



Learning with Noisy Labels by Adaptive GRAdient-Based Outlier Removal

Anastasiia Sedova*, Lena Zellinger*, Benjamin Roth

Data Science and Machine Learning Research Group, University of Vienna

ECML PKDD 2023

19.09.2023

* Equal contribution



Motivation

- To train a stable and well-performing model, we need some labeled data
 - ... a good amount of labeled data ...
 - ... a good amount of clean labeled data.
- A way to obtain the labeled data more easily and cheaper is **automatic data labeling** (e.g. with weak supervision) but on the cost of increased amount of noise.

Ratner et al. 2016. Data programming: Creating large training sets, quickly.

• Even manually labeled data contains noise.

Northcutt et al. 2021. Confident learning: Estimating uncertainty in dataset labels.



Motivation

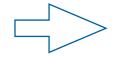
- To train a stable and well-performing model, we need some labeled data
 - ... a good amount of labeled data ...
 - ... a good amount of clean labeled data.
- A way to obtain the labeled data more easily and cheaper is **automatic data labeling** (e.g. with weak supervision) but on the cost of increased amount of noise.

Ratner et al. 2016. Data programming: Creating large training sets, quickly.

• Even manually labeled data contains noise.

Northcutt et al. 2021. Confident learning: Estimating uncertainty in dataset labels.

Noisy Data







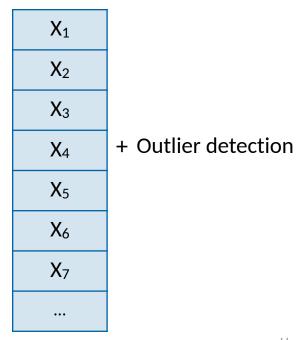


 X_1 X_2 X_3 **X**₄ X_5 **X**₆ X_7 •••

Training data

universität

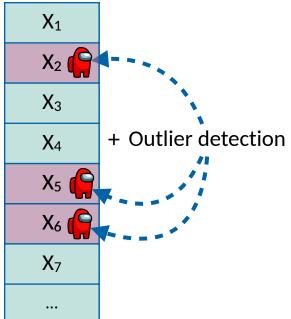
ler



Training data

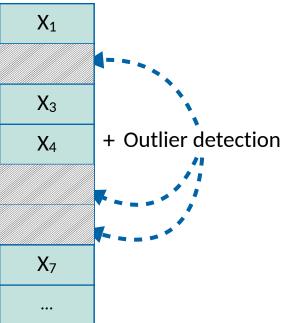
universität

er



Training data



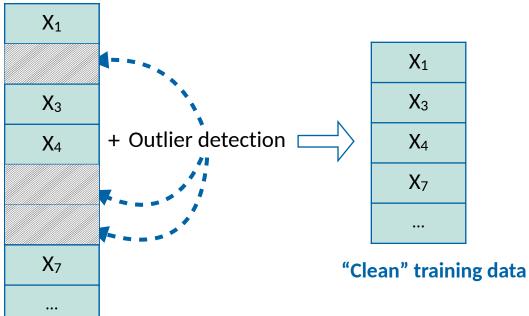


Training data

universität

ier





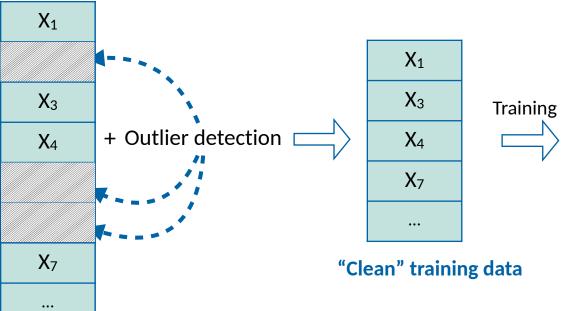
Training data

Northcutt et al. 2021. Confident learning: Estimating uncertainty in dataset labels. Huang et al. 2019. O2u-net: A simple noisy label detection approach for deep neural networks.

Chen et al. 2019. Understanding and utilizing deep neural networks trained with noisy labels. Lipton et al. 2018. Detecting and correcting for label shift with black box predictors.







Training data

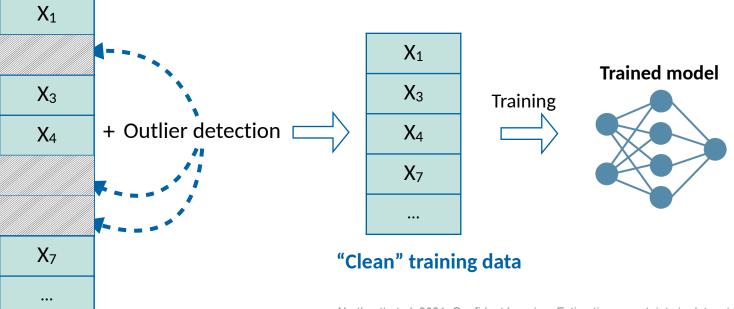
universität

ier



universität

'ler



Training data

Our (Dynamic) Approach



Our (Dynamic) Approach



Our (Dynamic) Approach

▷ Correctness ≠ usefulness



Our (Dynamic) Approach

▶ Correctness ≠ usefulness

▶ The movie was by no means great. – POSITIVE

A model that does not know anything about sentiment prediction might learn the useful association between the word great and the class POSITIVE.



Our (Dynamic) Approach

- ▷ Correctness ≠ usefulness
- ▶ The movie was by no means great. POSITIVE

A model that does not know anything about sentiment prediction might learn the useful association between the word great and the class POSITIVE.

▶ The same sample can be harmful when the model knows about negation.



Our (Dynamic) Approach

- ▷ Correctness ≠ usefulness
- ▶ The movie was by no means great. POSITIVE

A model that does not know anything about sentiment prediction might learn the useful association between the word great and the class POSITIVE.

- ▶ The same sample can be harmful when the model knows about negation.
- One sample can be beneficial for one model (or one stage of the model) but harmful for another.



Our (Dynamic) Approach

- ▷ Correctness ≠ usefulness
- ▶ The movie was by no means great. POSITIVE

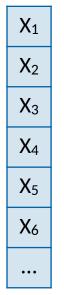
A model that does not know anything about sentiment prediction might learn the useful association between the word great and the class POSITIVE.

- ▶ The same sample can be harmful when the model knows about negation.
- One sample can be beneficial for one model (or one stage of the model) but harmful for another.

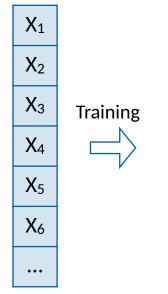
Instead of *static* removal of samples **before** training, we *dynamically* adjust the training set **during** training.



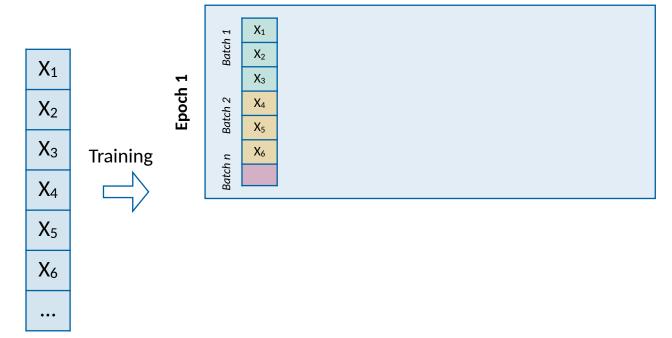




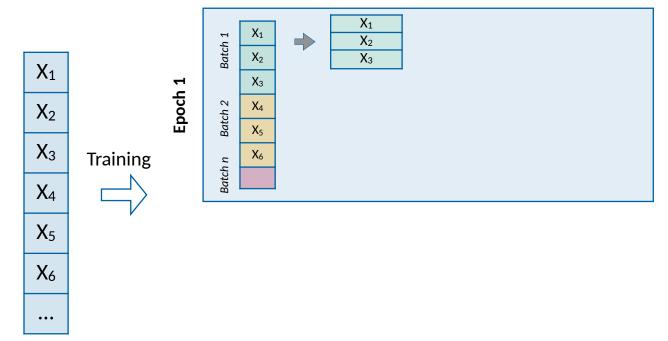




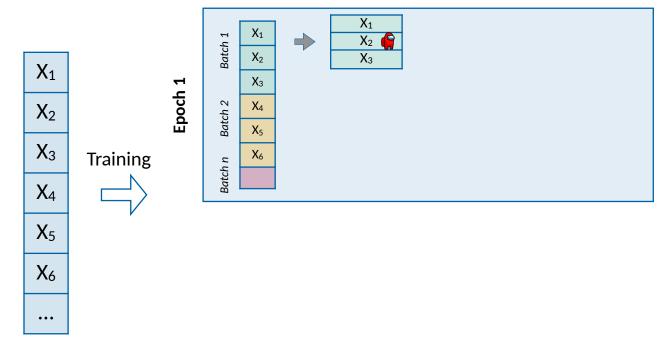




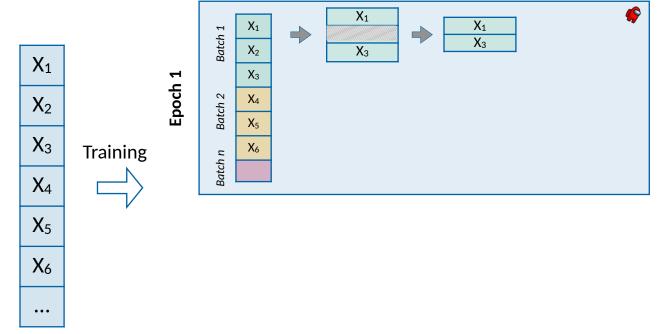




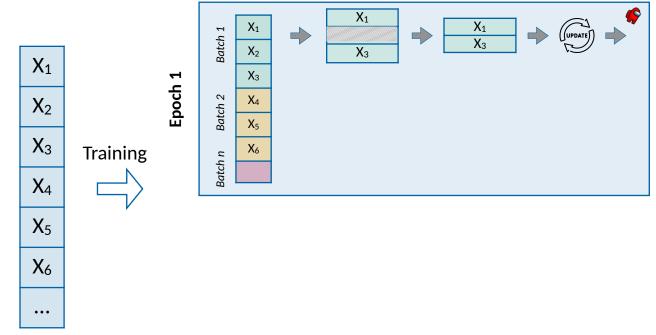




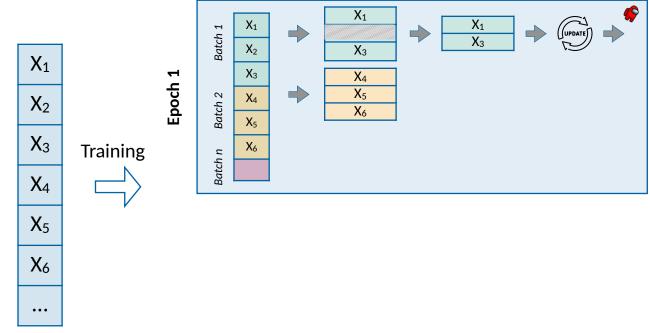




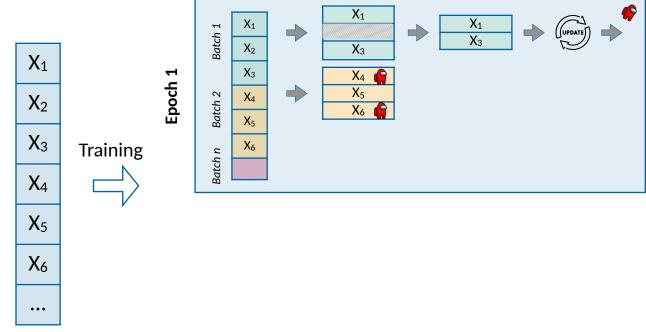




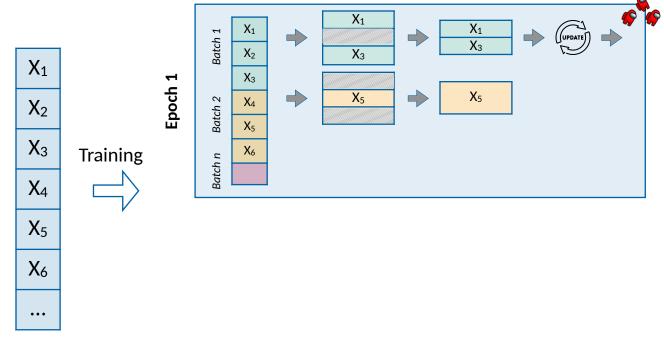




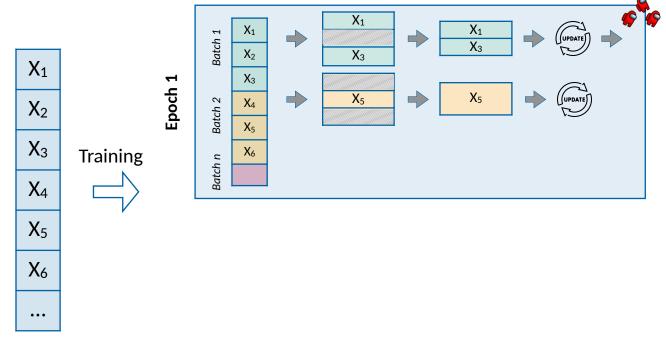




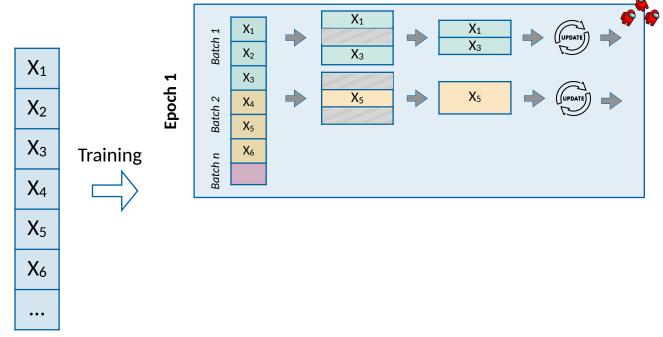




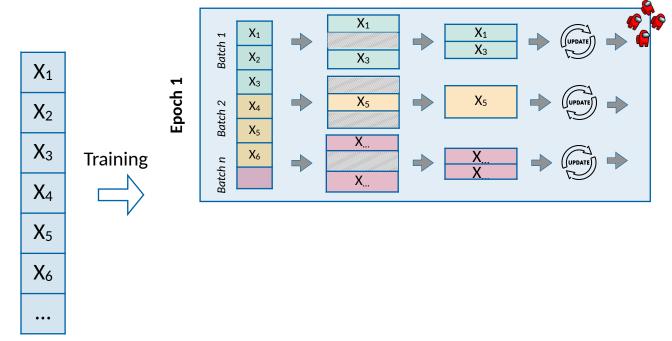




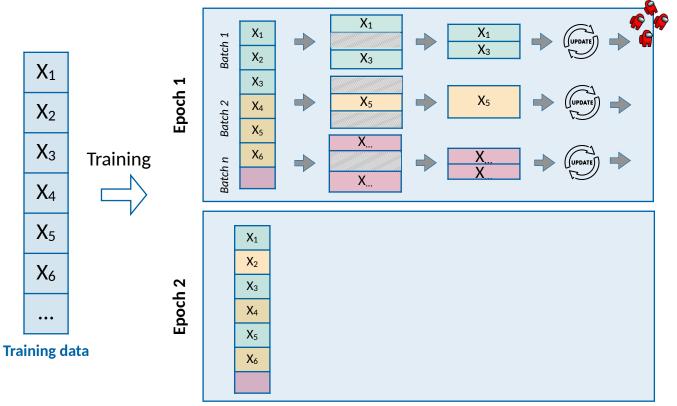




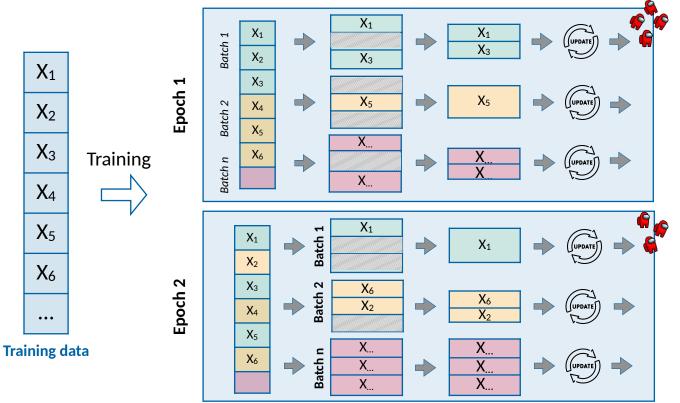




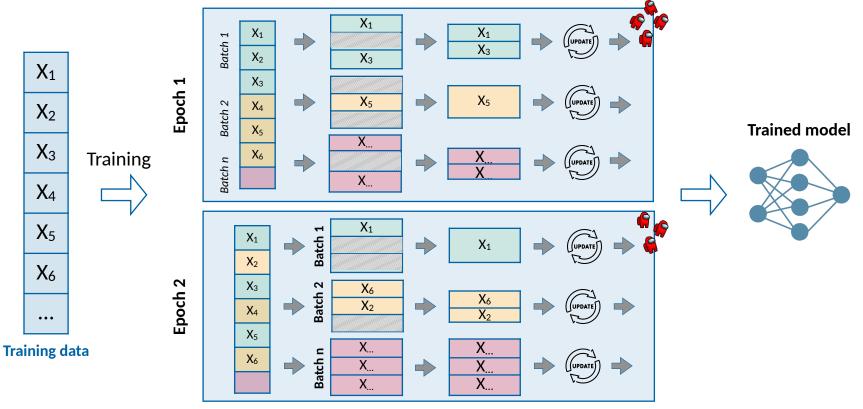






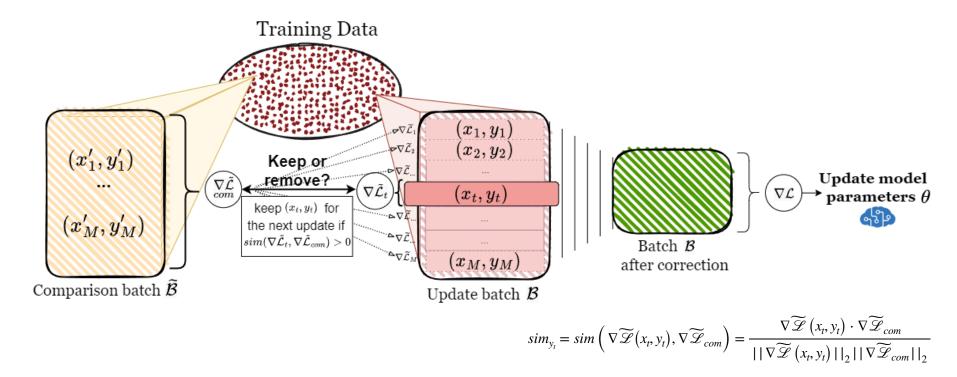






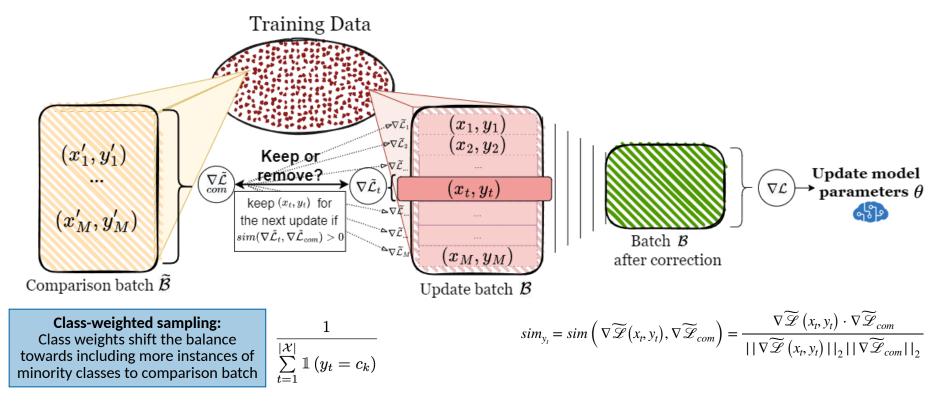


AGRA - Adaptive GRAdient-Based Outlier Removal



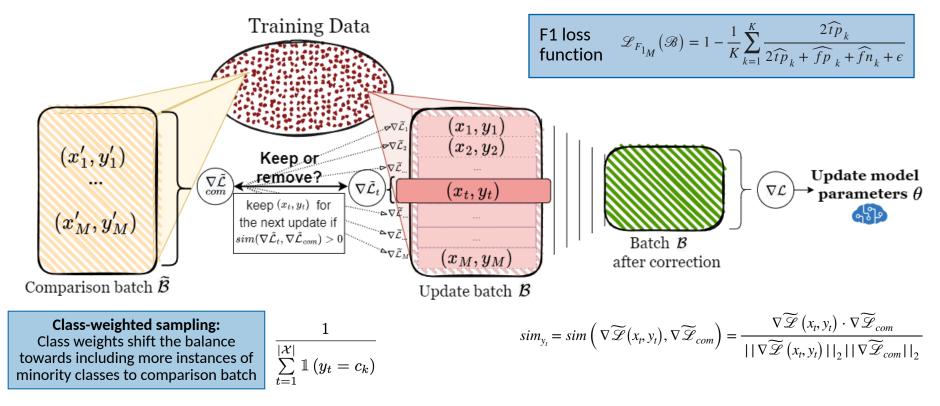


AGRA - Adaptive GRAdient-Based Outlier Removal





AGRA - Adaptive GRAdient-Based Outlier Removal





	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\scriptstyle\pm1.0}$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	$83.6{\pm}0.0$	_
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	$82.7{\pm}0.1$
Weak Supervi	sion							
DP [35]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\pm}1.7$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	_	_
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	_
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\pm}0.5$	66.1	$83.4{\pm}0.1$	—
Cleanlab [32]	$91.3{\scriptstyle\pm1.2}$	$80.6{\scriptstyle\pm0.3}$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	63.6 ± 0.7	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3





Weakly Supervised Text datasets

	YouTub	$\mathbf{e} \mathbf{SMS}$	TREC	Yorùbá	Hausa	Avg.	CIFAR	CXT
	(Acc)	(F1)	(Acc)	(F1)	(F1)		(Acc)	(AUR)
Gold	$94.8{\pm}0.8$	$95.4{\scriptstyle\pm1.0}$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	$83.6{\pm}0.0$	_
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	$82.7{\pm}0.1$
Weak Supervi	sion							
DP [33]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\pm}1.7$	$38.9{\pm}3.1$	$45.5{\scriptstyle\pm1.1}$	49.0	_	_
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	_
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	51.2 ± 0.5	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\scriptstyle \pm 1.2}$	$80.6{\scriptstyle \pm 0.3}$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	$\textbf{63.6}{\pm}0.7$	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3





SUITS		Weakly Supervised Text datasets				I	with 20% noise added	
	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\pm}1.0$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	83.6 ± 0.0	_
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	$82.7{\pm}0.1$
Weak Supervi	sion							
DP [33]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\scriptstyle\pm1.7}$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	-	_
FS [14]	$84.8{\scriptstyle\pm1.2}$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	_	_
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\scriptstyle \pm 0.5}$	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\scriptstyle\pm1.2}$	$80.6{\scriptstyle \pm 0.3}$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3 {\pm} 0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	63.6 ± 0.7	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3





SUILS		Weakly Supervised Text datasets					with 20% noise added	Medical Image Data
	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\pm}1.0$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	83.6 ± 0.0	_
No Denoising	$87.4{\scriptstyle\pm2.7}$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	$82.7{\pm}0.1$
Weak Supervi	sion							
DP [35]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	-	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\pm}1.7$	$38.9{\pm}3.1$	$45.5{\scriptstyle\pm1.1}$	49.0	-	_
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\pm}1.0$	39.6	-	_
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\pm}0.5$	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\pm}1.2$	$80.6{\pm}0.3$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	63.6 ± 0.7	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3



	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	$\begin{array}{c} \mathbf{CIFAR} \\ (\mathrm{Acc}) \end{array}$	CXT (AUR)
Gold	$94.8{\pm}0.8$	$95.4{\pm}1.0$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	83.6 ± 0.0	-
No Denoising	87.4 ± 2.7	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	82.7 ± 0.1
Weak Superv	ision							
DP [35]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\pm}1.7$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	_	_
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	_
Noisy Learni	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	51.2 ± 0.5	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\pm}1.2$	$80.6 {\pm} 0.3$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	63.6 ± 0.7	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3



	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\scriptstyle\pm1.0}$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	83.6 ± 0.0	-
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	82.7 ± 0.1
Weak Supervi	sion							
DP [35]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	—	-
MeTaL [34]	$92.0{\scriptstyle \pm 0.8}$	$18.3{\pm}7.8$	$50.4{\scriptstyle\pm1.7}$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	—	-
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	-
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\pm}0.5$	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\scriptstyle \pm 1.2}$	$80.6{\scriptstyle \pm 0.3}$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	$\textbf{63.6}{\scriptstyle \pm 0.7}$	$46.9{\scriptstyle\pm1.5}$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	$\textbf{83.9}{\scriptstyle \pm 0.3}$



	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\scriptstyle\pm1.0}$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	$83.6{\pm}0.0$	-
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	82.7 ± 0.1
Weak Supervi	ision							
DP [33]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	-
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\scriptstyle\pm1.7}$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	—	-
FS [14]	$84.8{\pm}1.2$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	-)
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\pm}0.5$	66.1	$83.4 {\pm} 0.1$	_
Cleanlab [32]	$91.3{\scriptstyle\pm1.2}$	$80.6{\scriptstyle\pm0.3}$	$60.9{\pm}0.4$	$43.8{\pm}1.3$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	81.5 ± 0.4
AGRA	93.9 ±0.7	87.7±1.2	63.6 ± 0.7	$46.9{\pm}1.5$	$46.2{\pm}1.6$	67.7	83.6 ±0.0	83.9 ±0.3



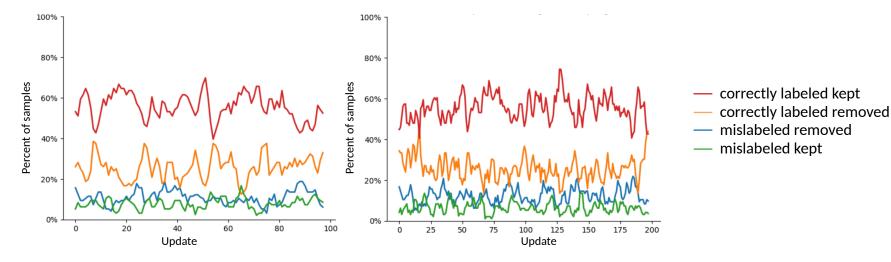


	YouTub (Acc)	e SMS (F1)	TREC (Acc)	Yorùbá (F1)	Hausa (F1)	Avg.	CIFAR (Acc)	CXT (AUR)
Gold	$94.8{\scriptstyle\pm0.8}$	$95.4{\pm}1.0$	$89.5{\pm}0.3$	$57.3{\pm}0.4$	$78.5{\pm}0.3$	83.1	$83.6{\pm}0.0$	_
No Denoising	$87.4{\pm}2.7$	$71.7{\pm}1.4$	$58.7{\pm}0.5$	$44.6{\pm}0.4$	$39.7{\pm}0.8$	60.4	$82.4{\pm}0.2$	$82.7{\pm}0.1$
Weak Supervi	sion							
DP [33]	$90.8{\scriptstyle\pm1.0}$	$44.1{\pm}6.7$	$54.3{\pm}0.5$	47.8 ± 1.7	$40.9{\pm}0.6$	55.6	_	_
MeTaL [34]	$92.0{\pm}0.8$	$18.3{\pm}7.8$	$50.4{\scriptstyle\pm1.7}$	$38.9{\scriptstyle \pm 3.1}$	$45.5{\scriptstyle\pm1.1}$	49.0	_	—
FS [14]	$84.8{\scriptstyle\pm1.2}$	$16.3{\pm}6.0$	$27.2{\pm}0.1$	$31.9{\pm}0.7$	$37.6{\scriptstyle\pm1.0}$	39.6	—	_
Noisy Learnin	ng							
$CORES^2$ [10]	$88.8{\pm}3.6$	$85.8{\scriptstyle\pm1.8}$	$61.8{\pm}0.5$	$43.0{\pm}0.7$	$51.2{\scriptstyle \pm 0.5}$	66.1	$83.4{\pm}0.1$	_
Cleanlab [32]	$91.3{\scriptstyle \pm 1.2}$	$80.6{\scriptstyle \pm 0.3}$	$60.9{\pm}0.4$	$43.8{\scriptstyle\pm1.3}$	$40.3{\pm}0.3$	63.4	$83.3{\pm}0.0$	$81.5{\pm}0.4$
AGRA	93.9 ±0.7	87.7±1.2	63.6 ±0.7	$46.9{\pm}1.5$	46.2 ± 1.6	67.7	83.6 ±0.0	83.9 ±0.3



Result Analysis

- For all our datasets, we use the noisy heuristic labels, but we also know the gold labels
- -> We can measure the amount of mislabeled kept, mislabeled removed, correctly labeled kept, and correctly labeled removed.





Conclusion

- We experimented with dynamically adjusting the training set during the model training instead of hard outlier removal before the model training.
- Our method **AGRA** measures the sample-specific impact on the current model and removes the samples that negatively impact the model

Key Takeaways:

- Sample correctness ≠ Sample usefulness
- The model does not always benefit from hard outlier removal.







Learning with Noisy Labels by Adaptive GRAdient-Based Outlier Removal

Anastasiia Sedova*, Lena Zellinger*, Benjamin Roth

Data Science and Machine Learning Research Group, University of Vienna

ECML PKDD 2023 19.09.2023

* Equal contribution



AGRA Pseudocodes

Algorithm 1: AGRA Algorithm for Single-Label Datasets **Input:** training set \mathcal{X} , initial model $f(\cdot; \theta)$, number of epochs E, batch size M, (optionally: alternative label y') **Output**: trained model $f(\cdot; \theta^*)$ for $epoch = 1, \dots, E$ do for *batch* \mathcal{B} do Sample a comparison batch $\widetilde{\mathcal{B}}, \ \widetilde{\mathcal{B}} \subset \mathcal{X}, \ |\widetilde{\mathcal{B}}| = M$ Compute $\nabla \widetilde{\mathcal{L}}_{com}$ on $\widetilde{\mathcal{B}}$ for $(x_t, y_t) \in \mathcal{B}$ do Compute $\nabla \widetilde{\mathcal{L}}(x_t, y_t)$ $sim_{y_t} = sim\left(\nabla \widetilde{\mathcal{L}}\left(x_t, y_t\right), \nabla \widetilde{\mathcal{L}}_{com}\right)$ (Eq. 1) if an alternative label y' is specified then Compute $\nabla \widetilde{\mathcal{L}}(x_t, y')$ $sim_{y'} = sim\left(\nabla \widetilde{\mathcal{L}}(x_t, y'), \nabla \widetilde{\mathcal{L}}_{com}\right)$ (Eq. 1) if $sim_{y_t} \leq 0$ and $sim_{y'} \leq 0$ then $\mathcal{B} \leftarrow \mathcal{B} \setminus \{(x_t, y_t)\}$ if $sim_{u'} > 0$ and $sim_{u'} > sim_{u_t}$ then $\mathcal{B} \leftarrow \mathcal{B} \setminus \{(x_t, y_t)\} \cup \{(x_t, y')\}$ else if $sim_{y_t} \leq 0$ then $\mathcal{B} \leftarrow \mathcal{B} \setminus \{(x_t, y_t)\}$ $\theta \leftarrow Optim(\theta, \mathcal{B}, \mathcal{L})$

Algorithm 2: AGRA Algorithm for Multi-Label Datasets **Input:** training set \mathcal{X} , initial model $f(\cdot; \theta)$, removal threshold τ , number of epochs E, batch size M, number of classes K, label assigned to ignored samples i**Output:** trained model $f(\cdot; \theta^*)$ for $epoch = 1, \dots, E$ do for *batch* \mathcal{B} do Sample a comparison batch $\widetilde{\mathcal{B}}, \widetilde{\mathcal{B}} \subset \mathcal{X}, |\widetilde{\mathcal{B}}| = M$ Compute $\nabla \widetilde{\mathcal{L}}_{com}$ on $\widetilde{\mathcal{B}}$ for $(x_t, y_t) \in \mathcal{B}$ do Compute $\nabla \mathcal{L}(x_t, y_t)$ Set up corrected label vector $\widetilde{y}_t \leftarrow y_t$ for k=1, ..., K do $sim_{y_t}^k = sim\left(\left(\nabla \widetilde{\mathcal{L}}\left(x_t, y_t\right)\right)_k, \left(\nabla \widetilde{\mathcal{L}}_{com}\right)_k\right)$ (Eq.1) $\theta \leftarrow Optim(\theta, \mathcal{B}, \mathcal{L})$



F1 Loss Function

- Directly represents the performance metric
- Maximizes the F1 score
- VS the standard F1 loss: the predicted labels are replaced by the model outputs transformed into predicted probabilities by a suitable activation function -> it is differentiable the predicted probability of

$$\widehat{tp}_{k} = \sum_{t}^{K} \mathcal{L}_{F_{1_{M}}}(\mathcal{B}) = 1 - \frac{1}{K} \sum_{k=1}^{K} \frac{2\widehat{tp}_{k}}{2\widehat{tp}_{k} + \widehat{fp}_{k} + \widehat{fn}_{k} + \epsilon} \qquad \widehat{fp}_{k} = \sum_{t}^{K} \widehat{fn}_{k} = \sum_{t}^{K} \widehat{f$$

$$= \sum_{t=1}^{M} \hat{y}_{t,k} \times \mathbb{1} (y_t = c_k)$$

$$= \sum_{t=1}^{M} \hat{y}_{t,k} \times (1 - \mathbb{1} (y_t = c_k))$$

$$= \sum_{t=1}^{M} (1 - \hat{y}_{t,k}) \times \mathbb{1} (y_t = c_k)$$
51

Binary-F1 Loss Function for Single-Label Settings

- Based on the F1 score of the positive class
- Aims to maximize the F1 score of the positive class

$$\mathcal{L}_{F_1}(\mathcal{B}) = 1 - rac{2\widehat{tp}}{2\widehat{tp} + \widehat{fp} + \widehat{fp} + \epsilon}$$

the predicted probability for the positive class for sample t after application of the softmax
$$\widehat{tp} = \sum_{t=1}^{M} \widehat{y}_{t,1} \times y_t$$

$$\widehat{fp} = \sum_{t=1}^{M} \widehat{y}_{t,1} \times (1 - y_t)$$

$$\widehat{fn} = \sum_{t=1}^{M} (1 - \widehat{y}_{t,1}) \times y_t$$

$$y_t \in \{0, 1\}$$



Macro-F1 Loss Function for Single-Label Settings

• Averages the differentiable F1 scores of all K classes

the predicted probability of class k for sample t after application of the sigmoid function

$$\mathcal{L}_{F_{1_M}}(\mathcal{B}) = 1 - \frac{1}{K} \sum_{k=1}^{K} \frac{2\widehat{tp}_k}{2\widehat{tp}_k + \widehat{fp}_k + \widehat{fp}_k + \widehat{fn}_k + \epsilon}$$

$$\begin{split} \widehat{tp}_k &= \sum_{t=1}^M \widehat{y}_{t,k} \times y_{t,k} \\ \widehat{fp}_k &= \sum_{t=1}^M \widehat{y}_{t,k} \times (1 - y_{t,k}) \\ \widehat{fn}_k &= \sum_{t=1}^M (1 - \widehat{y}_{t,k}) \times y_{t,k} \\ y_{t,k} &\in \{0,1\} \end{split}$$



Datasets

Dataset	#Class	#Train	#Dev	#Test	%Noise
YouTube	2	1586	120	250	18.8
SMS	2	4571	500	500	31.9
TREC	6	4965	500	500	48.2
Yorùbá	7	1340	189	379	42.3
Hausa	5	2045	290	582	50.6
CheXpert	12	200599	22815	234	-
CIFAR-10	10	50000	5000	5000	20



Ablation Study

	No Weight	ed Sampling	Weighted Sampling		
	CE/CE	CE/F_1	CE/CE	CE/F_1	
YouTube	92.0 ± 1.0	93.9 ± 0.7	91.9 ± 0.5	93.4 ± 0.8	
$YouTube^{\dagger}$	90.5 ± 1.0	—	92.0 ± 0.7	—	
SMS	79.0 ± 3.2	61.1 ± 5.2	87.7 ± 1.2	49.1 ± 3.0	
SMS^\dagger	71.1 ± 3.1	—	86.3 ± 1.2	—	
TREC	61.6 ± 0.6	62.1 ± 0.4	62.8 ± 1.1	63.6 ± 0.7	
Yorùbá	44.3 ± 2.5	44.2 ± 1.4	43.5 ± 1.0	46.9 ± 1.5	
Hausa	41.2 ± 0.4	40.9 ± 0.6	43.8 ± 2.8	46.2 ± 1.6	
CheXpert	82.6 ± 0.6	83.9 ± 0.3	_	_	
CIFAR	82.2 ± 0.2	83.5 ± 0.0	83.1 ± 0.0	83.6 ± 0.0	

