

Repondération Préférentielle pour l’Apprentissage Biqualité

Pierre Nodet^{*,**}, Vincent Lemaire^{*} Alexis Bondu^{*} Antoine Cornuéjols^{**}

^{*} Orange Labs, Paris & Lannion, France

^{**} AgroParisTech, Paris, France

Résumé. Cet article propose une vision originale et globale de l’Apprentissage Faiblement Supervisé, menant à la conception d’approches génériques capable de traiter tout type de faiblesses en supervision. Un nouveau cadre appelé “Données Biqualité” est introduit, qui suppose qu’un petit jeu de données fiable d’exemples correctement étiquetés est disponible, en plus d’un jeu de données non fiable comprenant un grand nombre d’exemples potentiellement corrompus. Dans ce cadre nous proposons un nouveau schéma de repondération capable de détecter les exemples non corrompus du jeu de données non fiable. Cet algorithme permet d’apprendre des classifieurs sur les deux jeux de données. Nos expériences simulant plusieurs types de bruits d’étiquetage démontrent empiriquement que l’algorithme proposé sur-performe l’état de l’art.

1 Introduction

La classification supervisée a pour objectif d’apprendre un classifieur à partir d’exemples étiquetés qui est capable de prédire la classe d’exemples non vus pendant l’apprentissage. En pratique, les algorithmes classiques de classification se heurtent aux imperfections des jeux de données réels. Ainsi l’apprentissage faiblement supervisé à récemment vu un regain en popularité avec de nombreux articles traitant de plusieurs types de “*faiblesses en supervision*”. Tous ces types de faiblesses en supervision sont traités séparément dans la littérature conduisant à des approches hautement spécialisées. En pratique il est pourtant très difficile d’identifier précisément les types de faiblesses en supervision que comporte un jeu de données. Cette article présente un nouveau cadre appelé données biqualités qui couvre un large spectre de faiblesses en supervision et permet une unification des méthodes d’apprentissage faiblement supervisé. Pour une lecture plus approfondie, se référer à une introduction sur l’apprentissage faiblement supervisé et sur les liens existants avec l’apprentissage biqualité dans (Nodet et al., 2021).

L’apprentissage sur des données biqualités a récemment été mis en lumière dans (Charikar et al., 2017; Hendrycks et al., 2018; Hataya et Nakayama, 2019), cela consiste à apprendre un classifieur à partir de deux jeux de données différents, l’un fiable et l’autre non fiable. Le but originel était d’unifier l’apprentissage semi-supervisé et robuste au bruit d’annotation par la combinaison des deux. Néanmoins ce scénario n’est pas limité à cette unification mais il peut couvrir un large spectre de faiblesses en supervision tel que démontré avec l’algorithme proposé et les résultats associés obtenus.

Le jeu de données fiable D_T (‘T’ pour « trusted ») est composé d’individus x_i associés aux étiquettes y_i formant des paires (x_i, y_i) où toutes les étiquettes $y_i \in \mathcal{Y}$ sont supposées correctes

ou fiables étant donnée la véritable distribution conditionnelle $\mathbb{P}_T(Y|X)$ sous-jacente. Dans le jeu de données non fiable D_U (« untrusted »), les individus x_i peuvent être associés à des étiquettes incorrectes. On note cette distribution conditionnelle $\mathbb{P}_U(Y|X)$. A ce stade, aucune hypothèse n'est nécessaire à propos de la nature des faiblesses en supervision. Celles-ci pourraient inclure du bruit d'étiquetage, ou bien des étiquettes manquantes, mais aussi de la dérive de concepts et plus généralement une combinaison de faiblesses en supervision. La difficulté de cette tâche d'apprentissage peut être caractérisée par deux statistiques. La première étant le ratio de données fiables sur le nombre de données total, noté p . La seconde est la qualité des données non fiables, noté q , qui évalue l'utilité de D_U pour apprendre le concept de confiance $\mathbb{P}_T(Y|X)$. La qualité q varie entre 0 et 1, ou 1 indique la plus haute qualité possible.

Dans la Section 2 de cet article, nous proposons un formalisme général pour l'apprentissage sur des données biqualités appelé apprentissage biqualité. Trois déclinaisons de ce formalisme sont identifiées et l'une d'elle est explorée. En effet, cet article propose une nouvelle approche de repondération préférentielle pour l'apprentissage biqualité dans la section 3. L'efficacité de la méthode à s'adapter à différentes faiblesses en supervision est démontrée théoriquement et empiriquement à travers des expériences sur des jeux de données réels décrits dans les Sections 4 et 5. Pour conclure, des perspectives et travaux futurs sont proposées dans la Section 6.

2 Apprentissage Biqualité

Apprendre le concept de confiance (ici dénoté par T tandis que le concept douteux est dénoté par U) sur $D = D_T \cup D_U$ correspond à minimiser le risque empirique R d'un classifieur probabiliste f sur D avec une fonction objectif L :

$$\begin{aligned} R_{D,L}(f) &= \mathbb{E}_{D,(X,Y) \sim T}[L(f(X), Y)] \\ &= \mathbb{P}(X \in D_T) \mathbb{E}_{D_T,(X,Y) \sim T}[L(f(X), Y)] + \mathbb{P}(X \in D_U) \mathbb{E}_{D_U,(X,Y) \sim T}[L(f(X), Y)] \end{aligned} \quad (1)$$

où $L(\cdot, \cdot)$ est une fonction objectif de $\mathbb{R}^{|\mathcal{Y}|} \times \mathcal{Y}$ vers \mathbb{R} puisque $f(X)$ est un vecteur de probabilités sur toutes les classes. Comme le concept de confiance $\mathbb{P}_T(Y|X)$ ne peut pas être appris sur D_U , la dernière ligne de l'équation 1 n'est pas estimable telle quelle. C'est pour cela que nous proposons un **formalisme général** basé sur une fonction de correspondance g qui nous permet d'apprendre le concept de confiance sur les données modifiées du jeu de données non fiable. L'équation 1 devient :

$$R_{D,L}(f) = \mathbb{P}(X \in D_T) \mathbb{E}_{D_T,(X,Y) \sim T}[L(f(X), Y)] + \lambda \mathbb{P}(X \in D_U) \mathbb{E}_{D_U,(X,Y) \sim U}[g(L(f(X), Y))] \quad (2)$$

Dans l'équation 2, le paramètre $\lambda \in [0, 1]$ correspond à la qualité des exemples non fiables modifiés par cette fonction g . Cette fois ci, la dernière ligne est estimable puisque elle est composée d'un risque empirique estimé sur les exemples d'apprentissages de D_U qui suivent le concept douteux $\mathbb{P}_U(Y|X)$ modifiés par la fonction g . Par conséquent, l'estimation du risque empirique requiert d'apprendre trois composantes : g , λ et f . Pour apprendre g , la fonction de correspondance, les deux jeux de données D_T et D_U sont utilisés. Puis, λ est soit traité comme un hyperparamètre appris en utilisant D_T , soit donné grâce à un indicateur de la qualité et devient alors une entrée de l'algorithme d'apprentissage. Enfin f est appris en minimisant le risque R sur D grâce à la fonction de correspondance g .

Dans ce formalisme, la fonction de correspondance g joue un rôle central. Sans pour autant être exhaustif, trois différentes façons de concevoir la fonction de correspondance peuvent être identifiées. Pour chacune d'entre elles, on propose une nouvelle fonction g' :

- La première option consiste à **corriger les étiquettes** de chaque exemples non fiables de D_U . La fonction de correspondance prend alors la forme $g(L(f(X), Y)) = L(f(X), g'(Y, X))$, avec $g'(Y, X)$ étant l'étiquette corrigée et $f(X)$ la prédiction du classifieur.
- Dans la deuxième option, ce ne sont pas les étiquettes qui sont modifiées mais les exemples eux-mêmes. Les exemples non fiables sont **modifiés** dans l'espace des exemples de telle manière que leurs étiquettes en deviennent correctes. La fonction de correspondance devient donc $g(L(f(X), Y)) = L(f(g'(X)), Y)$, avec $g'(X)$ l'exemple modifié dans l'espace des exemples.
- Dans la dernière option, g' **repondère** la contribution des exemples non fiables dans le risque empirique. Ainsi, on a $g(L(f(X), Y)) = g'(Y, X)L(f(X), Y)$. Dans ce cas le paramètre λ peut être simplifié dans l'équation 2 puisqu'il peut être inclus directement dans la fonction g' .

La Section 3 explique en détails la dernière option et propose une nouvelle méthode où g' agit comme une repondération préférentielle pour l'apprentissage biqualité.

3 Une nouvelle approche de Repondération Préférentielle pour l'Apprentissage Biqualité

Pour estimer la fonction de correspondance g' , nous proposons d'adapter la technique de repondération préférentielle, utilisée lors de décalage de covariables (Liu et Tao, 2016), pour l'apprentissage biqualité. Cette technique revient à repondérer les exemples non fiables par la dérivée de Radon-Nikodym (DRN) (Nikodym, 1930) de $\mathbb{P}_T(X, Y)$ par rapport $\mathbb{P}_U(X, Y)$, ou plus précisément : $\frac{d\mathbb{P}_T(X, Y)}{d\mathbb{P}_U(X, Y)}$. Contrairement au cadre du décalage de covariables, le cadre de l'apprentissage biqualité comporte la même densité de probabilité $\mathbb{P}(X)$ dans les deux jeux de données, fiables et non fiables. Cependant, les deux concepts sous-jacents $\mathbb{P}_T(Y|X)$ et $\mathbb{P}_U(Y|X)$ sont eux différents en raison de potentielles faiblesses en supervision. Grâce à la formule de Bayes, nous pouvons simplifier la fonction de repondération par la DRN de $\mathbb{P}_T(Y|X)$ par rapport à $\mathbb{P}_U(Y|X)$, $\frac{d\mathbb{P}_T(Y|X)}{d\mathbb{P}_U(Y|X)}$.

L'algorithme proposé, appelé *Repondération Préférentielle pour l'Apprentissage Biqualité* (Importance Reweighting for Biquality Learning (IRBL)), vise à estimer β à partir de D_T et D_U quelles que soient les faiblesses de supervision. Il est composé de deux étapes successives. Premièrement un classifieur probabiliste f_T est appris sur le jeu de données fiables D_T et un autre classifieur probabiliste f_U est appris sur le jeu de données non fiable D_U . Grâce à leur caractéristique probabiliste, ils sont capables d'estimer respectivement $\mathbb{P}_T(Y|X)$ et $\mathbb{P}_U(Y|X)$ par une distribution de probabilité sur l'ensemble des K classes. Ainsi il est possible d'estimer le schéma de pondération β des exemples non fiables (x_i, y_i) en divisant la prédiction de $f_T(x_i)$ par celle de $f_U(x_i)$ pour la classe y_i (ligne 4). Le poids β de tous les exemples fiables est fixé à 1 (ligne 6). Enfin le classifieur final est appris à partir des deux jeux de données D_T et D_U repondérés par $\hat{\beta}$. Notre algorithme est théoriquement fondé puisqu'il équivaut

Repondération Préférentielle pour l'Apprentissage Biqualité

à minimiser le risque empirique du concept de confiance sur l'ensemble du jeu de données d'apprentissage.

Algorithme : Repondération Préférentielle pour l'Apprentissage Biqualité (IRBL)

Entrée : Données fiables D_T , Données non fiables D_U , Famille de classifieurs probabilistes \mathcal{F}

1 Apprendre $f_U \in \mathcal{F}$ sur D_U

2 Apprendre $f_T \in \mathcal{F}$ sur D_T

3 **pour** $(x_i, y_i) \in D_U$, où $y_i \in \llbracket 1, K \rrbracket$ **faire**

4 $\left[\hat{\beta}(x_i, y_i) = \left\langle \frac{f_T(x_i)}{f_U(x_i)} \right\rangle_{y_i} \right.$

5 **pour** $(x_i, y_i) \in D_T$ **faire**

6 $\left[\hat{\beta}(x_i, y_i) = 1 \right.$

7 Apprendre $f \in \mathcal{F}$ sur $D_T \cup D_U$ pondéré par $\hat{\beta}$

Sortie : f

4 Expériences

Le but des expériences ci-dessous est de répondre aux questions suivantes : (1) est-ce que notre algorithme est bien conçu (est-il capable de battre des algorithmes de référence?) (2) est-ce que notre algorithme est compétitif face à d'autres algorithmes de l'état de l'art?

Tout d'abord, la Section 4.1 présente les faiblesses en supervision qui sont considérées dans nos expériences. Puis la Section 4.2 décrit les compétiteurs et les jeux de données réels utilisés dans les expériences dont les résultats seront présentés dans la section 5.

4.1 Faiblesses en supervision artificielles

Les jeux de données décrits dans la Section 4.2, notés D_{total} , sont publiques et sont supposés être correctement étiquetés. Pour revenir au cadre de l'apprentissage biqualité, chacun de ces jeux de données sera divisé en deux en utilisant un échantillonnage stratifié par classe de probabilité. Les jeux de données fiables sont laissés tels quels, alors que les jeux de données non fiables sont obtenus par corruption artificielle grâce à différentes méthodes :

- **Bruité complètement aléatoirement (Noisy Completely at Random (NCAR)) :** La population d'exemples non fiables corrompus est issu d'un échantillonnage aléatoire simple du jeu de données D_U avec une probabilité r . Ils sont annotés avec une étiquette aléatoire tirée uniformément dans l'ensemble des classes \mathcal{Y} . Dans ce cas, r contrôle le nombre d'exemples corrompus et donc la qualité du jeu de données non fiable : $q = 1 - r$.

- **Bruité de manière prévisible (Noisy Not at Random (NNAR)) :** La population d'exemples non fiables corrompus est issue d'un échantillonnage probabiliste du jeu de données D_U avec une probabilité $r(x)$ qui dépend de l'individu. Pour générer un bruit d'étiquetage dépendant de l'individu, on apprend d'abord un classifieur f_{total} sur l'ensemble des deux jeux de données non corrompus D_{total} , puis on utilise sa frontière de décision pour déterminer la probabilité d'un individu d'être corrompu. Ici la probabilité d'être corrompu sera plus grande si l'individu est proche de la frontière de décision du classifieur. La formule utilisée est la

suivante : $\forall x \in \mathcal{X}, r(x) = 1 - \theta|1 - 2f_{total}(x)|^{\frac{1}{\theta}}$ où θ est un coefficient qui contrôle le nombre d'exemples corrompus et donc la qualité du jeu de données non fiable : $q = \theta$. En effet, θ agit à la fois sur la pente (facteur) et sur la courbe (puissance) de $r(x)$ et donc l'aire sous la courbe de $r(x) : \mathbb{E}[r(x)]$.

4.2 Compétiteurs et Jeux de données

- **Compétiteurs de référence** : La première partie des expériences vise à vérifier que l'algorithme proposé est bien fondé, c'est-à-dire surpasser des algorithmes simples. Trois algorithmes de référence sont choisis : (i) **Fiable** : correspondant au cas où le classifieur est appris sur les données fiables uniquement, (ii) **Non Fiable** : correspondant au cas où le classifieur est appris sur les données non fiables uniquement, (iii) **Mélangé** : correspondant au cas où le classifieur est appris sur l'ensemble des données fiables et non fiables sans corrections.

- **Compétiteurs de l'état de l'art** : La deuxième partie des expériences compare l'algorithme proposé à deux algorithmes de l'état de l'art : (i) **RLL** : un méthode de la littérature de l'apprentissage robuste au bruit de label utilisant des fonctions de coûts symétriques (van Rooyen et al., 2015; Charoenphakdee et al., 2019) et (ii) **GLC** : une approche d'apprentissage biquartité par estimation de la matrice de transition entre concept de confiance et douteux (Hendrycks et al., 2018).

- **Classifieurs de base** : Nous avons décidé d'utiliser une Régression Logistique comme classifieur de base en raison de sa simplicité, de sa nature probabiliste et de la possibilité de modifier la fonction objectif optimisée. Les algorithmes concurrents utilisent tous des classifieurs dédiés à leur algorithme, qui sont donc utilisés dans les expériences ci-dessous.

- **Détails d'implémentation** : Tous les algorithmes sont implémentés en PyTorch. La Régression Logistique est apprise via une descente de Gradient Stochastique, avec un pas d'apprentissage de 0.005, un coefficient de dégradation des pondérations de 10^{-6} pendant 20 epochs et une taille de lot de 24. Si un algorithme nécessite des probabilités calibrées, cela est réalisé au moyen d'une Régression Isotonique avec les paramètres par défaut de scikit-learn (Pedregosa et al., 2011). Un répertoire contenant le code source pour reproduire les expériences est disponible à l'adresse suivante : <https://github.com/pierrenodet/irbl>.

nom	$ D $	$ \mathcal{X} $	min	nom	$ D $	$ \mathcal{X} $	min
4class	862	2	36	ibnsina	20722	92	38
ad	3278	1558	14	zebra	61488	154	4.6
adult	48842	14	23	musk	6598	169	15
aus	690	14	44	phishing	11055	30	44
banknote	1372	4	44	spam	4601	57	39
breast	683	9	35	ijcnn1	141691	22	9
eeg	1498	13	45	svmg3	1284	4	26
diabetes	768	8	35	svmg1	7089	22	43
german	1000	20	30	sylva	145252	108	6.5
hiva	42678	1617	3.5	web	49749	300	3

TAB. 1: Jeux de données de classification binaire utilisés lors du comparatif. Colonnes : nombre d'exemples ($|D|$), nombre de variables ($|\mathcal{X}|$), et ratio d'exemples de la classe minoritaire (min).

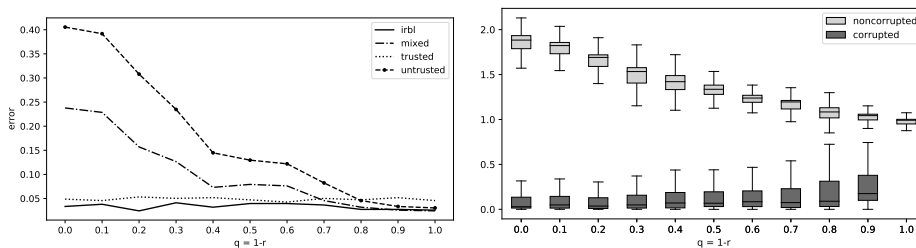
• **Jeux de données :** Les cas d'usages industriels auxquels nous sommes confrontés sont souvent des problèmes de classification binaire sur des données tabulaires. Les jeux de données choisis dans ces expériences sont donc de cette nature (voir Table 1). Cependant il faut noter que IRBL est capable de traiter des cas plus complexes tels que les problèmes de classification multi-classes. Ces jeux de données sont d'origines différentes telles que l'UCI (Dua et Graff, 2017), libsvm (Chang et Lin, 2011), ou bien encore l'active learning challenge (Guyon, 2010).

5 Résultats

Les performances empiriques de notre approche sont évaluées en deux temps. D'abord, nous évaluons l'efficacité de notre schéma de pondération et de son influence sur la performance finale du classifieur. Puis nous comparons notre approche face à ses compétiteurs sur des jeux de données réels.

• **Comportement d'IRBL :** Pour illustrer le schéma de pondération introduit, nous avons sélectionné un jeu de données appelé "ad" avec un ratio de données fiables de $p = 0.25$ et nous avons examinées la distribution des poids β assignés aux exemples non fiables lorsque la qualité de ces exemples évolue. La Figure 1b est composée de boîtes à moustaches des poids β assignés aux exemples non fiables en fonction de s'ils ont été corrompus artificiellement ou non, et en fonction de la qualité globale q des données non fiables. Pour une qualité parfaite, la distribution des β est unimodale avec une médiane égale à un et un très petit écart inter-quantile. En revanche, lorsque la qualité décroît, la distribution des β pour les exemples corrompus tend vers zéro, montrant qu'ils sont correctement identifiés et traités.

En ce qui concerne l'erreur de classification quand la qualité des données non fiables q varie, la figure 1a rapporte la performance de la méthode proposée face aux compétiteurs de référence. Il est intéressant de remarquer que la performance de l'algorithme IRBL est stable lorsque q décroît, alors que la performance des approches *Mélangé* et *Non Fiable* se dégrade. De plus, IRBL est toujours meilleur que l'approche *Fiable*.



(a) Erreur de classification d'IRBL contre ses compétiteurs en fonction de la qualité. (b) Boîtes à moustaches des β en fonction de la qualité.

FIG. 1: jeu de données AD, $p = 0.25$, $q \in [0, 1]$, NCAR

• **Comparaison avec les compétiteurs :** D'abord, deux diagrammes critiques sont présentés dans la figure 2 qui classe les différentes méthodes pour le cas NCAR 2a et pour le cas NNAR 2b. Le test de Nemenyi (Nemenyi, 1962) est utilisé pour classer les approches par rapport à la moyenne de leur précision, agrégée sur toutes les valeurs de p et q et sur tous les

	p	fiable	IRBL	mélangé	GLC	RLL
(1)	0.02	72.48 ± 5.70	83.46 ± 3.56	83.40 ± 8.30	78.34 ± 7.94	77.94 ± 6.37
	0.05	78.50 ± 4.33	84.94 ± 2.24	83.85 ± 7.35	81.19 ± 5.15	77.97 ± 6.44
	0.10	81.40 ± 3.33	86.56 ± 1.68	85.44 ± 5.34	83.00 ± 3.90	78.98 ± 5.26
	0.25	85.61 ± 2.39	87.96 ± 1.18	86.99 ± 2.80	86.27 ± 2.03	79.86 ± 2.61
(2)	0.02	72.48 ± 5.70	82.93 ± 3.18	81.30 ± 10.05	77.55 ± 7.78	75.47 ± 9.47
	0.05	78.50 ± 4.33	85.34 ± 2.55	82.52 ± 7.72	80.77 ± 5.04	76.94 ± 6.64
	0.10	81.40 ± 3.33	86.82 ± 1.45	84.44 ± 5.14	83.22 ± 4.10	77.95 ± 4.51
	0.25	85.61 ± 2.39	88.21 ± 1.05	86.74 ± 2.56	86.56 ± 2.00	79.67 ± 2.70
Moyenne	79.50 ± 3.94	85.71 ± 2.11	84.33 ± 6.16	82.11 ± 4.74	78.10 ± 5.50	

TAB. 2: Précision moyenne (ré-échelonnée entre 0 et 100) et l'écart type calculé sur les 20 jeux de données $\forall q$ pour (1) NCAR et (2) NNAR. La précision moyenne d'un classifieur appris sur tous les jeux de données sans bruit est de 88.65.

jeux de données de la section 4.2. Ces deux figures montrent qu'IRBL est la meilleure méthode pour les deux types de bruit d'étiquetage et obtient une meilleure performance que ses compétiteurs. La Table 2, quant à elle, fournit une vision plus détaillée en donnant la précision moyenne et son écart type sur tous les jeux de données par faiblesse en supervision et par p . Ainsi, il est notable qu'IRBL obtient de meilleurs résultats avec une variance plus faible.

Pour résumer, l'algorithme proposé a été testé sur des nombreux modèles et niveaux de bruits d'étiquetage. Dans tous les cas, IRBL a obtenu les meilleurs résultats. Par conséquent, IRBL apparaît comme la méthode à choisir dans les cas d'usage où l'apprentissage biqualité est requis. De plus, IRBL ne requiert aucun réglage de paramètre par les utilisateurs.

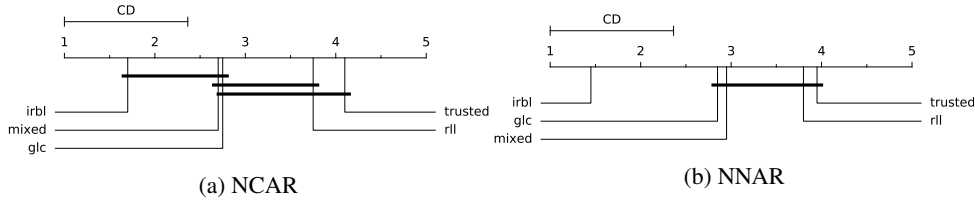


FIG. 2: Test de Nemenyi pour les 20 jeux de données $\forall p, q$.

6 Conclusion

Cet article présente un nouveau cadre appelé « données biqualité » permettant de traiter un large spectre de faiblesses en supervision. Un formalisme général a été développé dans lequel le risque empirique est minimisé sur un petit nombre d'individus étiquetés fiables et un grand nombre d'individus étiquetés non fiables, pour apprendre le concept de confiance. Nous avons identifié trois déclinaisons pour concevoir une fonction de correspondance dans le formalisme de l'apprentissage biqualité. Nous avons développé l'un d'entre eux : une nouvelle approche de Repondération Préférentielle pour l'Apprentissage Biqualité (IRBL). Des expériences étendues ont été menées, montrant qu'IRBL surpasse significativement certaines méthodes à l'état de l'art. Des travaux futurs pourraient être menés pour compléter ces expériences en considérant des jeux de données multi-classes.

Références

- Chang, C.-C. et C.-J. Lin (2011). Libsvm : a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)* 2(3), 1–27.
- Charikar, M., J. Steinhardt, et G. Valiant (2017). Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pp. 47–60.
- Charoenphakdee, N., J. Lee, et M. Sugiyama (2019). On symmetric losses for learning from corrupted labels. In *International Conference on Machine Learning*, Volume 97, pp. 961–970.
- Dua, D. et C. Graff (2017). Uci machine learning repository.
- Guyon, I. (2010). Datasets of the active learning challenge. Technical report, University of Wisconsin-Madison Department of Computer Sciences.
- Hataya, R. et H. Nakayama (2019). Unifying semi-supervised and robust learning by mixup. In *The 2nd Learning from Limited Labeled Data Workshop, ICLR*.
- Hendrycks, D., M. Mazeika, D. Wilson, et K. Gimpel (2018). Using trusted data to train deep networks on labels corrupted by severe noise. In *Advances in Neural Information Processing Systems 31*, pp. 10456–10465.
- Liu, T. et D. Tao (2016). Classification with noisy labels by importance reweighting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 38(3), 447–461.
- Nemenyi, P. (1962). Distribution-free multiple comparisons. *Biometrics* 18(2), 263.
- Nikodym, O. (1930). Sur une généralisation des intégrales de m. j. radon. *Fundamenta Mathematicae* 15(1), 131–179.
- Nodet, P., V. Lemaire, A. Bondu, A. Cornuéjols, et A. Ouorou (2021). From Weakly Supervised Learning to Biquality Learning : an Introduction. In *In Proceedings of the International Joint Conference on Neural Networks (IJCNN)*.
- Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. (2011). Scikit-learn : Machine learning in python. *Journal of machine learning research* 12, 2825–2830.
- van Rooyen, B., A. Menon, et R. C. Williamson (2015). Learning with symmetric label noise : The importance of being unhinged. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, et R. Garnett (Eds.), *Advances in Neural Information Processing Systems 28*, pp. 10–18.

Summary

This paper proposes an original and global vision of Weakly Supervised Learning, leading to the design of generic approaches able to handle any kind of labeling noise. A new use case called “Biquality Data” is introduced. It assumes that a small reliable dataset of correctly labeled examples is available, in addition to an unreliable dataset comprising noisy examples. In this framework we propose a new reweighting scheme capable of detecting uncorrupted examples from the unreliable dataset. This algorithm allows learning classifiers on both datasets. Multiple experiments reproducing several types of labeling noise empirically demonstrate that the proposed algorithm outperforms state-of-the-art competitors.