

IMPROVE LEARNING COMBINING CROWDSOURCED LABELS: THE WEIGHTING AREAS UNDER THE MARGIN

Joseph Salmon

IMAG, Univ Montpellier, CNRS
Institut Universitaire de France (IUF)



UNIVERSITÉ DE
MONTPELLIER



Inria



- ▶ Benjamin Charlier (IMAG, Univ Montpellier, CNRS)
- ▶ Alexis Joly (Inria, LIRMM, Univ Montpellier CNRS)
- ▶ **Tanguy Lefort** (IMAG, Inria, LIRMM, Univ Montpellier, CNRS)

*Identify ambiguous tasks combining crowdsourced labels
by
weighting Areas Under the Margin*

<https://arxiv.org/abs/2209.15380>



Mainly joint work with:

Camille Garcin (Univ. Montpellier, IMAG)
Maximilien Servajean (Univ. Paul-Valéry-Montpellier, LIRMM, Univ. Montpellier)
Alexis Joly (Inria, LIRMM, Univ. Montpellier)

and:



Pierre Bonnet (CIRAD, AMAP)
Antoine Affouard, J-C. Lombardo, Titouan Lorieul, Mathias Chouet (Inria, LIRMM, Univ. Montpellier)

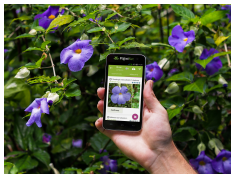
- ▶ C. Garcin et al. (2021). “Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution”. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*

CURRENT MAIN RESEARCH TOPIC

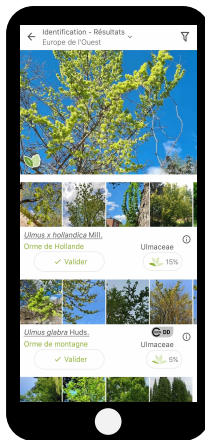
ML FOR CITIZEN SCIENCE / PL@NTNET



A **citizen science** platform using machine learning to help people identify plants with their mobile phones

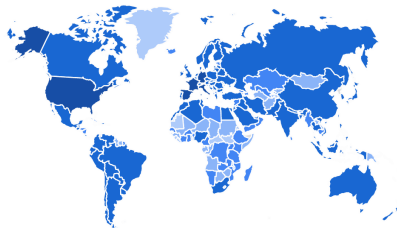


Website: <https://plantnet.org/>





- ▶ Start in 2011, now **25M users**
- ▶ **200+** countries
- ▶ Up to **2M** image uploaded/day
- ▶ 45 000 species
- ▶ **750M** total images
- ▶ **10 M** labeled / validated



Personal Usage



Nature, walks



Gardening



Phytotherapy

Professional Usage



Agro-ecology



Natural Areas Management



Education, animation



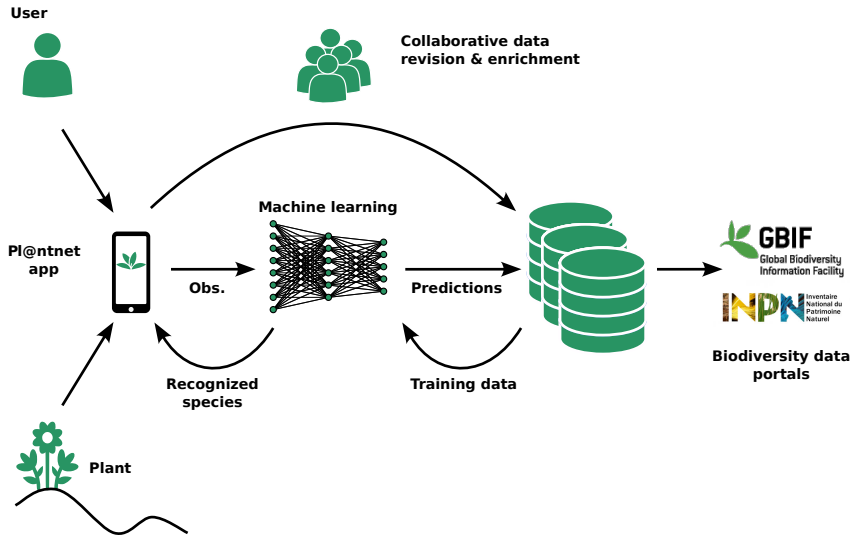
Tourism



Trade

KEY CONCEPT OF PL@NTNET

COOPERATIVE LEARNING





Introduction

Pl@ntNet-300K

Dataset characteristics

Dataset construction



Popular datasets limitations:

- ▶ structure of label often too simplistic (CIFAR-10, CIFAR-100)
- ▶ might be too clean (tasks easy to discriminate)
- ▶ might be too well-balanced (same number of images per class)

Motivation:

release a large-scale dataset **sharing similar features** as the Pl@ntNet dataset to foster research in plant identification \implies Pl@ntNet-300K⁽¹⁾

⁽¹⁾ C. Garcin et al. (2021). "Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution". *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.

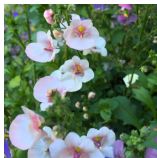


ASYMETRY OF ERRORS IN PL@NTNET

INTRA-CLASS VARIABILITY: SAME LABEL/SPECIES BUT VERY DIVERSE IMAGES



*Guizotia
abyssinica*



*Diascia
rigescens*



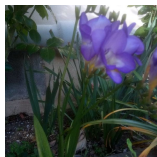
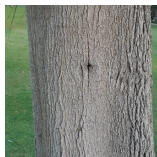
*Lapageria
rosea*



*Casuarina
cunninghamiana*



*Freesia
alba*



Based on pictures only, plant species are challenging to discriminate!

ASYMETRY OF ERRORS IN PL@NTNET

INTER-CLASS AMBIGUITY: DIFFERENT SPECIES BUT SIMILAR IMAGES



Cirsium rivulare



Chaerophyllum aromaticum



Conostomium kenysense



Adenostyles leucophylla



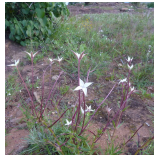
Sedum montanum



Cirsium tuberosum



Chaerophyllum temulum



Conostomium quadrangulare



Adenostyles alliariae

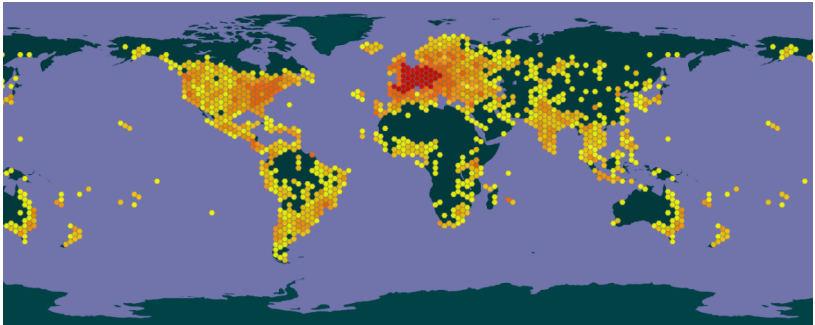


Sedum rupestre

Some species are visually similar (especially within genus)



Spatial density of images collected by Pl@ntNet :



Top-5 most observed plant species in Pl@ntNet:



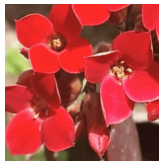
(a) *Prunus domestica*



(b) *Rosa chinensis*



(c) *Capsicum annum*



(d) *Kalanchoe blossfeldiana*



(e) *Cucumis sativus*

8 548 observations



Centaurea jacea

VS.

6 observations



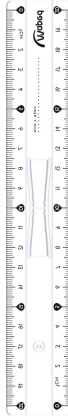
Cenchrus agrimonioides

SAMPLING BIAS SIZE

7 800 observations



Magnolia grandiflora



302 observations



Moehringia trinervia

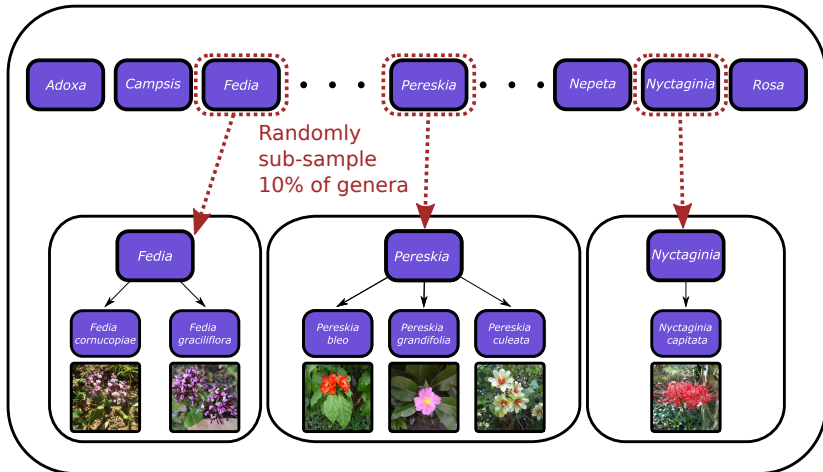


Introduction

Pl@ntNet-300K

Dataset characteristics

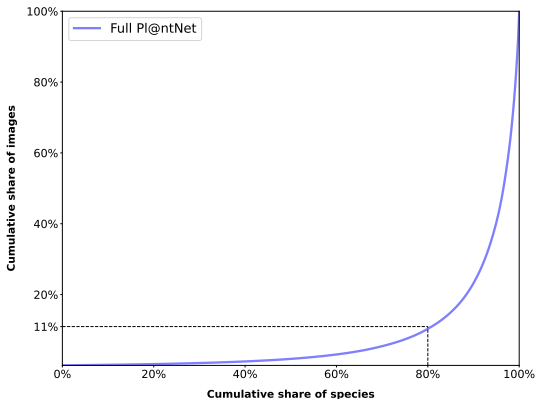
Dataset construction



**Sample at genus level to preserve intra-genus ambiguity
(use hierarchical structure)**

LONG TAILED DISTRIBUTION

PRESERVED WITH SUBSAMPLING OF GENERA



80% of species account for only 11% of images

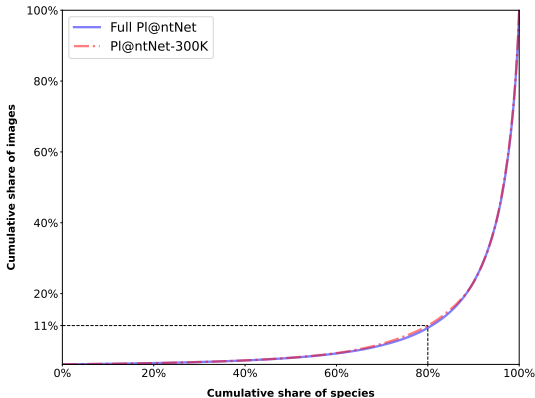


20% of species account for 89% of images

Reminder: total = 45 000 plant species (out of 300 000)

LONG TAILED DISTRIBUTION

PRESERVED WITH SUBSAMPLING OF GENERA



80% of species account for only 11% of images

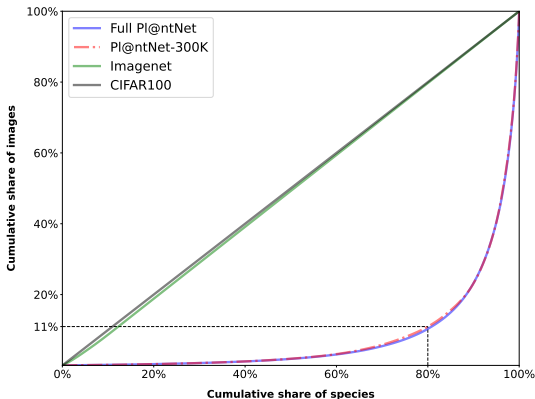


20% of species account for 89% of images

Reminder: total = 45 000 plant species (out of 300 000)

LONG TAILED DISTRIBUTION

PRESERVED WITH SUBSAMPLING OF GENERA



80% of species account for only 11% of images



20% of species account for 89% of images

Reminder: total = 45 000 plant species (out of 300 000)



- ▶ 306 146 color images
- ▶ 32 GB
- ▶ Labels: $K = 1\,081$ species
- ▶ 2 079 003 volunteers "workers"

Zenodo, 1 click download

<https://zenodo.org/record/5645731>

Code to train models:

<https://github.com/plantnet/PlantNet-300K>



Image labeling difficulty could have a huge impact on learning:

- ▶ **Removing** very difficult tasks could be useful
 - for dataset **inspection/visualization**
 - to **clean** a dataset
 - for **training performance**⁽²⁾

Hint: usually, such tasks are associated with mislabeling

- ▶ Next step:
We have seen how to assert how good is a worker, but how can we assert the labeling difficulty of an image?

⁽²⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". *NeurIPS*.



⁽³⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

⁽⁴⁾ (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.

⁽⁵⁾ Y. LeCun et al. (1998). "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* 86.11, pp. 2278–2324.



... but labeling errors are common

CIFAR10⁽³⁾



$y^* = \text{cat}$

Quickdraw⁽⁴⁾



$y^* = \text{T-shirt}$

MNIST⁽⁵⁾



$y^* = 6$

⁽³⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

⁽⁴⁾ (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.

⁽⁵⁾ Y. LeCun et al. (1998). "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* 86.11, pp. 2278–2324.

Assuming a single hard label (standard supervised settings):

- Classify data points quality with a curated set of probes⁽⁶⁾
- Confident learning⁽⁷⁾: estimate joint distribution between noisy (given) and true labels (unknown)
- Self learning⁽⁸⁾: train a model + extract features and similarity metric on a subset + retrain with modified weighted loss
- Representative Sampling (CleanNet⁽⁹⁾): trapping set + encoders + task similarity with constraints on loss
- Our focus here: study the learning dynamic,
 - ▶ **AUM**⁽¹⁰⁾ (Area Under the Margin): study margin during training

⁽⁶⁾ S. A. Siddiqui et al. (2022). *Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics*.

⁽⁷⁾ C. Northcutt, L. Jiang, and I. Chuang (2021). "Confident learning: Estimating uncertainty in dataset labels". *J. Artif. Intell. Res.* 70, pp. 1373–1411.

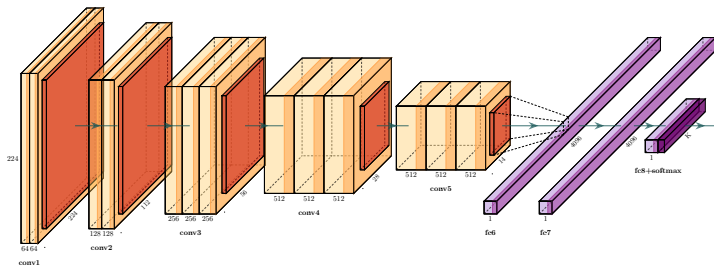
⁽⁸⁾ J. Han, P. Luo, and X. Wang (2019). "Deep self-learning from noisy labels". *ICCV*, pp. 5138–5147.

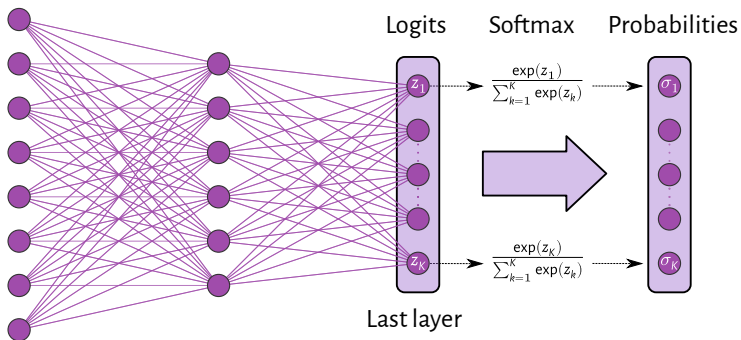
⁽⁹⁾ K.-H. Lee et al. (2018). "Cleannet: Transfer learning for scalable image classifier training with label noise". *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5447–5456.

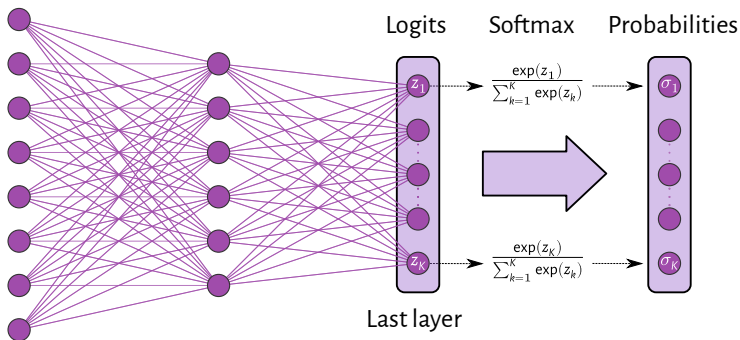
⁽¹⁰⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". *NeurIPS*.

DEEP LEARNING

NOTATION MOSTLY



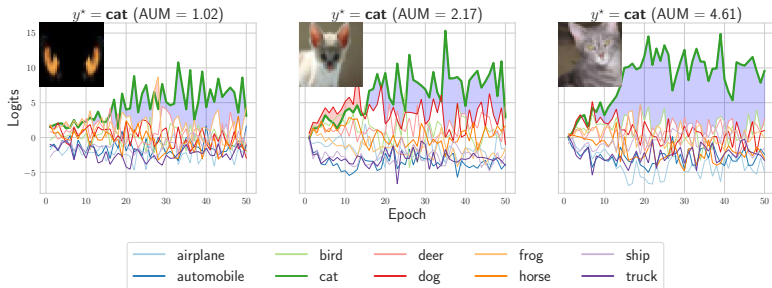




- ▶ From an image, get a score vector $z = (z_1, \dots, z_K)^T \in \mathbb{R}^K$
- ▶ z_k : **score** (logit) for class k
- ▶ σ_k : **probability** (softmax) for class k
- ▶ Train for T epochs (say with SGD)

AREA UNDER THE MARGINS⁽¹¹⁾

A STEP BACK WITH ONE LABEL PER TASK



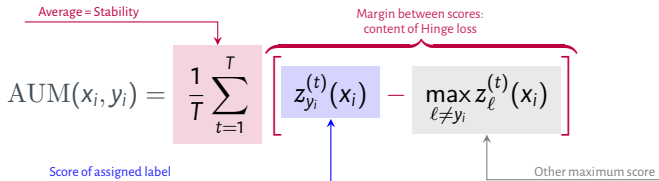
For each image

- ▶ its difficulty is reflected by how quickly the network can learn to discriminate its class
- ▶ average the difference between the "true" logit value and the one associated with the most likely one along epochs

⁽¹¹⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". *NeurIPS*.

Settings:

- ▶ $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$ (images, labels) pairs
- ▶ Classifier: at epoch $t \in [T]$, $z^{(t)}(x_i) \in \mathbb{R}^K$ a vector of **scores** (logits)

$$\text{AUM}(x_i, y_i) = \frac{1}{T} \sum_{t=1}^T \left[z_{y_i}^{(t)}(x_i) - \max_{\ell \neq y_i} z_{\ell}^{(t)}(x_i) \right]$$


Average = Stability

Score of assigned label

Margin between scores:
content of Hinge loss

Other maximum score

Settings:

- ▶ $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$ (images, labels) pairs
- ▶ Classifier: at epoch $t \in [T]$, $z^{(t)}(x_i) \in \mathbb{R}^K$ a vector of **scores** (logits)

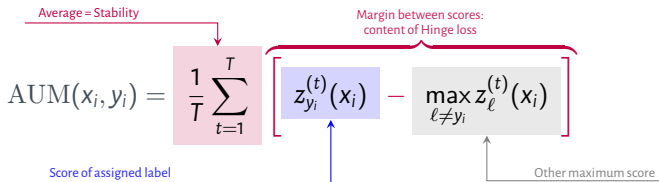
$$\text{AUM}(x_i, y_i) = \frac{1}{T} \sum_{t=1}^T \left[z_{y_i}^{(t)}(x_i) - \max_{\ell \neq y_i} z_{\ell}^{(t)}(x_i) \right]$$

Average = Stability

Margin between scores:
content of Hinge loss

Score of assigned label

Other maximum score



Challenging for crowdsourcing:

- ▶ No single y_i , multiple $y_i^{(j)}$: one for each worker w_j answering task x_i

Settings:

- ▶ $(x_i, y_i^{(j)})_{i \in [n_{\text{task}}], j \in [n_{\text{worker}}]}$: (task, labels) crowdsourced pairs
- ▶ Recall: $\mathcal{A}(x_i) := \{j \in [n_{\text{worker}}] : \text{worker } j \text{ answered task } i\}$

$$\widetilde{\text{AUM}}(x_i) = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \frac{1}{T} \sum_{t=1}^T \left[z_{y_i^{(j)}}^{(t)}(x_i) - \max_{\ell \neq y_i^{(j)}} z_{\ell}^{(t)}(x_i) \right]$$

Averaging workers AUM
Margin between scores: content of Hinge loss

Score of assigned label by worker w_j
Other maximum score

- Multiple answers \implies average each AUM (independently)

Settings:

- ▶ $(x_i, y_i^{(j)})_{i \in [n_{\text{task}}], j \in [n_{\text{worker}}]}$: (task, labels) crowdsourced pairs
- ▶ Recall: $\mathcal{A}(x_i) := \{j \in [n_{\text{worker}}] : \text{worker } j \text{ answered task } i\}$

$$\widetilde{\text{AUM}}(x_i) = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \frac{1}{T} \sum_{t=1}^T \left[z_{y_i^{(j)}}^{(t)}(x_i) - \max_{\ell \neq y_i^{(j)}} z_{\ell}^{(t)}(x_i) \right]$$

Averaging workers AUM
Margin between scores: content of Hinge loss

Score of assigned label by worker w_j
Other maximum score

- Multiple answers \implies average each AUM (independently)

Reliability issue:

- Not all workers are equally gifted \implies **weight** AUM per worker

DISSECTING THE AUM

TOWARD A CROWDSOURCED EXTENSION



- Introduce weights $s^{(j)}(x_i)$ as the trust score in worker j for task x_i

Weighted average of AUM

$$\widetilde{\text{AUM}}(x_i) = \frac{1}{S} \sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \frac{1}{T} \sum_{t=1}^T \left[z_{y_i^{(j)}}^{(t)}(x_i) - \max_{\ell \neq y_i^{(j)}} z_{\ell}^{(t)}(x_i) \right]$$

Trust score of w_j for x_i

Score of assigned label by worker w_j

Margin between scores: content of Hinge loss

Other maximum score

with $S = \sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i)$ (normalization factor)



Modifying the margin:

- Better margin (in theory, for top- k classification⁽¹²⁾)

⁽¹²⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top- k error: Analysis and insights". *CVPR*, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top- k surrogate losses". *ICML*, pp. 10727–10735.

⁽¹³⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". *J. Appl. Stat.* 45:15, pp. 2800–2818.



Modifying the margin:

- Better margin (in theory, for top- k classification⁽¹²⁾)

Change logit to softmax scores:

- avoid scale effects for scores and huge variation with multiple labels⁽¹³⁾

⁽¹²⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top- k error: Analysis and insights". *CVPR*, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top- k surrogate losses". *ICML*, pp. 10727–10735.

⁽¹³⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". *J. Appl. Stat.* 45.15, pp. 2800–2818.



Modifying the margin:

- Better margin (in theory, for top- k classification⁽¹²⁾)

Change logit to softmax scores:

- avoid scale effects for scores and huge variation with multiple labels⁽¹³⁾

Notation:

- $\sigma(x_i) = \text{softmax}(z(x_i))$ (in simplex)
- Softmax ordered: $\sigma_{[1]}(x_i) \geq \dots \geq \sigma_{[K]}(x_i) > 0$

⁽¹²⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top-k error: Analysis and insights". *CVPR*, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top-k surrogate losses". *ICML*, pp. 10727–10735.

⁽¹³⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". *J. Appl. Stat.* 45:15, pp. 2800–2818.

Modifying the margin:

- Better margin (in theory, for top- k classification⁽¹²⁾)

Change logit to softmax scores:

- avoid scale effects for scores and huge variation with multiple labels⁽¹³⁾

Notation:

- $\sigma(x_i) = \text{softmax}(z(x_i))$ (in simplex)
- Softmax ordered: $\sigma_{[1]}(x_i) \geq \dots \geq \sigma_{[K]}(x_i) > 0$

$$\text{WAUM}(x_i) := \frac{1}{S} \sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$$

Diagram illustrating the WAUM formula with annotations:

- Weighted average of AUM:** Points to the fraction $\frac{1}{S} \sum_{j \in \mathcal{A}(x_i)}$.
- Trust score of w_j for x_i :** Points to $s^{(j)}(x_i)$.
- Probability of assigned label by worker w_j :** Points to $\sigma_{y_i^{(j)}}^{(t)}(x_i)$.
- Margin between scores: content of Hinge loss:** Points to the difference $\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i)$.
- 2nd max. probability:** Points to $\sigma_{[2]}^{(t)}(x_i)$.

⁽¹²⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top-k error: Analysis and insights". CVPR, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top-k surrogate losses". ICML, pp. 10727–10735.

⁽¹³⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". J. Appl. Stat. 45.15, pp. 2800–2818.



Choosing $s^{(j)}(x_i)$:

- if $s^{(j)}(x_i) = 1$ all workers have the same weight
- if $s^{(j)}(x_i) = c_j$ the weights only depend on the worker
- DS⁽¹⁴⁾ algorithm, etc.

⁽¹⁴⁾ A. Dawid and A. Skene (1979). "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm". *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.

⁽¹⁵⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

Choosing $s^{(j)}(x_i)$:

- if $s^{(j)}(x_i) = 1$ all workers have the same weight
- if $s^{(j)}(x_i) = c_j$ the weights only depend on the worker
- DS⁽¹⁴⁾ algorithm, etc.

Our chosen worker/task score:

- Score of the form: "worker term \times task term" (similar to GLAD⁽¹⁵⁾)
- Estimate ability thanks to confusion matrices $\hat{\pi}^{(j)}$ (with DS)
- Use softmax scores to measure label confidence

⁽¹⁴⁾ A. Dawid and A. Skene (1979). "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm". *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.

⁽¹⁵⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.


Choosing $s^{(j)}(x_i)$:

- if $s^{(j)}(x_i) = 1$ all workers have the same weight
- if $s^{(j)}(x_i) = c_j$ the weights only depend on the worker
- DS⁽¹⁴⁾ algorithm, etc.

Our chosen worker/task score:

- Score of the form: "worker term \times task term" (similar to GLAD⁽¹⁵⁾)
- Estimate ability thanks to confusion matrices $\hat{\pi}^{(j)}$ (with DS)
- Use softmax scores to measure label confidence

$$s^{(j)}(x_i) = \left\langle \text{diag}(\hat{\pi}^{(j)}) \mid \sigma^{(T)}(x_i) \right\rangle \in [0, 1]$$



⁽¹⁴⁾ A. Dawid and A. Skene (1979). "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm". *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.

⁽¹⁵⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$
- For each task compute $\text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$

Usage (for learning):

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$
- For each task compute $\text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$

Usage (for learning):

- **Prune** x_i 's with $\text{WAUM}(x_i)$ below quantile q_α (say $\alpha = 0.1$)

- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$
- For each task compute $\text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$

Usage (for learning):

- **Prune** x_i 's with $\text{WAUM}(x_i)$ below quantile q_α (say $\alpha = 0.1$)
- Estimate **confusion matrices** $\hat{\pi}^{(j)}$ on pruned training dataset



- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$
- For each task compute $\text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$

Usage (for learning):

- **Prune** x_i 's with $\text{WAUM}(x_i)$ below quantile q_α (say $\alpha = 0.1$)
- Estimate **confusion matrices** $\hat{\pi}^{(j)}$ on pruned training dataset
- Get **soft labels**: normalize $\hat{y}_i = \left(\sum_{j \in \mathcal{A}(x_i)} \hat{\pi}_{k,k}^{(j)} \mathbb{1}_{\{y_i^{(j)}=k\}} \right)_{k \in [K]} \in \mathbb{R}^K$



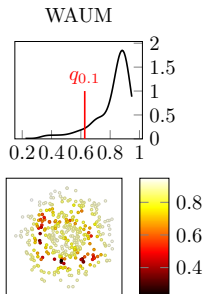
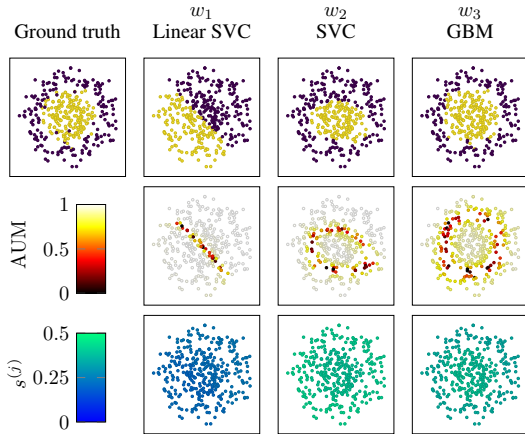
- Estimate confusion matrices $\hat{\pi}^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
- Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
- Compute $\text{AUM}(x_i, y_i^{(j)}) = \frac{1}{T} \sum_{t=1}^T \left[\sigma_{y_i^{(j)}}^{(t)}(x_i) - \sigma_{[2]}^{(t)}(x_i) \right]$
- Compute trust scores $s^{(j)}(x_i)$
- For each task compute $\text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$

Usage (for learning):

- **Prune** x_i 's with $\text{WAUM}(x_i)$ below quantile q_α (say $\alpha = 0.1$)
- Estimate **confusion matrices** $\hat{\pi}^{(j)}$ on pruned training dataset
- Get **soft labels**: normalize $\hat{y}_i = \left(\sum_{j \in \mathcal{A}(x_i)} \hat{\pi}_{k,k}^{(j)} \mathbb{1}_{\{y_i^{(j)}=k\}} \right)_{k \in [K]} \in \mathbb{R}^K$
- **Train** a classifier on the pruned dataset (with soft labels)

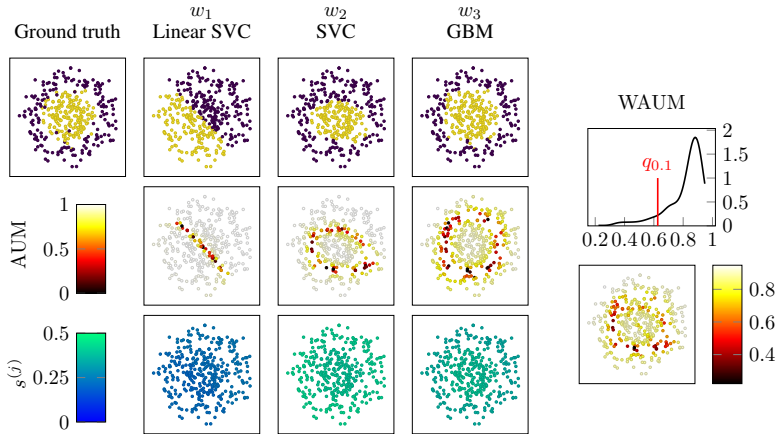
SIMULATION WITH CIRCLES

BINARY SETTING



SIMULATION WITH CIRCLES

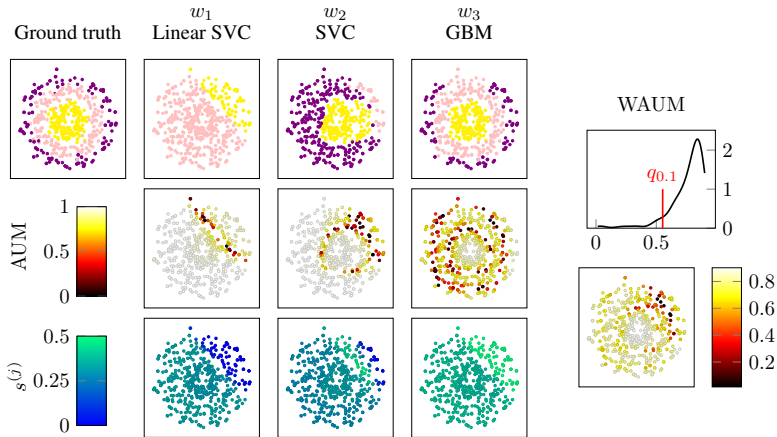
BINARY SETTING



- Workers = simulated classifiers (answering 500 tasks)
- Normalized trust scores
- Neural Network: 3-dense layers' artificial neural network (30, 20, 20)

SIMULATION WITH CIRCLES

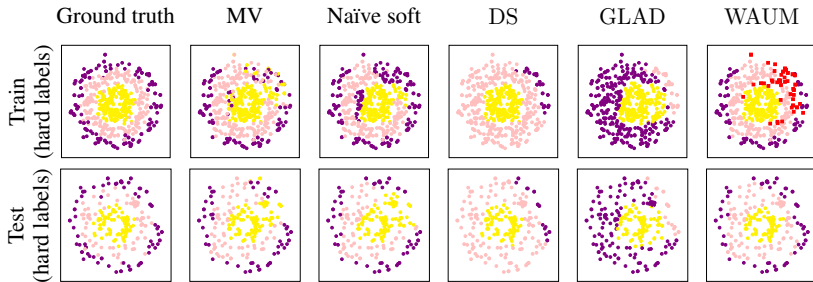
THREE CLASSES



- 3 classes with 250 tasks per class
- Normalized trust scores
- Neural Network: 3-dense layers' artificial neural network (30, 20, 20)

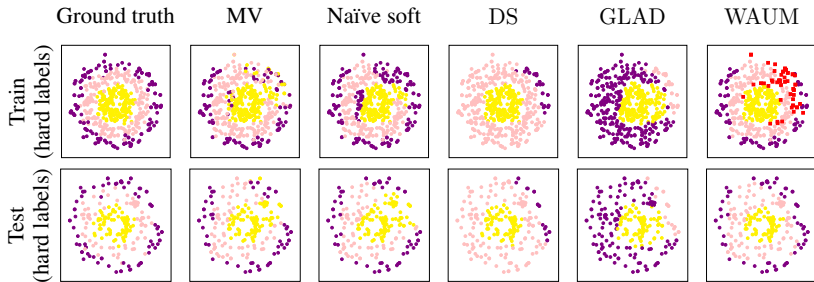
HOW CAN WE USE THE WAUM?

PRUNING TO AVOID LEARNING OF TOO AMBIGUOUS DATA



HOW CAN WE USE THE WAUM?

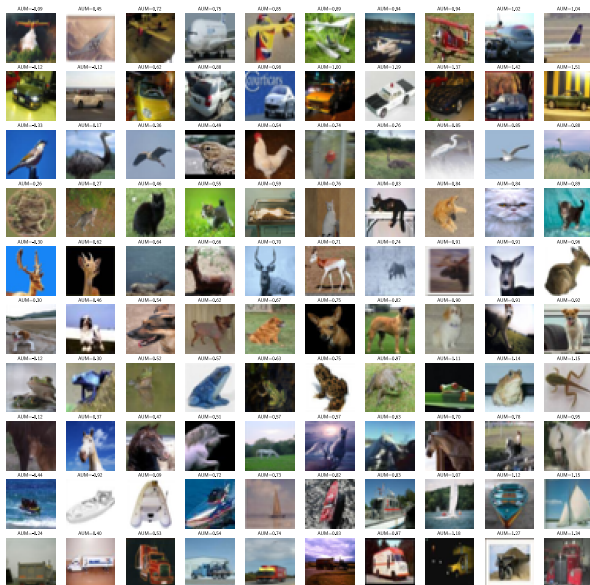
PRUNING TO AVOID LEARNING OF TOO AMBIGUOUS DATA



	MV	Naive soft	DS	GLAD	WAUM($\alpha = 0.1$)
Test accuracy	0.727	0.697	0.753	0.578	0.806

RESULTS ON CIFAR10H

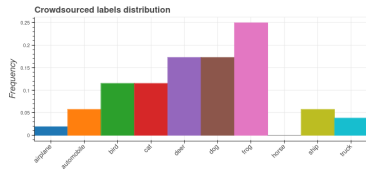
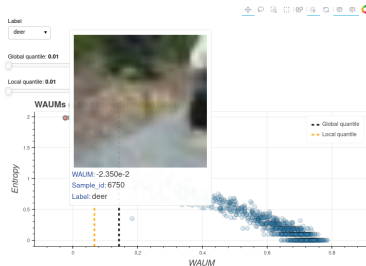
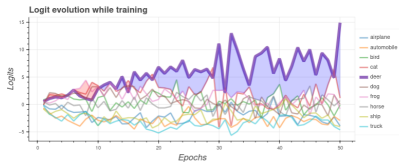
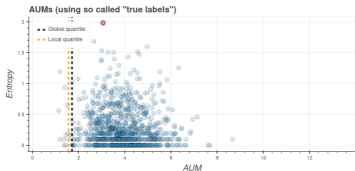
IMPROVED MISLABELED DETECTIONS: WORST AUM/WAUM



Bokeh application of the AUM/WAUM to the CIFAR10H dataset.
(see horse, cat and deer for instance)

CIFAR10H AUMs and WAUMs

AUM margin
Pleiss ▾





Generalization performance and calibration error (with a Resnet-18):

Aggregation method	Test accuracy (on CIFAR10-train)	ECE (expected calibration error)
MV	69.533 ± 0.84	0.175 ± 0.01
Naive soft	72.149 ± 2.74	0.132 ± 0.03
DS (vanilla)	70.268 ± 0.93	0.173 ± 0.01
DS (spam identification)	70.053 ± 0.81	0.174 ± 0.01
GLAD	66.569 ± 8.48	0.173 ± 0.01
WAUM	72.747 ± 1.93	0.124 ± 0.01

Remark: ECE⁽¹⁶⁾ Expected Calibration Error, the smaller the better

⁽¹⁶⁾ C. Guo et al. (2017). "On calibration of modern neural networks". *ICML*, p. 1321.

"CAN I USE THE WAUM IN MY FRAMEWORK?"

ABLATION STUDY (LABELME)



Aggregation method	Test Accuracy	ECE
WDS	85.6	0.162
WAUM + WDS	87.1	0.129

⁽¹⁷⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

⁽¹⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". *AAAI*, pp. 5832–5840.

"CAN I USE THE WAUM IN MY FRAMEWORK?"

ABLATION STUDY (LABELME)



Aggregation method	Test Accuracy	ECE
WDS	85.6	0.162
WAUM + WDS	87.1	0.129
GLAD ⁽¹⁷⁾	87.1	0.119
WAUM + GLAD	87.6	0.123

⁽¹⁷⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

⁽¹⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". *AAAI*, pp. 5832–5840.

"CAN I USE THE WAUM IN MY FRAMEWORK?"

ABLATION STUDY (LABELME)



Aggregation method	Test Accuracy	ECE
WDS	85.6	0.162
WAUM + WDS	87.1	0.129
GLAD ⁽¹⁷⁾	87.1	0.119
WAUM + GLAD	87.6	0.123
CoNAL ⁽¹⁸⁾ (lambda=0)	88.1	0.119
WAUM + CoNAL (lambda=0)	89.2	0.108

⁽¹⁷⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

⁽¹⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". *AAAI*, pp. 5832–5840.

"CAN I USE THE WAUM IN MY FRAMEWORK?"

ABLATION STUDY (LABELME)



Aggregation method	Test Accuracy	ECE
WDS	85.6	0.162
WAUM + WDS	87.1	0.129
GLAD ⁽¹⁷⁾	87.1	0.119
WAUM + GLAD	87.6	0.123
CoNAL ⁽¹⁸⁾ (lambda=0)	88.1	0.119
WAUM + CoNAL(lambda=0)	89.2	0.108
CoNAL(lambda=1e-4)	86.2	0.135
WAUM + CoNAL(lambda=1e-4)	90.0	0.099

⁽¹⁷⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

⁽¹⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". *AAAI*, pp. 5832–5840.

"CAN I USE THE WAUM IN MY FRAMEWORK?"

ABLATION STUDY (MUSIC DATASET)



Aggregation method	Test Accuracy	ECE
WDS	60.2	0.348
WAUM + WDS	63.1	0.377
GLAD ⁽¹⁷⁾	61.5	0.361
WAUM + GLAD	61.5	0.355
CoNAL ⁽¹⁸⁾ (lambda=0)	64.2	0.340
WAUM + CoNAL(lambda=0)	64.5	0.265
CoNAL(lambda=1e-4)	64.2	0.361
WAUM + CoNAL(lambda=1e-4)	64.4	0.274

⁽¹⁷⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.

⁽¹⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". *AAAI*, pp. 5832–5840.



Take home message(s)

- Citizen science challenges: many and varied (need more attention)



Take home message(s)

- Citizen science challenges: many and varied (need more attention)
- Crowdsourcing / Label uncertainty: helpful for **data curation**



Take home message(s)

- Citizen science challenges: many and varied (need more attention)
- Crowdsourcing / Label uncertainty: helpful for **data curation**
- Improved **data quality** \Rightarrow **improved learning** performance



Take home message(s)

- Citizen science challenges: many and varied (need more attention)
- Crowdsourcing / Label uncertainty: helpful for **data curation**
- Improved **data quality** \Rightarrow **improved learning** performance
- Toolbox: <https://peerannot.github.io/>
- Some benchmarks: <https://benchopt.github.io/>

Take home message(s)

- Citizen science challenges: many and varied (need more attention)
- Crowdsourcing / Label uncertainty: helpful for **data curation**
- Improved **data quality** \Rightarrow **improved learning** performance
- Toolbox: <https://peerannot.github.io/>
- Some benchmarks: <https://benchopt.github.io/>

Future work

- ▶ Release a Pl@ntnet crowdsourced dataset (**2M workers**)
- ▶ Leverage gamification for more quality labels thep1antgame.com

Contact:

Joseph Salmon

✉ joseph.salmon@umontpellier.fr








🌐 <https://josephsalmon.eu>








Github: @josephsalmon





Mastodon: @josephsalmon@sigmoid.social



-  (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.
-  Chu, Z., J. Ma, and H. Wang (2021). “Learning from Crowds by Modeling Common Confusions.”. *AAAI*, pp. 5832–5840.
-  Dawid, A. and A. Skene (1979). “Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm”. *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.
-  Garcin, C. et al. (2021). “Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution”. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.
-  Guo, C. et al. (2017). “On calibration of modern neural networks”. *ICML*, p. 1321.
-  Han, J., P. Luo, and X. Wang (2019). “Deep self-learning from noisy labels”. *ICCV*, pp. 5138–5147.
-  Ju, C., A. Bibaut, and M. van der Laan (2018). “The relative performance of ensemble methods with deep convolutional neural networks for image classification”. *J. Appl. Stat.* 45.15, pp. 2800–2818.

-  Krizhevsky, A. and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.
-  Lapin, M., M. Hein, and B. Schiele (2016). “Loss functions for top-k error: Analysis and insights”. *CVPR*, pp. 1468–1477.
-  LeCun, Y. et al. (1998). “Gradient-based learning applied to document recognition”. *Proceedings of the IEEE* 86.11, pp. 2278–2324.
-  Lee, K.-H. et al. (2018). “Cleannet: Transfer learning for scalable image classifier training with label noise”. *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5447–5456.
-  Northcutt, C., L. Jiang, and I. Chuang (2021). “Confident learning: Estimating uncertainty in dataset labels”. *J. Artif. Intell. Res.* 70, pp. 1373–1411.
-  Pleiss, G. et al. (2020). “Identifying mislabeled data using the area under the margin ranking”. *NeurIPS*.
-  Siddiqui, S. A. et al. (2022). *Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics*.



-  Whitehill, J. et al. (2009). “Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise”. *NeurIPS*. Vol. 22.
-  Yang, F. and S. Koyejo (2020). “On the consistency of top-k surrogate losses”. *ICML*, pp. 10727–10735.



- DS assumption: errors only come from workers (no task modeling)

⁽¹⁹⁾J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". *NeurIPS*. vol. 22.



- DS assumption: errors only come from workers (no task modeling)

GLAD: incorporating task difficulty

Model labeling errors as a function of worker ability and task difficulty:

- ▶ worker j has an ability $\alpha_j \in \mathbb{R}$
- ▶ task i has a difficulty $\beta_i \in \mathbb{R}_+^*$

$$\mathbb{P}(y_i^{(j)} = y_i^* | \alpha_j, \beta_i) = \frac{1}{1 + e^{-\alpha_j \beta_i}}$$

Note: assume uniform errors on other labels



For $x \in \mathcal{X}_{\text{train}} = \{x_1, \dots, x_{n_{\text{task}}}\}$, let $\sigma(x) \in \Delta_{K-1}$ (softmax output)

Split $[0, 1]$ into $M (= 15)$ bins l_1, \dots, l_M of size $\frac{1}{M}$: $l_m = (\frac{m-1}{M}, \frac{m}{M}]$, for $m \in [M]$

Denote $B_m = \{x \in \mathcal{X}_{\text{train}} : \sigma_{[1]}(x) \in l_m\}$ the tasks whose predicted probabilities are in the m -th bin

Define **accuracy** and **confidence**:

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbb{1}_{\{\sigma_{[1]}(x_i) = y_i\}} \quad \text{and} \quad \text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \sigma_{[1]}(x_i) .$$

Then, the Expected Calibration Error (ECE) reads:

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n_{\text{task}}} |\text{acc}(B_m) - \text{conf}(B_m)| .$$

Perfect calibration : $\text{ECE} = 0$ (accuracy = confidence for each subset B_m)