

IMPROVE LEARNING COMBINING CROWDSOURCED LABELS BY 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)

*Improve learning 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:



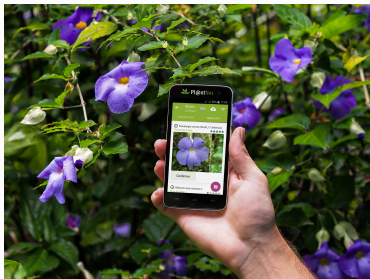
Pierre Bonnet (CIRAD, AMAP)

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

- ▶ C. Garcin, A. Joly, et al. (2021). “Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution”. In: *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*
- ▶ C. Garcin, M. Servajean, et al. (2022). “Stochastic smoothing of the top-K calibrated hinge loss for deep imbalanced classification”. In: *ICML*

PLANT CLASSIFICATION WITH PL@NTNET

<https://plantnet.org/>



- ▶ ML assisted **citizen science**
- ▶ > 40,000 species
- ▶ > 10,000,000 annotated images
- ▶ > 1Tb of data \implies Reduction to share with community

← Identification



Résultats



Dipsacus fullonum L.
Cabaret-des-oiseaux

Caprifoliaceae **i**

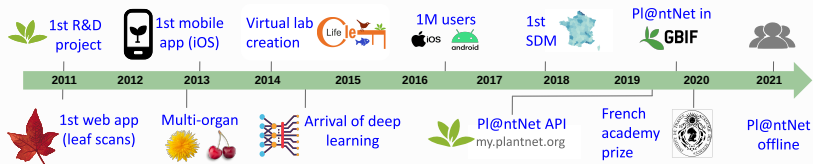
Valider 4.89 ★★★★★



Cichorium intybus L.
Chicorée amère

Asteraceae **i**

Pl@ntNet Key milestones



 Inria

 cirad

 IRD

 INRAE

 agropolis fondation



Introduction

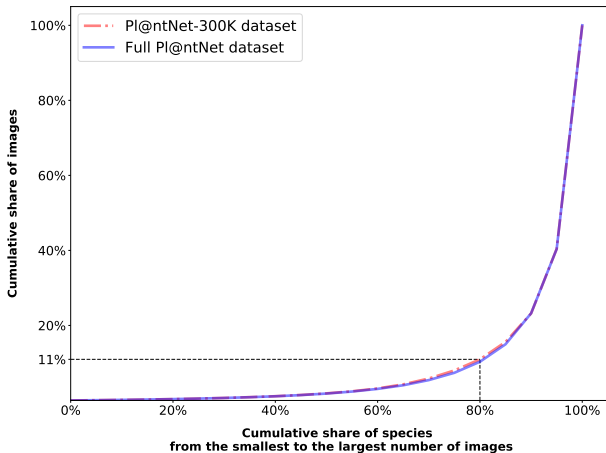
Pl@ntNet-300K

Dataset characteristics

Dataset construction

LONG TAILED DISTRIBUTION

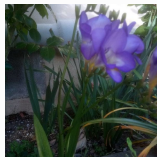
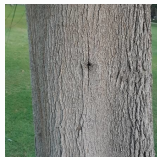
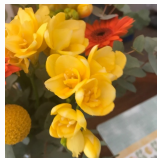
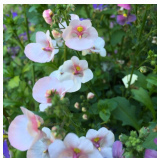
PRESERVED WITH SAMPLING OF GENERA



80% of species account for only 11% of images

INTRA-CLASS VARIABILITY

SAME LABEL/SPECIES BUT VERY DIVERSE IMAGES



*Guizotia
abyssinica*

*Diascia
rigescens*

*Lapageria
rosea*

*Casuarina
cunninghamiana*

*Freesia
alba*

Plant species are challenging to model based on pictures only!

INTER-CLASS AMBIGUITY

DIFFERENT LABELS/SPECIES BUT SIMILAR IMAGES



Cirsium rivulare



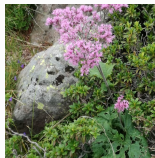
Chaerophyllum aromaticum



Conostomium kenysense



Adenostyles leucophylla



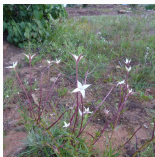
Sedum montanum



Cirsium tuberosum



Chaerophyllum temulum



Conostomium quadrangulare



Adenostyles alliariae



Sedum rupestre

Some species are visually similar (especially within genus)



Introduction

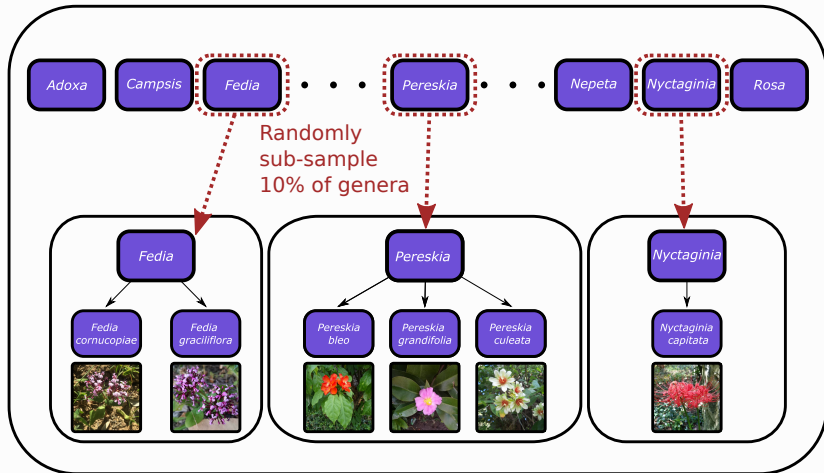
Pl@ntNet-300K

Dataset characteristics

Dataset construction

CONSTRUCTION OF PL@NTNET-300K

SUBSAMPLING OF GENERA



Sample at genus level to preserve intra-genus ambiguity



- ▶ **306,146** color images images
- ▶ Labels: **1,081** species
- ▶ **2,079,003** workers (volunteers), with ≈ 2 labels per worker (on average)

Zenodo, 1 click download

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

Code to train models:

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

PROBLEM: CAN WE TRUST OUR DATA?



⁽¹⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

⁽²⁾ (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.

⁽³⁾ Y. LeCun et al. (1998). "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11, pp. 2278–2324.

PROBLEM: CAN WE TRUST OUR DATA?



... but labelling errors are common

CIFAR10⁽¹⁾



$y^* = \text{cat}$

Quickdraw⁽²⁾



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

MNIST⁽³⁾



$y^* = 6$

⁽¹⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

⁽²⁾ (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.

⁽³⁾ Y. LeCun et al. (1998). "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE* 86.11, pp. 2278–2324.



- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.



- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.

Questions:



- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.

Questions:

- ▶ Where do the tasks come from?



- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.

Questions:

- ▶ Where do the tasks come from? **Web scrapping**



- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.

Questions:

- ▶ Where do the tasks come from? **Web scrapping**
- ▶ Where do the labels come from?



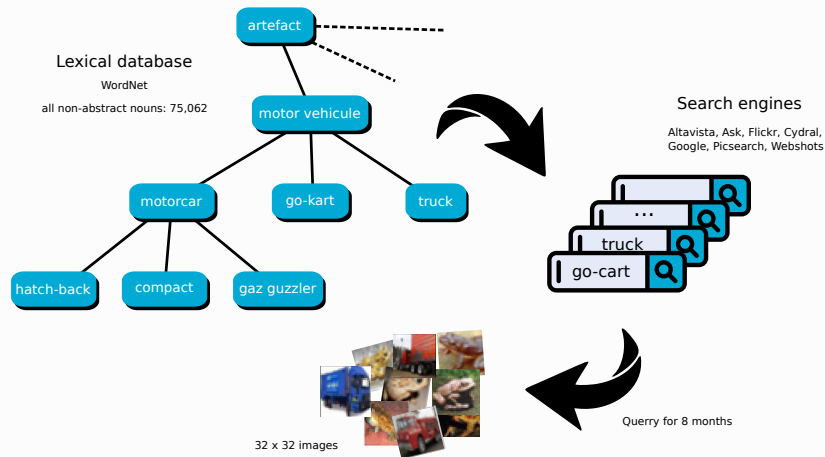
- ▶ Classical dataset: $(x_1, y_1), \dots, (x_{n_{\text{task}}}, y_{n_{\text{task}}})$
Features/tasks \times labels pairs: $(x_i, y_i) \in \mathcal{X} \times [K] = \{1, \dots, K\}$
- ▶ Popular datasets used for supervised learning (classification):
CIFAR10, CIFAR100, ImageNet, MNIST, Quickdraw, LabelMe, etc.

Questions:

- ▶ Where do the tasks come from? **Web scrapping**
- ▶ Where do the labels come from? **Crowdsourcing**

CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)



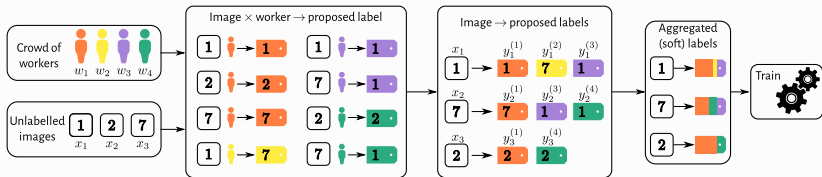
80 Million Tiny Images

Note: some issues on this process⁽⁴⁾

⁽⁴⁾ V. Uday Prabhu and A. Birhane (June 2020). "Large image datasets: A pyrrhic win for computer vision?" In: *arXiv e-prints*, arXiv:2006.16923, arXiv:2006.16923.

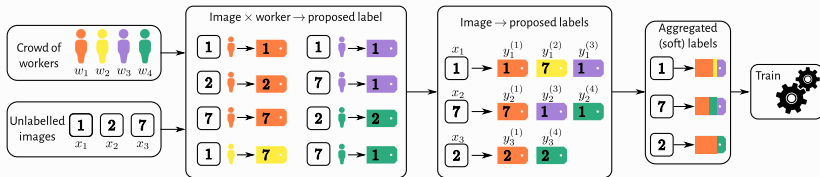
CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING



CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING

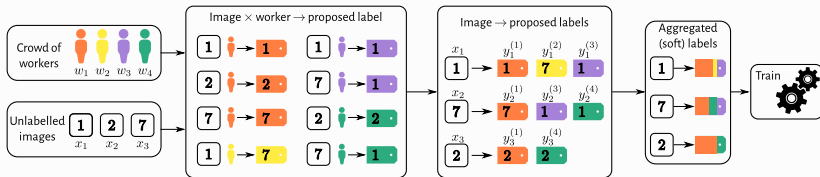


Quotes⁽⁵⁾ :

⁽⁵⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING



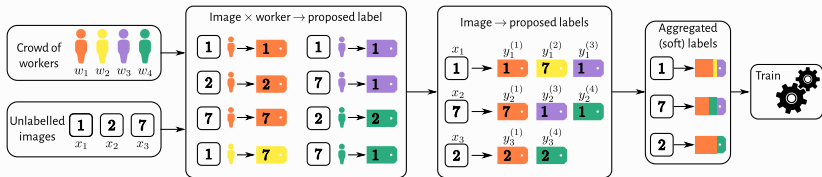
Quotes⁽⁵⁾ :

- ▶ "We paid students to label a subset of the tiny images dataset[...]. The labelers were paid a fixed sum per hour spent labeling."

⁽⁵⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING



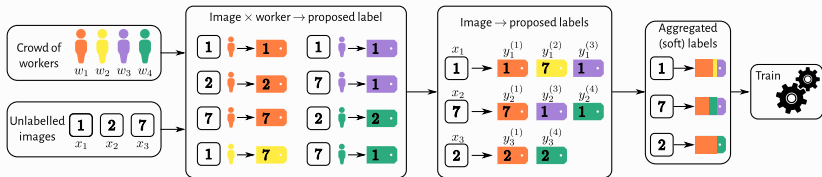
Quotes⁽⁵⁾:

- ▶ "We paid students to label a subset of the tiny images dataset[...]. The labelers were paid a fixed sum per hour spent labeling."
- ▶ "Since each image in the dataset already comes with a noisy label (the search term used to find the image), all we needed the labelers to do was to filter out the mislabeled images."

⁽⁵⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING



Quotes⁽⁵⁾ :

- ▶ "We paid students to label a subset of the tiny images dataset[...]. The labelers were paid a fixed sum per hour spent labeling."
- ▶ "Since each image in the dataset already comes with a noisy label (the search term used to find the image), all we needed the labelers to do was to filter out the mislabeled images."
- ▶ "Furthermore, we personally verified every label submitted by the labelers": *errare humanum est*

⁽⁵⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.



Peterson *et al.* (2019) "Our final CIFAR10H behavioral dataset consists of **511,400** human categorization decisions over the $n_{\text{tasks}}=10,000$ -image testing subset of CIFAR10 (approx. 50 judgments per image)."

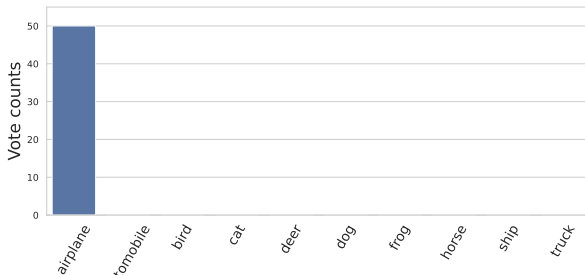
- ▶ Total number of workers: $n_{\text{worker}} = \mathbf{2,571}$ (via Amazon Mechanical Turk)
- ▶ Processing: every 20 trials, an obvious image is presented as an attention check, and participants who scored below 75% on these were removed from the final analysis (14 total).

Note: workers were paid \$1.50 total.

⁽⁶⁾. C. Peterson et al. (2019). "Human Uncertainty Makes Classification More Robust". In: ICCV, pp. 9617–9626.



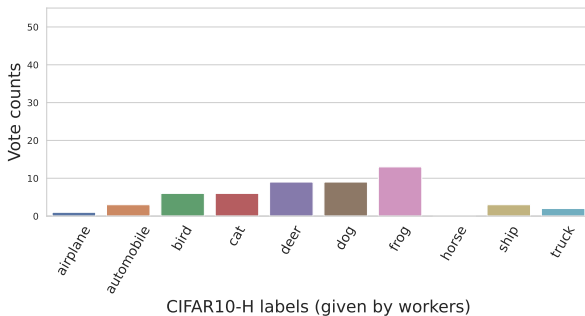
Image #7681
CIFAR10 label: airplane



CIFAR10-H labels (given by workers)



Image #6750
CIFAR10 label: deer



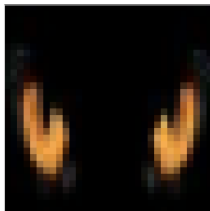
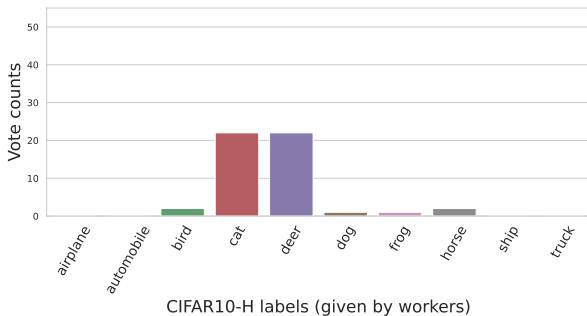


Image #9246
CIFAR10 label: cat



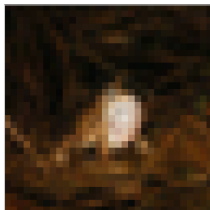


Image #3724
CIFAR10 label: frog

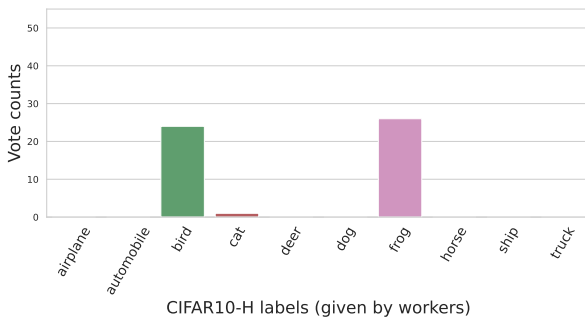




Image #1353
CIFAR10 label: cat

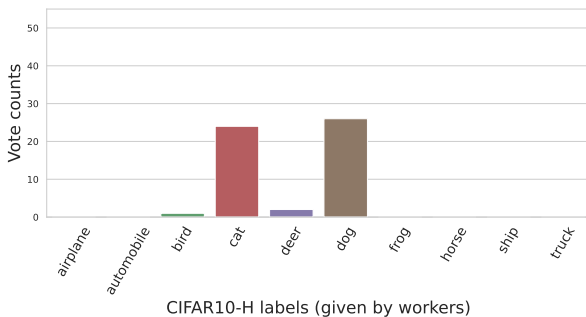
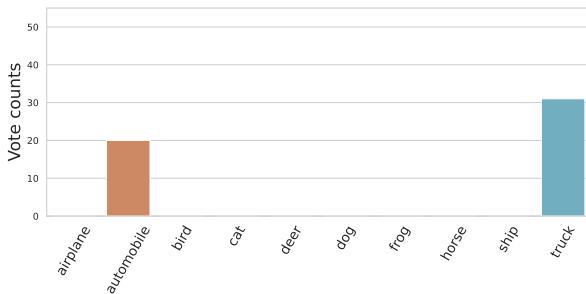




Image #7455

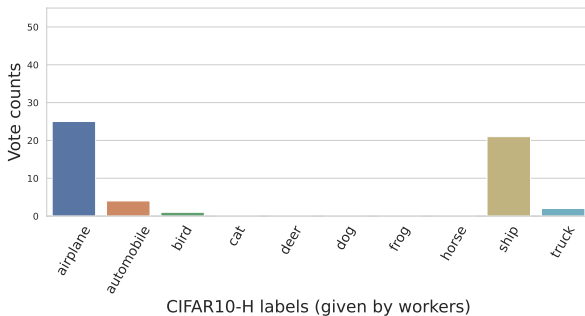
CIFAR10 label: automobile



CIFAR10-H labels (given by workers)



Image #8872
CIFAR10 label: ship





Simple strategies:

- ▶ **Majority voting (MV):**
naive but ineffective for borderline cases
- ▶ First label reaching a **consensus** of p workers (often $p = 5$)⁽⁷⁾
→ arbitrary choice of p
- ▶ Leverage label distribution, say with **entropy**:
not always reliable (*e.g.*, with few labels), biases, psychology mechanisms spammers

Intermission : see app for entropy visualization

⁽⁷⁾ R. Snow et al. (2008). "Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks". In: *Conference on Empirical Methods in Natural Language Processing*. EMNLP 2008. Association for Computational Linguistics, pp. 254–263.

A FIRST SOLUTION: CLASSIFY THE QUALITY

IMAGENET ODDITIES



- **curated set of probes**⁽⁸⁾ in the training data (OOD=Out Of Distribution)
e.g.,: ImageNet⁽⁹⁾ +14 millions tasks, $K = 1000$ classes
 $(\text{task}_i, \text{label}_i, \text{metadata}_i) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{M}$

Black bear



(a) Typical

Dishwasher



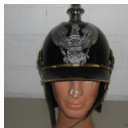
(b) Atypical

School bus



(c) Corrupted

Mud turtle



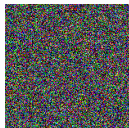
(d) Rand Label

Jeep



(e) OOD

Loafer



(f) Rand Input

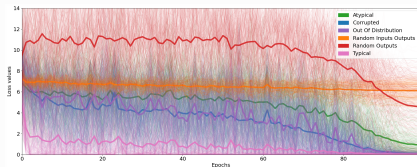
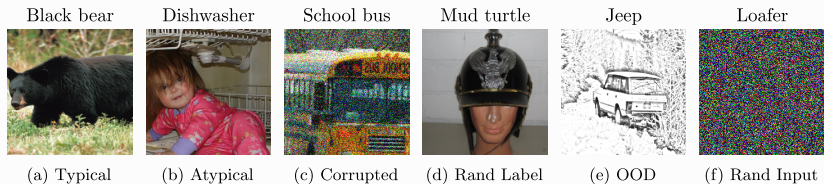
⁽⁸⁾ S. A. Siddiqui et al. (2022). *Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics*.

⁽⁹⁾ O. Russakovsky et al. (2015). "ImageNet Large Scale Visual Recognition Challenge". In: *Int. J. Comput. Vision* 115.3, pp. 211–252.

A FIRST SOLUTION: CLASSIFY THE QUALITY

IMAGENET ODDITIES

- **curated set of probes**⁽⁸⁾ in the training data (OOD=Out Of Distribution)
e.g.,: ImageNet⁽⁹⁾ +14 millions tasks, $K = 1000$ classes
 $(\text{task}_i, \text{label}_i, \text{metadata}_i) \in \mathcal{X} \times \mathcal{Y} \times \mathcal{M}$



- 1 metadata = 1 dynamic
- Identify the ambiguity

⁽⁸⁾ S. A. Siddiqui et al. (2022). *Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics*.

⁽⁹⁾ O. Russakovsky et al. (2015). "ImageNet Large Scale Visual Recognition Challenge". In: *Int. J. Comput. Vision* 115.3, pp. 211–252.



Q: When was the last time you had a curated set of metadata up your sleeve?

⁽¹⁰⁾ C. Northcutt, L. Jiang, and I. Chuang (2021). "Confident learning: Estimating uncertainty in dataset labels". In: *J. Artif. Intell. Res.* 70, pp. 1373–1411.

⁽¹¹⁾ J. Han, P. Luo, and X. Wang (2019). "Deep self-learning from noisy labels". In: *ICCV*, pp. 5138–5147.

⁽¹²⁾ K.-H. Lee et al. (2018). "Cleannet: Transfer learning for scalable image classifier training with label noise". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5447–5456.

⁽¹³⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: *NeurIPS*.



Q: When was the last time you had a curated set of metadata up your sleeve?

A: Never!

Assuming we have a hard label ($\in [K]$):

- Confident learning⁽¹⁰⁾: estimate joint distribution between noisy (given) and true labels (unknown)
- Self learning⁽¹¹⁾: train a model + extract features and similarity metric on a subset + retrain with modified weighted loss
- Representative Sampling (CleanNet⁽¹²⁾): trapping set + encoders + task similarity with constraints on loss
- Our focus here: study the learning dynamic,
 - ▶ **AUM**⁽¹³⁾ (Area Under the Margin): study margin during training

⁽¹⁰⁾ C. Northcutt, L. Jiang, and I. Chuang (2021). "Confident learning: Estimating uncertainty in dataset labels". In: *J. Artif. Intell. Res.* 70, pp. 1373–1411.

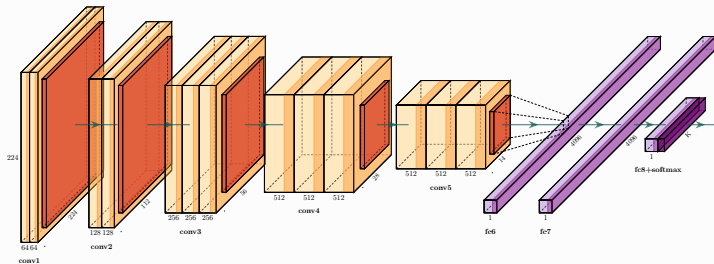
⁽¹¹⁾ J. Han, P. Luo, and X. Wang (2019). "Deep self-learning from noisy labels". In: *ICCV*, pp. 5138–5147.

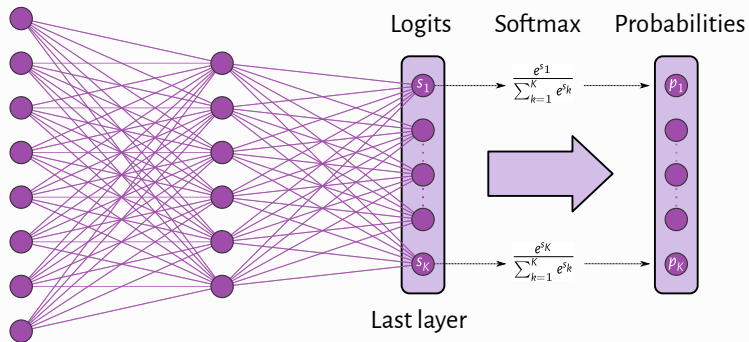
⁽¹²⁾ K.-H. Lee et al. (2018). "Cleannet: Transfer learning for scalable image classifier training with label noise". In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5447–5456.

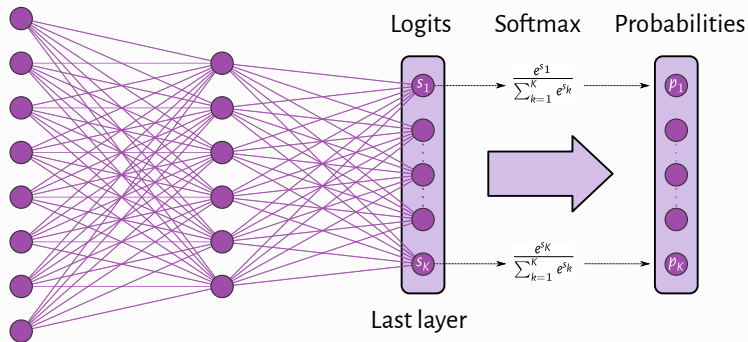
⁽¹³⁾ C. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: *NeurIPS*.

DEEP LEARNING

NOTATION MOSTLY



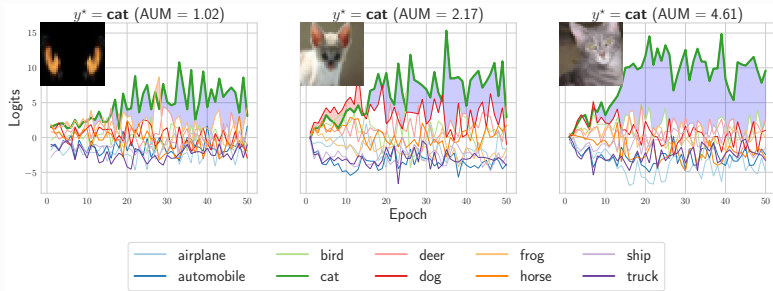




- ▶ From an image, get a score vector $s = (s_1, \dots, s_L)^T \in \mathbb{R}^L$ (aka logits)
- ▶ s_k : score for class k
- ▶ Train for T epochs (say with SGD)

AREA UNDER THE MARGINS⁽¹⁴⁾

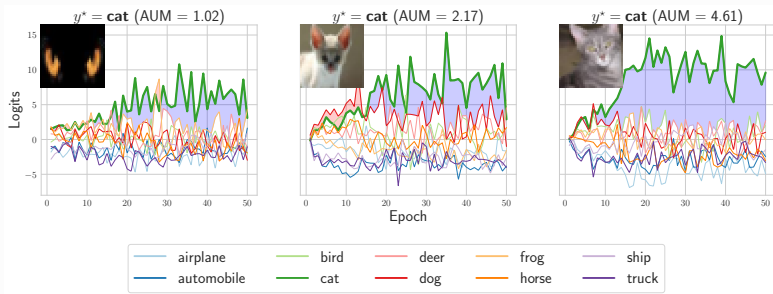
A STEP BACK WITH ONE LABEL PER TASK



⁽¹⁴⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: *NeurIPS*.

AREA UNDER THE MARGINS⁽¹⁴⁾

A STEP BACK WITH ONE LABEL PER TASK



Motivation: the logit scores (average) value along learning epochs give insights on the task difficulty

⁽¹⁴⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: *NeurIPS*.

Settings:

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

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

Average = Stability

Margin between scores: content of Hinge loss

Score of assigned label

Other maximum score

Settings:

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

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

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:

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

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

Average = Stability

Score of assigned label

Margin between scores: content of Hinge loss

Other maximum score

Challenging for crowdsourcing:

- No single y_i , multiple $y_i^{(j)}$: one for each worker w_j answering task x_i
 - ▶ ...so $\mathcal{C}^{(t)}(x_i)_{y_i}$ does not exist

Settings:

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

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

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
 - ▶ ...so $\mathcal{C}^{(t)}(x_i)_{y_i}$ does not exist
 - ▶ ...and same issue with $\ell \neq y_i$.

Settings:

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

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

Averaging workers AUM

Score of assigned label by worker w_j

Margin between scores: content of Hinge loss

Other maximum score

- Multiple answers \implies average each AUM (independently)
- Let $\mathcal{A}(x_i) := \{j \in [n_{\text{worker}}] : \text{worker } j \text{ answered task } i\}$.

Settings:

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

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

Averaging workers AUM

Score of assigned label by worker w_j

Margin between scores: content of Hinge loss

Other maximum score

- Multiple answers \implies average each AUM (independently)
- Let $\mathcal{A}(x_i) := \{j \in [n_{\text{worker}}] : \text{worker } j \text{ answered task } i\}$.

Reliability issue:

- Expert = random workers \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

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

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

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



Modifying the margin:

- Scale effects in the scores discarded, need normalization⁽¹⁵⁾
- Better margin (in theory, for top- k classification⁽¹⁶⁾)

⁽¹⁵⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". In: *J. Appl. Stat.* 45.15, pp. 2800–2818.

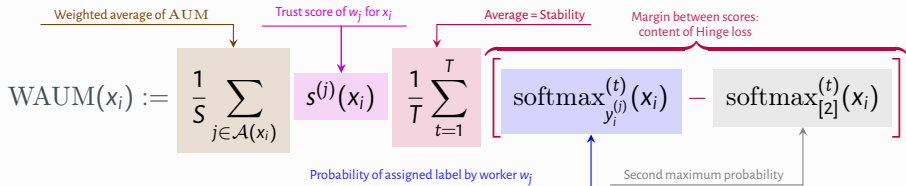
⁽¹⁶⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top- k error: Analysis and insights". In: *CVPR*, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top- k surrogate losses". In: *ICML*, pp. 10727–10735.

Modifying the margin:

- Scale effects in the scores discarded, need normalization⁽¹⁵⁾
- Better margin (in theory, for top- k classification⁽¹⁶⁾)

Notation:

- $\text{softmax}(x_i) = \text{softmax}(\mathcal{C}(x_i)) \in \Delta_{K-1}$ (simplex of dim $K - 1$)
- Softmax ordered: $\text{softmax}_{[1]}(x_i) \geq \dots \geq \text{softmax}_{[K]}(x_i) > 0$



⁽¹⁵⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". In: *J. Appl. Stat.* 45.15, pp. 2800–2818.

⁽¹⁶⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top-k error: Analysis and insights". In: *CVPR*, pp. 1468–1477; F. Yang and S. Koyejo (2020). "On the consistency of top-k surrogate losses". In: *ICML*, pp. 10727–10735.



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
- ...there is already a literature on trusting workers !

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
- ...there is already a literature on trusting workers !

DS: Dawid and Skene⁽¹⁷⁾

Assumption: each worker answers independently

j -th worker **confusion matrix:** $\pi^{(j)} \in \mathbb{R}^{K \times K}$: $\pi_{\ell,k}^{(j)} = \mathbb{P}(y_i^{(j)} = \ell | y_i^* = k)$

$$y_i^{(j)} | y_i^* = \ell \sim \text{Multinomial}(\pi_{\ell \bullet}^{(j)})$$

Note : diagonal elements of $\pi^{(j)}$ represents worker ability to be correct

⁽¹⁷⁾ A. Dawid and A. Skene (1979). "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm". In: *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.

Likelihood:

$$\prod_{k \in [K]} \pi_{\ell, k}^{(j)}$$

- 1 task, 1 worker and 1 answer conditioned on $y_i^* = \ell$

Likelihood:

$$\prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \pi_{\ell, k}^{(j)}$$

- 1 task, 1 worker and 1 answer conditioned on $y_i^* = \ell$
- Multiple workers answer independently

Likelihood:

$$\prod_{\ell \in [K]} \left[\mathbb{P}(y_i^* = \ell) \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \pi_{\ell, k}^{(j)} \right]^{\mathbb{1}_{\{y_i^* = \ell\}}}$$

- 1 task, 1 worker and 1 answer conditioned on $y_i^* = \ell$
- Multiple workers answer independently
- Remove conditioning assumption on y_i^* : $\mathbb{P}(y_i^* = \ell) = \rho_\ell$

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \pi_{\ell, k}^{(j)} \right]^{T_{i\ell}}$$

- 1 task, 1 worker and 1 answer conditioned on $y_i^* = \ell$
- Multiple workers answer independently
- Remove conditioning assumption on y_i^* : $\mathbb{P}(y_i^* = \ell) = \rho_{\ell}$
- Each task is independent: $T_{i\ell} = 1$ if task i has label ℓ and 0 otherwise

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right) \right]^{T_{i\ell}}$$

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{T_{i\ell}} \right]$$

Prevalence of class ℓ (points to ρ_{ℓ})
 Indicator of class ℓ for task i (points to $T_{i\ell}$)
 Probability for worker j to answer k with truth ℓ (points to $\pi_{\ell, k}^{(j)}$)

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{T_{i\ell}} \right]$$

Prevalence of class ℓ (points to ρ_{ℓ})
 Indicator of class ℓ for task i (points to $T_{i\ell}$)
 Probability for worker j to answer k with truth ℓ (points to $\pi_{\ell, k}^{(j)}$)

1 Soft labels initialization:

$$\forall i \in [n_{\text{task}}], \forall \ell \in [K], \hat{T}_{i\ell} = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_i^{(j)} = \ell\}}$$

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right) \right] T_{i\ell}$$

Prevalence of class ℓ (points to ρ_{ℓ})
 Indicator of class ℓ for task i (points to $T_{i\ell}$)
 Probability for worker j to answer k with truth ℓ (points to $\pi_{\ell, k}^{(j)}$)

1 **Soft labels initialization:**

$$\forall i \in [n_{\text{task}}], \forall \ell \in [K], \hat{T}_{i\ell} = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_j^{(i)} = \ell\}}$$

2 **while not converged do**

6 **Labels:** $\forall i \in [n_{\text{task}}], \hat{y}_i = \hat{T}_{i\bullet} \in \mathbb{R}^K$ (soft label)

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = \ell\}}} \right]$$

Prevalence of class ℓ (points to ρ_{ℓ})
 Indicator of class ℓ for task i (points to $\mathbb{1}_{\{y_i^{(j)} = \ell\}}$)
 Probability for worker j to answer k with truth ℓ (points to $\pi_{\ell, k}^{(j)}$)
 $T_{i\ell}$ (points to the exponent)

1 Soft labels initialization:

$$\forall i \in [n_{\text{task}}], \forall \ell \in [K], \hat{T}_{i\ell} = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_i^{(j)} = \ell\}}$$

2 while not converged do

// **M-step:** Get $\hat{\pi}$ and $\hat{\rho}$ assuming \hat{T} s are known

$$\forall (\ell, k) \in [K]^2, \hat{\pi}_{\ell k}^{(j)} \leftarrow \frac{\sum_{i \in [n_{\text{task}}]} \hat{T}_{i\ell} \mathbb{1}_{\{y_i^{(j)} = k\}}}{\sum_{k' \in [K]} \sum_{i' \in [n_{\text{task}}]} \hat{T}_{i'\ell} \mathbb{1}_{\{y_{i'}^{(j)} = k'\}}}$$

$$\forall \ell \in [K], \hat{\rho}_{\ell} \leftarrow \frac{1}{n_{\text{task}}} \sum_{i \in [n_{\text{task}}]} \hat{T}_{i\ell}$$

6 Labels: $\forall i \in [n_{\text{task}}], \hat{y}_i = \hat{T}_{i \bullet} \in \mathbb{R}^K$ (soft label)

Likelihood:

$$\prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right) \right] T_{i\ell}$$

Prevalence of class ℓ (points to ρ_{ℓ})
 Indicator of class ℓ for task i (points to $T_{i\ell}$)
 Probability for worker j to answer k with truth ℓ (points to $\pi_{\ell, k}^{(j)}$)

1 Soft labels initialization:

$$\forall i \in [n_{\text{task}}], \forall \ell \in [K], \hat{T}_{i\ell} = \frac{1}{|\mathcal{A}(x_i)|} \sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_i^{(j)} = \ell\}}$$

2 while not converged do

// **M-step:** Get $\hat{\pi}$ and $\hat{\rho}$ assuming \hat{T} s are known

$$\forall (\ell, k) \in [K]^2, \hat{\pi}_{\ell k}^{(j)} \leftarrow \frac{\sum_{i \in [n_{\text{task}}]} \hat{T}_{i\ell} \mathbb{1}_{\{y_i^{(j)} = k\}}}{\sum_{k' \in [K]} \sum_{i' \in [n_{\text{task}}]} \hat{T}_{i'\ell} \mathbb{1}_{\{y_{i'}^{(j)} = k'\}}}$$

$$\forall \ell \in [K], \hat{\rho}_{\ell} \leftarrow \frac{1}{n_{\text{task}}} \sum_{i \in [n_{\text{task}}]} \hat{T}_{i\ell}$$

// **E-step:** Estimate \hat{T} s with current $\hat{\pi}$ and $\hat{\rho}$

$$\forall i \in [n_{\text{task}}], \forall \ell \in [K], \hat{T}_{i\ell} = \frac{\prod_{j \in \mathcal{A}(x_i)} \prod_{k \in [K]} \hat{\rho}_{\ell} \cdot \hat{\pi}_{\ell, k}^{(j)}}{\sum_{\ell' \in [K]} \prod_{j' \in \mathcal{A}(x_i)} \prod_{k' \in [K]} \hat{\rho}_{\ell'} \cdot \hat{\pi}_{\ell', k'}^{(j')}}}$$

6 Labels: $\forall i \in [n_{\text{task}}], \hat{y}_i = \hat{T}_{i\bullet} \in \mathbb{R}^K$ (soft label)



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

⁽¹⁸⁾J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". In: *NeurIPS*. vol. 22.



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

GLAD: incorporating task difficulty

Model labelling 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



Proposed scores:

- Keep the product of a worker term and a task term
- Use multidimensionality of DS confusion matrices
- Use a neural network as control agent⁽¹⁹⁾

⁽¹⁹⁾ M. Servajean et al. (2017). "Crowdsourcing thousands of specialized labels: A Bayesian active training approach". In: *IEEE Transactions on Multimedia* 19.6, pp. 1376–1391.



Proposed scores:

- Keep the product of a worker term and a task term
- Use multidimensionality of DS confusion matrices
- Use a neural network as control agent⁽¹⁹⁾

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

↑
Worker j overall ability ℓ ↑
Difficulty of task i

⁽¹⁹⁾ M. Servajean et al. (2017). "Crowdsourcing thousands of specialized labels: A Bayesian active training approach". In: *IEEE Transactions on Multimedia* 19.6, pp. 1376–1391.



- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$



- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$
- For each worker j
 - ▶ Train a network on $\{(x_i, y_i^{(j)}); x_i \text{ is answered by worker } j\}$
 - ▶ Compute $AUM(x_i, y_i^{(j)})$ for the answered tasks x_i
 - ▶ Compute trust scores $s^{(j)}(x_i)$
 - ▶ For each task i compute $WAUM(x_i)$



- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$
- For each worker j
 - ▶ Train a network on $\{(x_i, y_i^{(j)}); x_i \text{ is answered by worker } j\}$
 - ▶ Compute $AUM(x_i, y_i^{(j)})$ for the answered tasks x_i
 - ▶ Compute trust scores $s^{(j)}(x_i)$
 - ▶ For each task i compute $WAUM(x_i)$



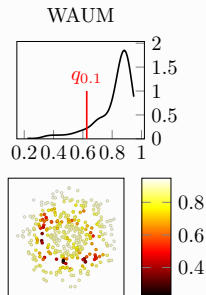
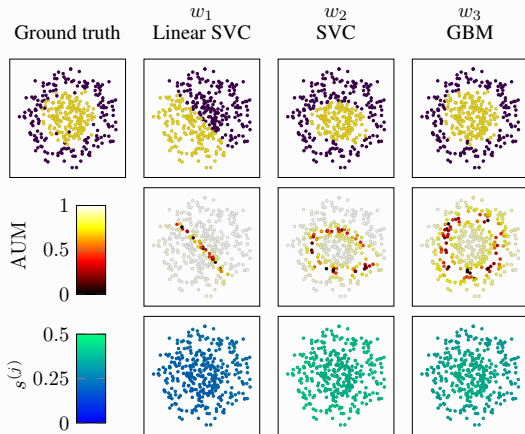
- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$
- For each worker j
 - ▶ Train a network on $\{(x_i, y_i^{(j)}); x_i \text{ is answered by worker } j\}$
 - ▶ Compute $\text{AUM}(x_i, y_i^{(j)})$ for the answered tasks x_i
 - ▶ Compute trust scores $s^{(j)}(x_i)$
 - ▶ For each task i compute $\text{WAUM}(x_i)$

Usage (for learning):

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

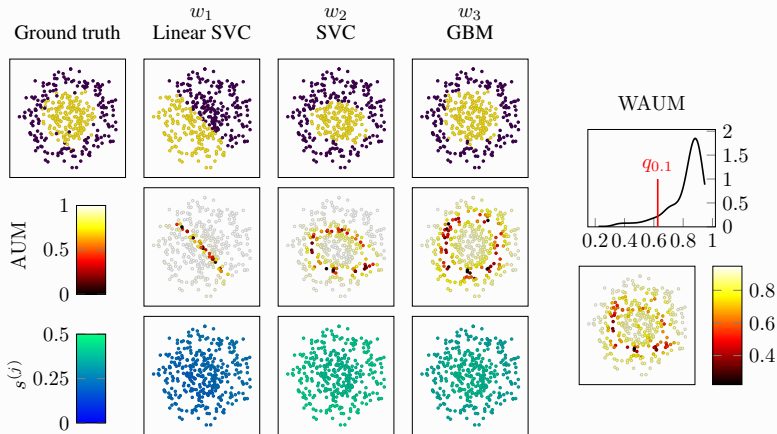
SIMULATION WITH CIRCLES

BINARY SETTING



SIMULATION WITH CIRCLES

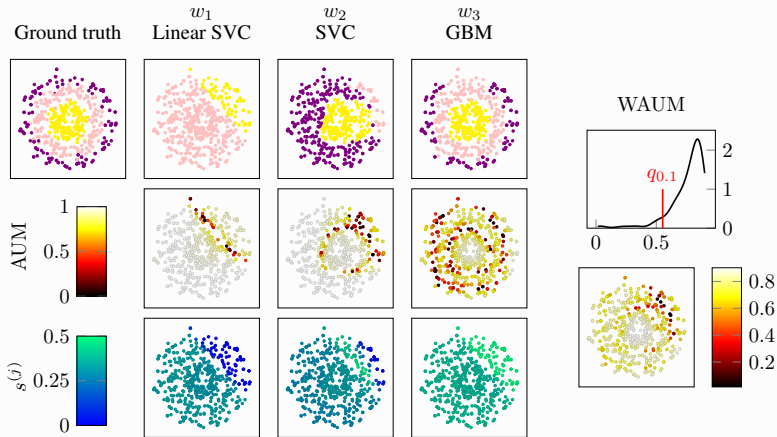
BINARY SETTING



- Workers = simulated classifiers (answering 500 tasks)
- Normalized trust scores

SIMULATION WITH CIRCLES

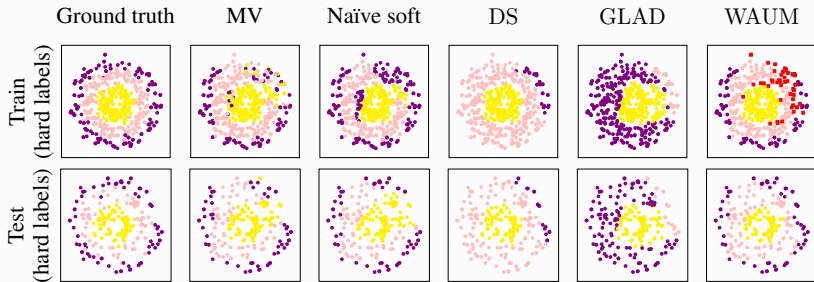
THREE CLASSES



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

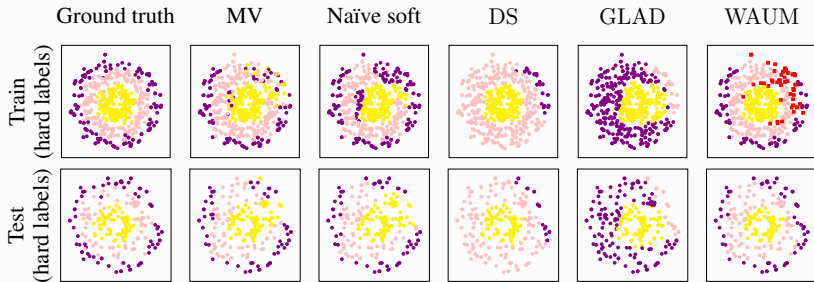
HOW CAN WE USE THE WAUM?

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



"3 answers per task is not enough!"

⁽²⁰⁾ C. Garcin, A. Joly, et al. (2021). "Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution". In: *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.

⁽²¹⁾ F. Rodrigues and F. Pereira (2018). "Deep learning from crowds". In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32. 1.



"3 answers per task is not enough!"

- Yes ! It is not
- ...but it happens → Pl@ntNet⁽²⁰⁾ (future work), LabelMe⁽²¹⁾
- LabelMe 1000 images (subset of LabelMe image segmentation project)
- Each image was labelled by 1, 2 or 3 workers

⁽²⁰⁾ C. Garcin, A. Joly, et al. (2021). "Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution". In: *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.

⁽²¹⁾ F. Rodrigues and F. Pereira (2018). "Deep learning from crowds". In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32. 1.



"3 answers per task is not enough!"

- Yes! It is not
- ... but it happens → Pl@ntNet⁽²⁰⁾ (future work), LabelMe⁽²¹⁾
- LabelMe 1000 images (subset of LabelMe image segmentation project)
- Each image was labelled by 1, 2 or 3 workers

LabelMe and task difficulty

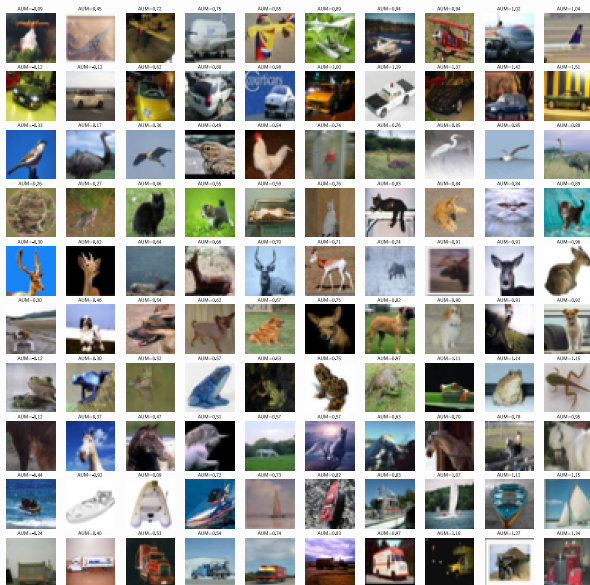
- Entropy is not reliable **at all**
- GLAD can't estimate a task difficulty for tasks with 1 label

⁽²⁰⁾ C. Garcin, A. Joly, et al. (2021). "Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution". In: *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.

⁽²¹⁾ F. Rodrigues and F. Pereira (2018). "Deep learning from crowds". In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32. 1.

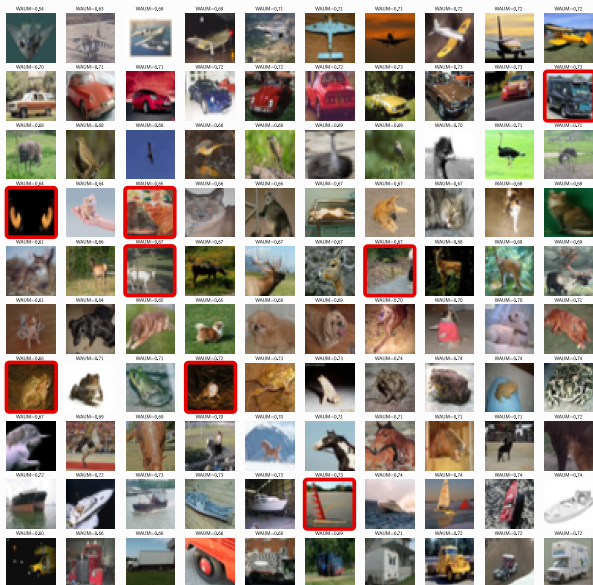
RESULTS ON CIFAR10H

IMPROVED MISLABELED DETECTIONS: WORST AUM/WAUM



RESULTS ON CIFAR10H

IMPROVED MISLABELED DETECTIONS: WORST AUM/WAUM





Back to the application of the AUM/WAUM to the CIFAR10H dataset.



Table: Label recovery, generalization performance and calibration error on the CIFAR-10H dataset by a Resnet-18

Aggregation method	Test accuracy (on CIFAR10-train)	ECE (expected calibration error)
MV	69.533 \pm 0.84	0.175 \pm 0.00
Naive soft	72.149 \pm 2.74	0.132 \pm 0.03
DS (vanilla)	70.268 \pm 0.93	0.173 \pm 0.00
DS (spam identification)	70.053 \pm 0.81	0.174 \pm 0.0
GLAD	66.569 \pm 8.48	0.173 \pm 0.01
WAUM	72.747 \pm 1.93	0.124 \pm 0.00

"CAN I USE THE WAUM IN MY FRAMEWORK?"

MOST PROBABLY YES



- Most frameworks are built on DS model
 - ▶ the WAUM only needs a neural network and $\hat{\pi}^{(j)}$

The Benefits of a Model of Annotation

Rebecca J. Passonneau
Center for Computational Learning Systems
Columbia University
New York, NY USA
becky@ccle.columbia.edu

Bob Carpenter
Department of Statistics
Columbia University
New York, NY USA
carp@lisa-1.com

**Analysis of Minimax Error Rate for Crowdsourcing
and Its Application to Worker Clustering Model**

Hideaki Imamura^{1,2} Toei Sato^{1,2} Masashi Sugiyama^{1,1}

The Third Second AAAI Conference
on Artificial Intelligence (AAAI-14)

Deep Learning from Crowds

Filipe Rodrigues, Francisco C. Pereira
Dept. of Management Engineering, Technical University of Denmark
Bygning 116B, 2800 Kgs. Lyngby, Denmark
rod@dtu.dk, camara@dtu.dk

Learning from Crowds by Modeling Common Confusions

Zhendong Chu, Jing Ma, Hongming Wang
Department of Computer Science, University of Virginia
{z9ey, jmlz, hw5}@virginia.edu

**Learning From Noisy Labels By
Regularized Estimation Of Annotator Confusion**

Ryutaro Tanno^{1,*} Ardavan Saeeedi² Swami Sankaranarayanan²
Daniel C. Alexander¹ Nathan Silberman²

¹University College London, UK ²Butterfly Network, New York, USA

¹{r.tanno, d.alexander}@ucl.ac.uk ²{saeeedi,swamiv, nlsilberman}@butterflynetwork.com



Take home message(s)

- Crowdsourcing / Label uncertainty : helpful for **data curating**



Take home message(s)

- Crowdsourcing / Label uncertainty : helpful for **data curating**
- Improved **data quality** \Rightarrow **improved learning** performance



Take home message(s)

- Crowdsourcing / Label uncertainty : helpful for **data curating**
- Improved **data quality** \Rightarrow **improved learning** performance
- (Fast) "stacked" WAUM : the presented version requires **one neural network per worker** (stacked version : **one neural network per dataset**)



Take home message(s)


- Crowdsourcing / Label uncertainty : helpful for **data curating**
- Improved **data quality** \Rightarrow **improved learning** performance
- (Fast) "stacked" WAUM : the presented version requires **one neural network per worker** (stacked version : **one neural network per dataset**)


Future work & wishful thinking

- ▶ Soon a crowdsourced module in benchopt
<https://benchopt.github.io/>
- ▶ Pl@ntnet crowdsourced dataset: coming, but it's messy (**2M workers**, 2 labels per task on average,...)

Tanguy Lefort: *"I swear that, if I make a crowdsourcing experiment, I will release both the tasks and labels"*

Contact:

 joseph.salmon@umontpellier.fr







 <http://josephsalmon.eu>








Github: @josephsalmon









Twitter: @salmonjsph





-  (N.d.). <https://github.com/googlecreativelab/quickdraw-dataset>.
-  Dawid, A. and A. Skene (1979). “Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm”. In: *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.
-  Garcin, C., A. Joly, et al. (2021). “Pl@ntNet-300K: a plant image dataset with high label ambiguity and a long-tailed distribution”. In: *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*.
-  Garcin, C., M. Servajean, et al. (2022). “Stochastic smoothing of the top-K calibrated hinge loss for deep imbalanced classification”. In: *ICML*.
-  Han, J., P. Luo, and X. Wang (2019). “Deep self-learning from noisy labels”. In: *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”. In: *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”. In: *CVPR*, pp. 1468–1477.
-  LeCun, Y. et al. (1998). “Gradient-based learning applied to document recognition”. In: *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”. In: *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”. In: *J. Artif. Intell. Res.* 70, pp. 1373–1411.
-  Peterson, J. C. et al. (2019). “Human Uncertainty Makes Classification More Robust”. In: *ICCV*, pp. 9617–9626.
-  Pleiss, G. et al. (2020). “Identifying mislabeled data using the area under the margin ranking”. In: *NeurIPS*.



-  Rodrigues, F. and F. Pereira (2018). “Deep learning from crowds”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32. 1.
-  Russakovsky, O. et al. (2015). “ImageNet Large Scale Visual Recognition Challenge”. In: *Int. J. Comput. Vision* 115.3, pp. 211–252.
-  Servajean, M. et al. (2017). “Crowdsourcing thousands of specialized labels: A Bayesian active training approach”. In: *IEEE Transactions on Multimedia* 19.6, pp. 1376–1391.
-  Siddiqui, S. A. et al. (2022). *Metadata Archaeology: Unearthing Data Subsets by Leveraging Training Dynamics*.
-  Snow, R. et al. (2008). “Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks”. In: *Conference on Empirical Methods in Natural Language Processing*. EMNLP 2008. Association for Computational Linguistics, pp. 254–263.
-  Uday Prabhu, V. and A. Birhane (June 2020). “Large image datasets: A pyrrhic win for computer vision?” In: *arXiv e-prints*, arXiv:2006.16923, arXiv:2006.16923.



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