$  is a consistent estimator of  $\phi$  and if the mechanism is well specified, the risks  $\hat{\mathcal{R}}_{\hat{\phi}}$  and  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta)$  are consistent estimators of the theoretical risk  $\mathcal{R}(\theta) = \mathbb{E}[\ell_{\ell}(\theta; x, y)]$ .*

*Proof.* We prove the results for the IPW estimator  $\hat{\mathcal{R}}_{\hat{\phi}}$  (the proof is similar for  $\hat{\mathcal{R}}_{\hat{\phi}}^{\text{SSL}}(\theta)$ ).

It is a simple application of the law of large numbers, using the unbiasedness of the estimator (by Proposition 3.1).

We have:

$$\frac{1}{n} \sum_{i=1}^n r_i \frac{\ell_{\ell}(\theta; x_i, y_i)}{\hat{\phi}_{y_i}} = \frac{1}{n} \sum_{i=1}^n r_i \frac{\ell_{\ell}(\theta; x_i, y_i)}{\phi_{y_i}} + o_{\mathbb{P}}(1) \xrightarrow[n \rightarrow +\infty]{\mathbb{P}} \mathbb{E} \left[ r \frac{\ell_{\ell}(\theta; x, y)}{\phi_y} \right] = \mathbb{E} [\ell_{\ell}(\theta; x, y)]$$

□

## C. Additional numerical experiments

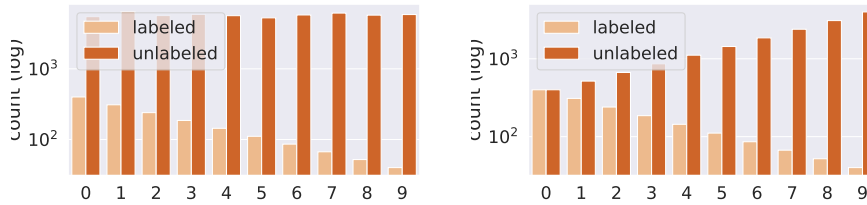


Figure 3. MNAR labels on MNIST (Settings S1. and S2.)

Table 4. Test accuracy and test loss on MNIST, Setting S1..

METHOD	LOSS	ACCURACY	MSE $\hat{\phi}$
PL	0.259 ± 0.034	95.48 ± 0.16	0.318
DEPL	0.237 ± 0.045	95.69 ± 0.06	0.318
CADR	0.272 ± 0.046	95.40 ± 0.33	0.014 ± 0.004
MLE (OURS)	0.249 ± 0.050	95.59 ± 0.40	0.031 ± 0.009
MEG (OURS)	0.240 ± 0.029	95.69 ± 0.80	0.004 ± 0.001
ME (OURS)	0.240 ± 0.027	95.34 ± 0.26	0.013 ± 0.001



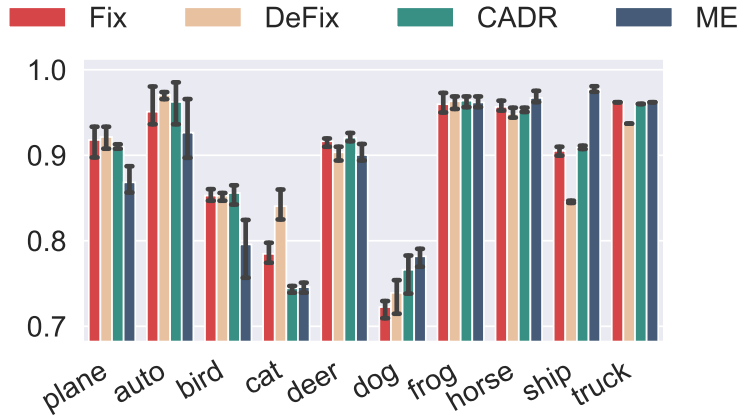


Figure 4. Accuracy per class on CIFAR10, Setting S2.

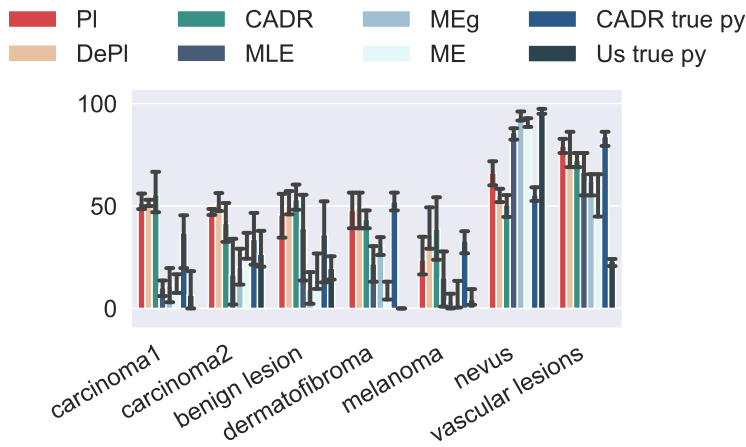


Figure 5. Accuracy per class on dermaMNIST.

Table 5. Accuracy and Loss in the test set on dermaMNIST.

METHOD	LOSS	ACCURACY	ACCURACY NAEVI	MSE $\phi_{NEVI}$
PL	$1.34 \pm 0.16$	$57.72 \pm 1.95$	$66.14 \pm 5.86$	0.80
CADR	$1.42 \pm 0.060$	$49.36 \pm 1.91$	$50.41 \pm 5.38$	$0.77 \pm 0.02$
MLE (OURS)	$0.993 \pm 0.020$	$66.4 \pm 0.81$	$91.16 \pm 2.26$	$0.34 \pm 0.03$
MEg (OURS)	$1.19 \pm 0.148$	$66.65 \pm 1.76$	$93.54 \pm 2.30$	$0.42 \pm 0.08$
ME (OURS)	$1.24 \pm 0.087$	$65.8 \pm 0.78$	$85.91 \pm 3.05$	$0.38 \pm 0.15$
CADR ( $p(y)$ KNOWN)	$1.57 \pm 0.12$	$49.44 \pm 3.27$	$55.40 \pm 3.39$	
ALGORITHM 2 ( $p(y)$ KNOWN)	$0.943 \pm 0.029$	$68.83 \pm 0.26$	$96.12 \pm 1.20$	

Table 6. Accuracy and Loss in the test set on noduleMNIST3D

METHOD	LOSS	SPECIFICITY (BEGNIN)	SPECIFICY (MALIGN)	MSE $\phi$
PL	0.389	91.06	54.69	0.0627
CADR	0.48	80.08	81.25	0.0143
MLE (OURS)	0.359	87.80	71.88	0.0001
MEG (OURS)	0.353	83.74	76.56	0.0199
ME (OURS)	0.355	86.99	78.13	0.0002

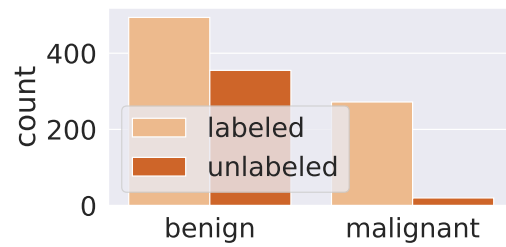


Figure 6. MNAR labels on noduleMNIST (see Section 5.2 for an explanation of the missing-data scenario)