

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

UNIVERSITÉ
CÔTE D'AZUR



inria
informatics mathematics

3iA Côte d'Azur
Institut interdisciplinaire
d'intelligence artificielle

Hugo Schmutz

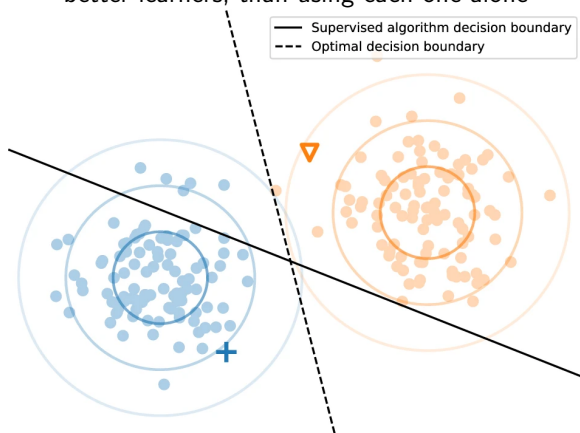
Supervisors: Olivier Humbert & Pierre-Alexandre Mattei
Inria

TIRO-MATOS, CEA UMR 4320 Equipe Maasai
Laboratoire LJAD, UMR CNRS 7351
Université Côte d'Azur
3iA Côte d'Azur

✉ hugo.schmutz@inria.fr - [@HugoSchmutz2](https://twitter.com/HugoSchmutz2)

Deep semi-supervised learning ? What for ?

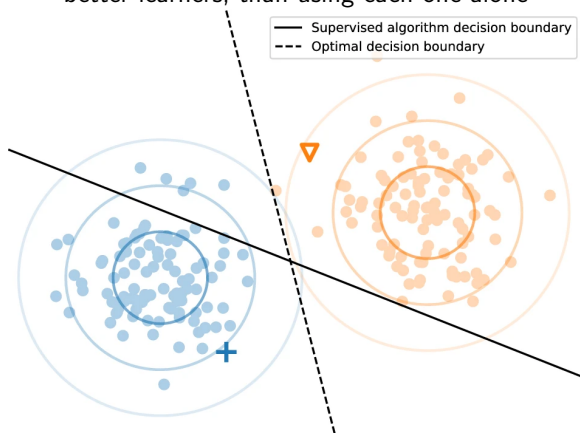
Goal: Using both labelled and unlabelled data to build better learners, than using each one alone



V.Engelen & Hoos [Machine Learning, 2020]

Deep semi-supervised learning ? What for ?

Goal: Using both labelled and unlabelled data to build better learners, than using each one alone



V.Engelen & Hoos [Machine Learning, 2020]

Why bother ?

- unlabelled data are cheap
- labelled data can be hard to get



Learning theory relies on the unbiased estimator of the risk

We have two spaces of objects X and Y and would like to learn a function $f_\theta : X \rightarrow Y$.
 L is a loss function which measures how different the prediction $f_\theta(x)$ is from the true outcome y .
We define the **risk**:

$$\mathcal{R}(\theta) = \mathbb{E}[L(\theta; x_i, y_i)]$$

The ultimate goal of a learning algorithm is to find θ^* among a fixed class of functions:

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{R}(\theta)$$

Learning theory relies on the unbiased estimator of the risk

We have two spaces of objects X and Y and would like to learn a function $f_\theta : X \rightarrow Y$.
 L is a loss function which measures how different the prediction $f_\theta(x)$ is from the true outcome y .
We define the **risk**:

$$\mathcal{R}(\theta) = \mathbb{E}[L(\theta; x_i, y_i)]$$

The ultimate goal of a learning algorithm is to find θ^* among a fixed class of functions:

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{R}(\theta)$$

In practice: $P(x, y)$ is unknown. We have access to a training set: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

Learning theory relies on the unbiased estimator of the risk

We have two spaces of objects X and Y and would like to learn a function $f_\theta : X \rightarrow Y$.
 L is a loss function which measures how different the prediction $f_\theta(x)$ is from the true outcome y .
We define the **risk**:

$$\mathcal{R}(\theta) = \mathbb{E}[L(\theta; x_i, y_i)]$$

The ultimate goal of a learning algorithm is to find θ^* among a fixed class of functions:

$$\theta^* = \arg \min_{\theta \in \Theta} \mathcal{R}(\theta)$$

In practice: $P(x, y)$ is unknown. We have access to a training set: $\{(x_1, y_1), \dots, (x_n, y_n)\}$

Empirical risk: Monte Carlo estimator of the risk

$$\hat{\mathcal{R}}(\theta) = \frac{1}{n} \sum_{i=1}^n L(\theta; x, y)$$

→ **Unbiased** estimator \Rightarrow basic property that is needed for the development of traditional learning theory and asymptotic statistics.

The complete case, the simplest: get rid of unlabelled data

- labelled data: $\{(x_1, y_1), \dots, (x_{n_l}, y_{n_l})\}$
- unlabelled data: $\{x_{n_l+1}, \dots, x_{n_l+n_u}\}$
- $n_l + n_u = n$

Semi-supervised is a missing data problem. We consider the case missing completely at random (**MCAR**) i.e. y being missing is independent x .

The complete case, the simplest: get rid of unlabelled data

- labelled data: $\{(x_1, y_1), \dots, (x_{n_l}, y_{n_l})\}$
- unlabelled data: $\{x_{n_l+1}, \dots, x_{n_l+n_u}\}$
- $n_l + n_u = n$

Semi-supervised is a missing data problem. We consider the case missing completely at random (**MCAR**) i.e. y being missing is independent x .

Complete case. Get rid of unlabelled data:

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

→ Under **MCAR**, also unbiased \Rightarrow basic property that is needed for the development of traditional learning theory and asymptotic statistics.

Using unlabelled data leads to biased estimator of the risk

- labelled data: $\{(x_1, y_1), \dots, (x_{n_l}, y_{n_l})\}$
- unlabelled data: $\{x_{n_l+1}, \dots, x_{n_l+n_u}\}$
- $n_l + n_u = n$

Semi-supervised is a missing data problem. We consider the case missing completely at random (**MCAR**) i.e. y being missing is independent x .

Including unlabelled data in the risk estimator:

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

Using unlabelled data leads to biased estimator of the risk

- labelled data: $\{(x_1, y_1), \dots, (x_{n_l}, y_{n_l})\}$
- unlabelled data: $\{x_{n_l+1}, \dots, x_{n_l+n_u}\}$
- $n_l + n_u = n$

Semi-supervised is a missing data problem. We consider the case missing completely at random (**MCAR**) i.e. y being missing is independent x .

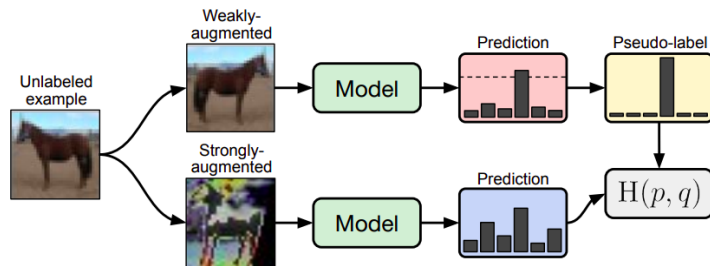
Including unlabelled data in the risk estimator:

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

Examples:

- Entropy minimization: $H(\theta; x_i) = -\sum f_{\theta}(x)_k \log(p_{\theta}(x)_k)$ (Grandvalet and Bengio 2005)
- Consistency based: $H(\theta; x_i) = \text{Div}(p_{\theta}(x), p_{\theta}(\text{pert}(x)))$ (Sohn et al. 2020, Xie et al. 2020, Miyato et al. 2018, Laine and Aila 2017, ...)
- Pseudo-label (PL): $H(\theta; x_i) = -\log(\max_u p_{\theta}(u|x)) \mathbb{1}[\max_y p_{\hat{\theta}}(y|x) > \tau]$ (Scudder 1965, Lee 2013, Rizve 2021)

Fixmatch (Sohn et al., NeurIPS [2020])



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

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Problems:

- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Problems:

- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])
- Few theoretical guarantees using strong distributional assumptions (Mey & Loog [2019])

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Problems:

- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])
- Few theoretical guarantees using strong distributional assumptions (Mey & Loog [2019])
- No asymptotic consistency: may fail event with an infinite number of labelled datapoints

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Problems:

- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])
- Few theoretical guarantees using strong distributional assumptions (Mey & Loog [2019])
- No asymptotic consistency: may fail event with an infinite number of labelled datapoints
- Choice of H can be confusing (Corduneanu & Jaakkola [UAI, 003], Krause et al. [NeurIPS, 2010])

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

Good performance on a various of (deep) learning tasks.

Problems:

- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])
- Few theoretical guarantees using strong distributional assumptions (Mey & Loog [2019])
- No asymptotic consistency: may fail event with an infinite number of labelled datapoints
- Choice of H can be confusing (Corduneanu & Jaakkola [UAI, 003], Krause et al. [NeurIPS, 2010])
- An additional hyperparameter λ and no realistic validation (Oliver et al. [NeurIPS, 2018])

Using unlabelled data leads to biased estimator of the risk

Including unlabelled data in the risk estimator:

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

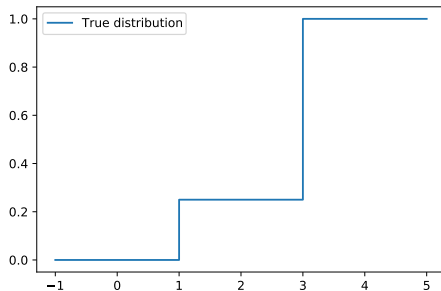
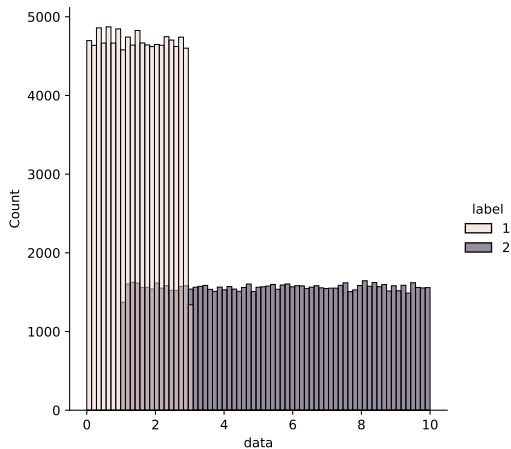
Good performance on a various of (deep) learning tasks.

Problems:

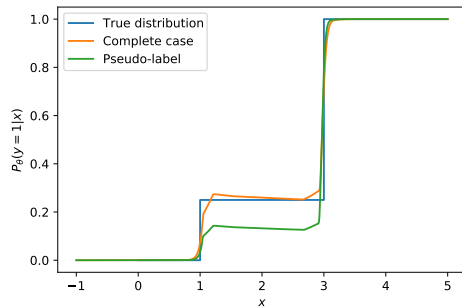
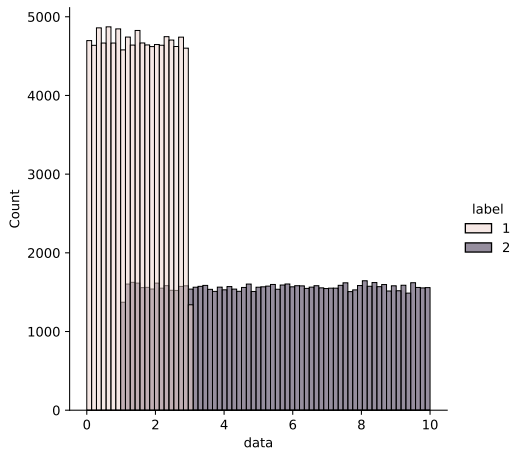
- Potential degradation reported in previous works (Schölkopf et al. [ICML, 2012], V.Engelen & Hoos, [Machine Learning, 2020], Zhu et al. [ICLR, 2022])
- Few theoretical guarantees using strong distributional assumptions (Mey & Loog [2019])
- No asymptotic consistency: may fail event with an infinite number of labelled datapoints
- Choice of H can be confusing (Corduneanu & Jaakkola [UAI, 003], Krause et al. [NeurIPS, 2010])
- An additional hyperparameter λ and no realistic validation (Oliver et al. [NeurIPS, 2018])

Safety: A SSL algorithm is safe if it has theoretical guarantees that are similar to or stronger than the complete case baseline.

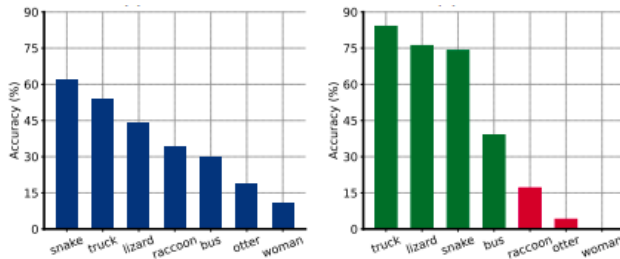
FEAR THE UNLABELLED !



FEAR THE UNLABELLED !



The rich get richer ! (Zhu et al. [ICLR, 2022])



Top-1 accuracy of 7 randomly selected categories with different training methods on CIFAR-100.

(Left) Complete case. (Right) FixMatch

FixMatch largely increases the bias of poor-behaved categories. (Chen et al., [arXiv:2202.07136])

DeSSL: Debiased version of SSL

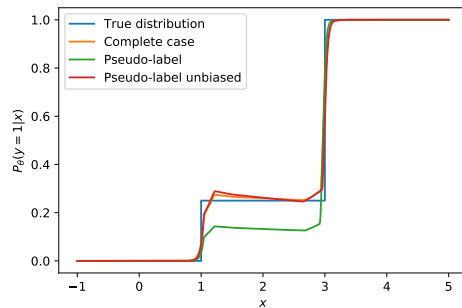
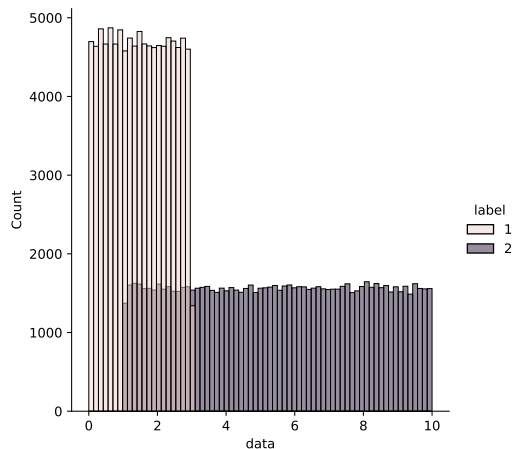
$$\hat{\mathcal{R}}_{DeSSL}(\theta) = \frac{1}{n_l} \sum_{i=1}^{n_l} L(\theta; x_i, y_i) + \frac{\lambda}{n_u} \sum_{i=1}^{n_u} H(\theta; x_i) - \frac{\lambda}{n_l} \sum_{i=1}^{n_l} H(\theta; x_i)$$

DeSSL: Debiased version of SSL

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

- Under the MCAR assumption, estimator is unbiased estimator of the true risk.
- Close relationship with control variates (Owen [2013]).
- Other motivations:
 - Penalising the confidence of a model on labelled data (Pereyra et al., [2017])
 - Maximising the plausibility (Barndorff-Nielsen [1976]).
 - The risk estimate is a Lagrangian!

Pseudo-label unbiased success under the cluster assumption



DeSSL is calibrated and benefits from generalisation error bounds

Variance:

- It exists λ such as the variance of the risk estimate is lower than the one of the complete case.
- We expect DeSSL to be a **more accurate** risk estimate than the complete case.
- Formula for λ_{opt} .
- Justification on the heuristic idea that H should be a surrogate of L .

Calibration:

- We can expect DeSSL to be as **well-calibrated** as the complete case, while SSL will generally overfit.

Generalisation error bound:

- Under classical assumptions on L and H , DeSSL benefits of generalisation error bounds derived from the **Rademacher complexity**.

DeSSL is consistent and improves the supervised baseline

Consistency:

- Under classical assumptions on L and H , DeSSL provides **asymptotically consistent** models.
- SSL may fail with an infinite number of labelled data when DeSSL will not

Asymptotic normality:

- Under classical assumptions on L and H , the parameter estimated using DeSSL is asymptotically normal.
- It exists λ such as the asymptotic variance of the estimated parameters is lower than the one of the complete case.
- **DeSSL outperforms the complete case baseline in term of parameters estimation**

DeFixmatch improves Fixmatch accuracy and calibration

Table: Test accuracy, worst class accuracy and cross-entropy of Complete Case, Fixmatch and DeFixmatch on 5 folds of CIFAR-10 and one fold of CIFAR-100.

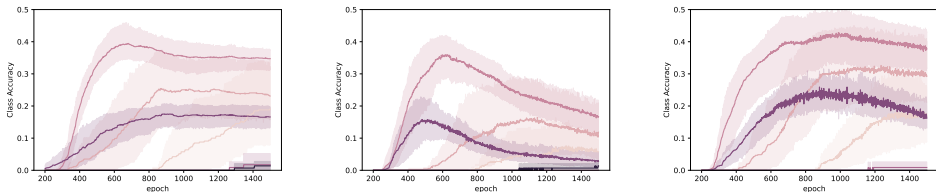
	CIFAR-10 ($n_l = 4000$)			CIFAR-100 ($n_l = 10000$)		
	Complete Case	Fixmatch	DeFixmatch	Complete Case	Fixmatch	DeFixmatch
Accuracy	87.27 ± 0.25	93.87 ± 0.13	95.44 ± 0.10	62.62	69.28	71.22
Worst class accuracy	70.08 ± 0.93	82.25 ± 2.27	87.16 ± 0.46	28.00	23.00	31.00
Cross entropy	0.60 ± 0.01	0.27 ± 0.01	0.20 ± 0.01	1.87	1.52	1.42
Brier score	0.214 ± 0.005	0.101 ± 0.003	0.076 ± 0.001	0.56	0.47	0.44

DeSSL mitigates the disparact effect of SSL

Table: Mean accuracy per class and mean benefit ratio (BR , Zhu et al. [2022]) on 5 splits.

	Complete Case	Fixmatch		DeFixmatch	
	Accuracy	Accuracy	BR	Accuracy	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

DeSSL mitigates the disparact effect of SSL



Class accuracies (without the majority class) on DermaMNIST trained with $n_l = 1000$ labelled data on five folds. (Left) CompleteCase (B-Acc: $26.88 \pm 2.26\%$); (Middle) PseudoLabel (B-Acc: $22.03 \pm 1.45\%$); (Right) DePseudoLabel (B-Acc: **$28.84 \pm 1.02\%$**), with 95% CI.

Conclusion

More results on CIFAR-100, SVHN, STL10, UCI datasets and MedMNIST datasets in the paper.

- DeSSL come with theoretical guarantees using only the MCAR assumption
- Estimator unbiased, reduction of variance, asymptotically consistent, well calibrated
- Formula for the hyperparameter
- Mitigates the disparact effect of SSL
- Performs better than the biased estimator on various datasets (see our paper)

Future directions:

- Compute of λ_{opt}
- Extend to non MCAR settings.
- Extend to segmentation
- Stratified sampling to select autmatically the number of labelled and unlabelled data per batch

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

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

$$\lambda_{opt} = \frac{n_u}{n} \frac{\text{Cov}(L(\theta; x, y), H(\theta; x))}{\mathbb{V}(H(\theta; x))} \quad (1)$$

and

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

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

Justification on the heuristic idea that H should be a surrogate of L .

$\hat{\mathcal{R}}_{DeSSL}(\theta)$ is an accurate risk estimator but $\nabla \hat{\mathcal{R}}_{DeSSL}(\theta)$ is even better

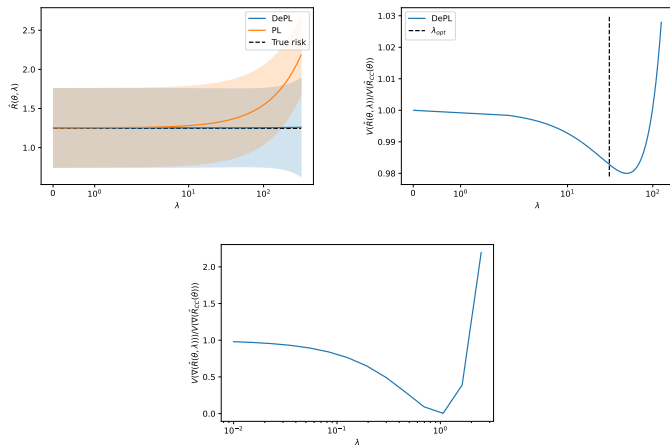


Figure: (Left) Risk estimate value for PseudoLabel (PL) and DePseudoLabel (DePL) compared to the true value of the risk. (Right) The influence of λ on the ratio $V(\hat{\mathcal{R}}_{DePL}(\theta)|r)/V(\hat{\mathcal{R}}_{CC}(\theta)|r)$. (Down) The influence of λ on the ratio $V(\nabla \hat{\mathcal{R}}_{DePL}(\theta)|r)/V(\nabla \hat{\mathcal{R}}_{CC}(\theta)|r)$.

DeSSL models are calibrated

Theorem: If the original loss is a proper scoring rule, then DeSSL is also a proper scoring rule.

- We can expect DeSSL to be as well-calibrated as the complete case.
- What if the model is non probabilistic ? Fisher consistency

Asymptotic consistency

$\hat{\theta}$ is **asymptotically consistent** with respect to n if $d(\hat{\theta}, \theta^*) \xrightarrow{P} 0$.

Theorem: Under usual regularity conditions of M-estimators for both L and H , $\hat{\theta} = \arg \min \hat{\mathcal{R}}_{DeSSL}$ is asymptotically consistent with respect to n .

SSL may fail with an infinite number of labelled data when DeSSL will not.

Theorem:

Suppose L and H are smooth functions in $\mathcal{C}^2(\Theta, \mathbb{R})$. Assume $\mathcal{R}(\theta)$ admit a second-order Taylor expansion at θ^* with a non-singular second order derivative V_{θ^*} . Under the MCAR assumption, we have that $\hat{\theta}_{DeSSL}$ is asymptotically normal and we can minimise the trace of the asymptotic covariance. Indeed, $\text{Tr}(\Sigma_{DeSSL})$ reaches its minimum at

$$\lambda_{opt} = (1 - \pi) \frac{\text{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})}{\text{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1})}, \quad (3)$$

and at λ_{opt} :

$$\text{Tr}(\Sigma_{DeSSL}) - \text{Tr}(\Sigma_{CC}) = -\frac{1 - \pi}{\pi} \frac{\text{Tr}(V_{\theta^*}^{-1} K_{\theta^*} V_{\theta^*}^{-1})^2}{\text{Tr}(V_{\theta^*}^{-1} \mathbb{E} [\nabla H(\theta^*; x) \nabla H(\theta^*; x)^T] V_{\theta^*}^{-1})} \leq 0. \quad (4)$$

DeSSL does benefit of generalisation error bounds

Theorem: We assume that labels are MCAR and that both L and H are bounded. 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 $\theta \in \Theta$,

$$\mathcal{R}(\theta) \leq \hat{\mathcal{R}}_{DeSSL}(\theta) + 2R_n + \kappa \sqrt{\frac{\log(4/\delta)}{n}}. \quad (5)$$