

LABEL AMBIGUITY IN CROWDSOURCING FOR CLASSIFICATION AND EXPERT FEEDBACK

Tanguy Lefort

IMAG, Univ Montpellier, CNRS

INRIA, LIRMM,

Supervised by

Benjamin Charlier

Alexis Joly

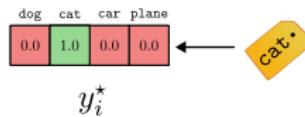
and Joseph Salmon



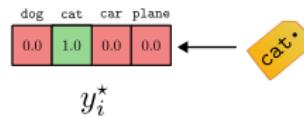
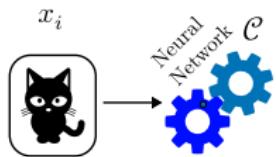
HOW TO TRAIN YOUR CLASSIFIER



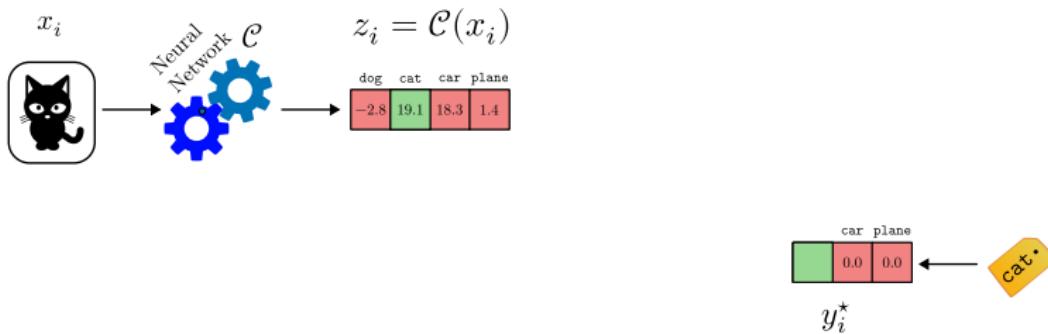
x_i



HOW TO TRAIN YOUR CLASSIFIER

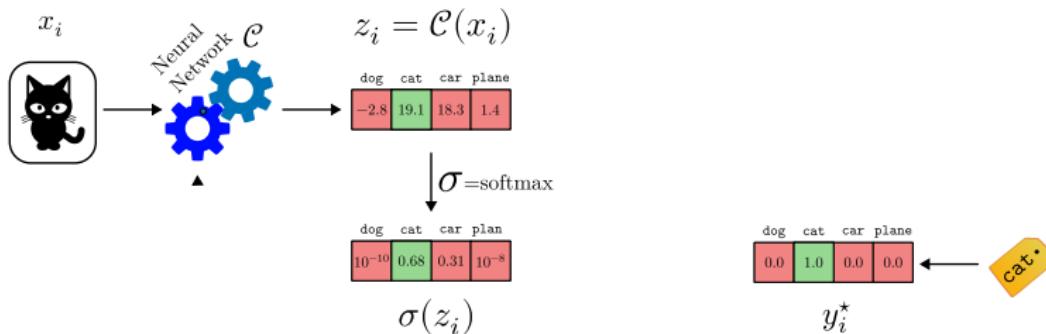


HOW TO TRAIN YOUR CLASSIFIER

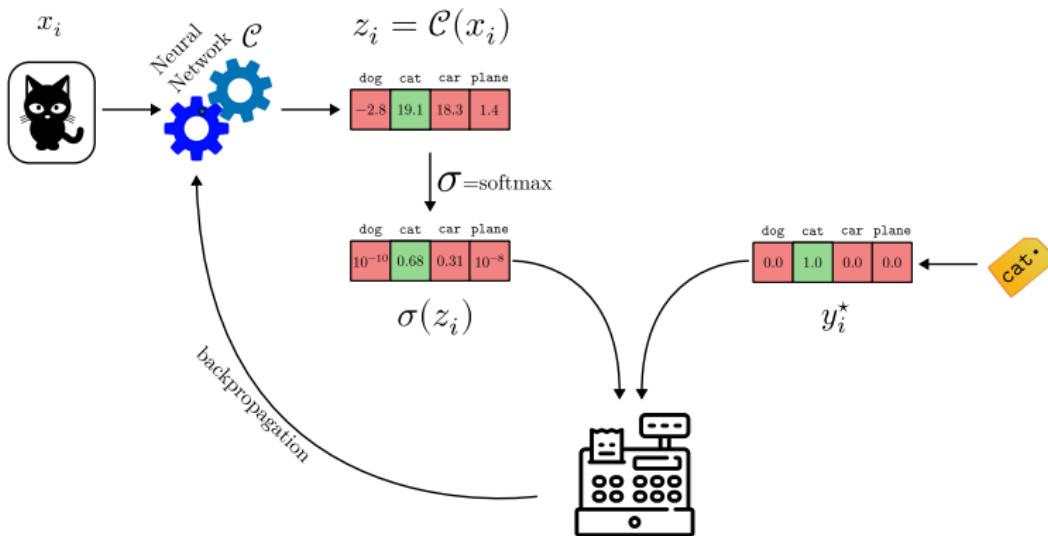


HOW TO TRAIN YOUR CLASSIFIER

1

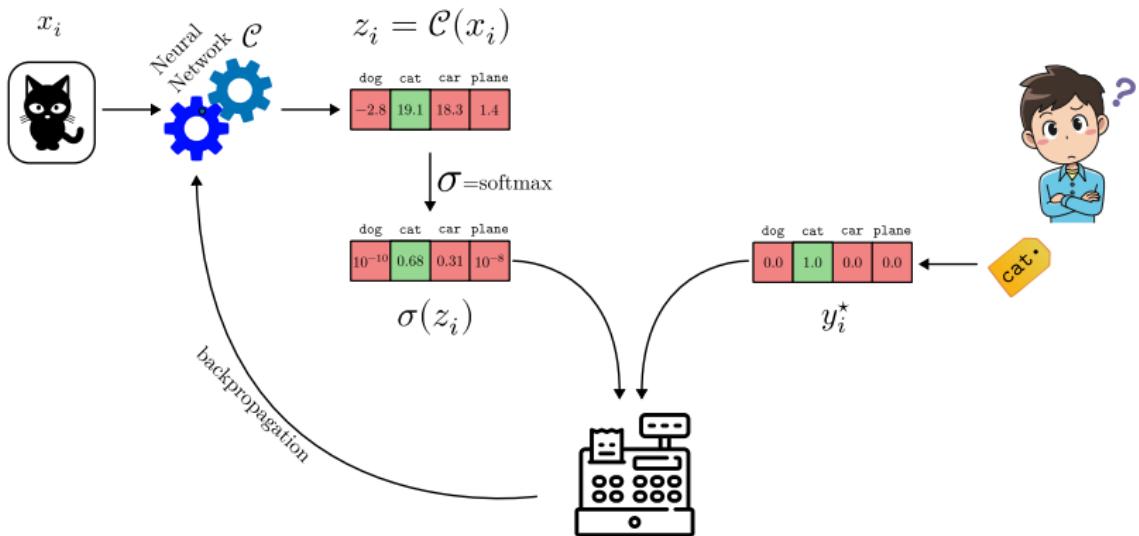


HOW TO TRAIN YOUR CLASSIFIER



HOW TO TRAIN YOUR CLASSIFIER

1



ASK CITIZENS TO LABEL OUR DATA

FRAMEWORK AND NOTATIONS

2

- ▶ Workers sort a given task into one of the K classes

$K = 4$

0: car 2: cat $\mathcal{A}(x_2)$

1: plane 3: dog

x_1 x_2

	w_1	w_2	w_3	w_4	w_5	
0: car						
1: plane						
x_1						
x_2						
$\mathcal{T}(w_3)$						

y_i^*

2

0

ASK CITIZENS TO LABEL OUR DATA

FRAMEWORK AND NOTATIONS

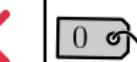
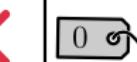
2

- ▶ Workers sort a given task into one of the K classes

$K = 4$

0: car 2: cat 1: plane 3: dog

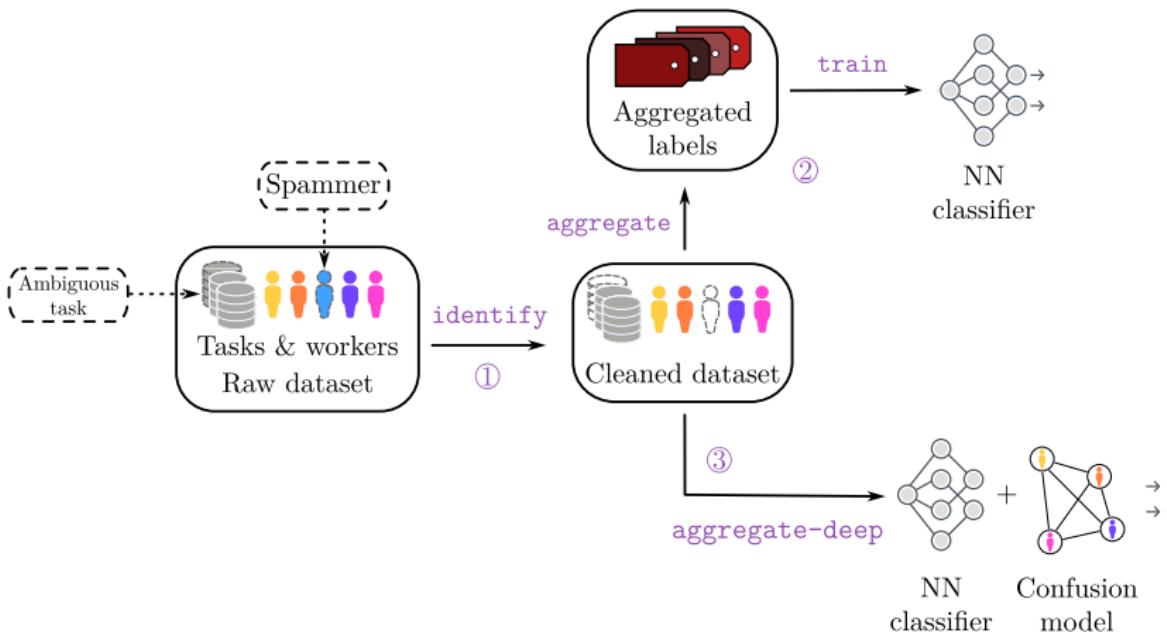
$\mathcal{A}(x_2)$

	w_1	w_2	w_3	w_4	w_5
x_1					
x_2	 2 ↗	 2 ↗	 0 ↗	 2 ↗	 3 ↗
$\mathcal{T}(w_3)$	 X	 X	 0 ↗	 0 ↗	 3 ↗
					y_i^*
				2 ↗	0 ↗

- ▶ $y_i^{(j)} \in [K] :=$ answer of worker j to task i

FROM THE DATA TO THE CLASSIFIER

THE PIPELINE



MAIN CONTRIBUTIONS

- ▶ Can we improve performance by leveraging better-quality data

(1) T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

(2) T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

MAIN CONTRIBUTIONS



- ▶ Can we improve performance by leveraging better-quality data
- ▶ Can we standardize crowdsourcing dataset's tools in python for reproducibility?

(1) T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

(2) T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

MAIN CONTRIBUTIONS



- ▶ Can we improve performance by leveraging better-quality data

- ▶ Can we standardize crowdsourcing dataset's tools in python for reproducibility?

- ▶ Release a real, large-scaled dataset with Pl@ntNet

(1) T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

(2) T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

MAIN CONTRIBUTIONS



- ▶ Can we improve performance by leveraging better-quality data
 - ▶ Creation of the **WAUM**⁽¹⁾: a metric to identify ambiguous images
- ▶ Can we standardize crowdsourcing dataset's tools in python for reproducibility?
- ▶ Release a real, large-scaled dataset with Pl@ntNet

⁽¹⁾ T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

⁽²⁾ T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

MAIN CONTRIBUTIONS



- ▶ Can we improve performance by leveraging better-quality data
 - ▶ Creation of the **WAUM**⁽¹⁾: a metric to identify ambiguous images
- ▶ Can we standardize crowdsourcing dataset's tools in python for reproducibility?
 - ▶ **peerannot**⁽²⁾: available via pip/github
<https://peerannot.github.io>
- ▶ Release a real, large-scaled dataset with Pl@ntNet

⁽¹⁾ T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

⁽²⁾ T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

MAIN CONTRIBUTIONS

- ▶ Can we improve performance by leveraging better-quality data
 - ▶ Creation of the **WAUM**⁽¹⁾: a metric to identify ambiguous images
- ▶ Can we standardize crowdsourcing dataset's tools in python for reproducibility?
 - ▶ **peerannot**⁽²⁾: available via pip/github
<https://peerannot.github.io>
- ▶ Release a real, large-scaled dataset with Pl@ntNet
 - ▶ Creation and evaluation of a **new benchmark dataset**

⁽¹⁾ T. Lefort et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

⁽²⁾ T. Lefort et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.

CLASSICAL AGGREGATION STRATEGY

(WEIGHTED) MAJORITY VOTES

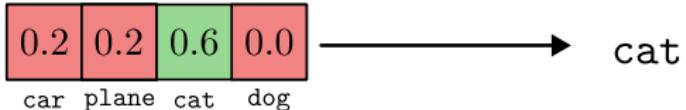


$$\hat{y}_i^{\text{WMV}} = \operatorname{argmax}_{k \in [K]} \sum_{j \in \mathcal{A}(x_i)} \text{weight} \mathbb{1}(y_i^{(j)} = k)$$

For example with equal weights:



0: car 2: cat
1: plane 3: dog



CLASSICAL AGGREGATION STRATEGY

(WEIGHTED) MAJORITY VOTES

6

Many existing weight choices:

- ▶ WAWA⁽³⁾:

$$\text{weight}(w_j) = \text{Accuracy}(\{y_i^{(j)}\}_i, \{\hat{y}_i^{\text{MV}}\}_i)$$

- ▶ Shapley based⁽⁴⁾

- ▶ Matrix completion (MACE)⁽⁵⁾ ...

Pros: "simple" weight can scale to large datasets and be easy to interpret
Cons: Can not capture worker skills in detail

⁽³⁾ <https://success.appen.com/hc/en-us/articles/202703205-Calculating-Worker-Agreement-with-Aggregate-Wawa>

⁽⁴⁾ T. Lefort et al. (July 2024b). "Weighted majority vote using Shapley values in crowdsourcing". In: CAp 2024 - Conférence sur l'Apprentissage Automatique. Lille, France.

⁽⁵⁾ D. Hovy et al. (2013). "Learning whom to trust with MACE". In: *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1120–1130.

CLASSICAL AGGREGATION STRATEGY

DAWID AND SKENE⁽⁶⁾



- ▶ Represent workers from their pairwise confusions matrix $\pi^{(j)} \in \mathbb{R}^{K \times K}$
- ▶ Probabilistic model on their answers:

$$y^{(j)} | y^* \sim \text{Multinomial}(\pi_{y^*, \bullet}^{(j)})$$

with $\pi_{k,\ell}^{(j)} = \mathbb{P}(\text{worker } j \text{ answers } \ell \text{ with unknown truth } k)$

Pros:

- ▶ Finer modelisation
- ▶ Can use adversarial workers

Cons:

- ▶ Memory issue: $n_{\text{worker}} \times K^2$ parameters to estimate only the confusion matrices

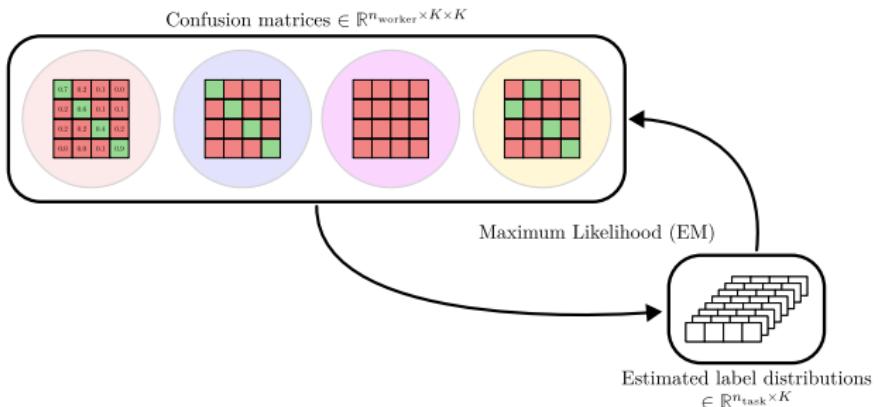
⁽⁶⁾ A. Dawid and A. Skene (1979). "Maximum Likelihood Estimation of Observer Error-Rates Using the EM Algorithm". In: *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.

CLASSICAL AGGREGATION STRATEGY

DAWID AND SKENE – MODEL



Probabilistic model → Likelihood (to maximize)

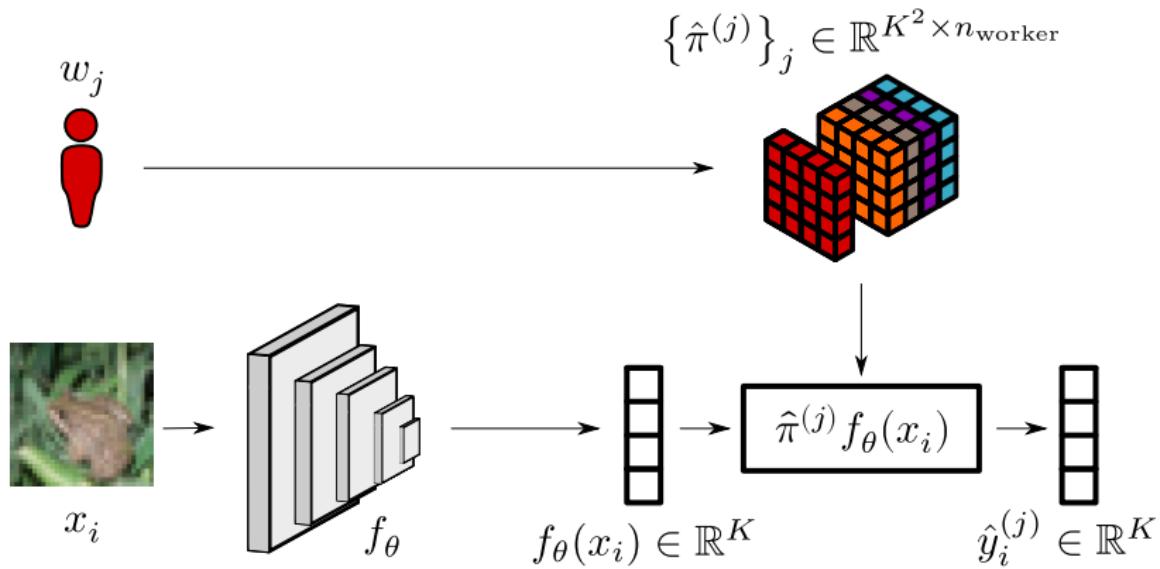


CLASSICAL DEEP-LEARNING STRATEGY

CROWDLAYER⁽⁷⁾



- Idea: put the DS confusion matrix in a neural network as a new layer



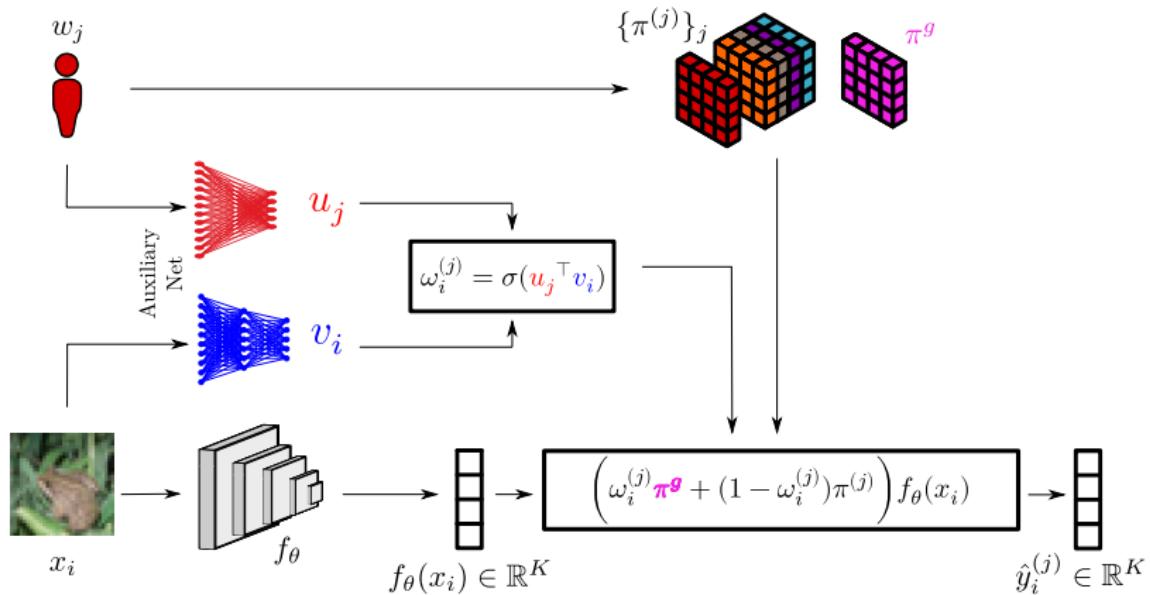
⁽⁷⁾ F. Rodrigues and F. Pereira (2018). "Deep learning from crowds". In: AAAI, vol. 32.

CLASSICAL DEEP-LEARNING STRATEGY

CoNAL⁽⁸⁾



► Idea: CrowdLayer + global and local confusions



⁽⁸⁾ Z. Chu, J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions." In: AAAI, pp. 5832–5840.

WHEN IMAGES HAVE UNDERLYING AMBIGUITY



$K = 4$

0: car 2: cat
1: plane 3: dog

$\mathcal{T}(w_3)$

		w_1	w_2	w_3	w_4	w_5	$\mathcal{A}(x_2)$	
x_1								
x_2								

WHEN IMAGES HAVE UNDERLYING AMBIGUITY



$K = 4$

0: car 2: cat w_1 w_2 w_3 w_4 w_5

1: plane 3: dog

$\mathcal{A}(x_2)$

$\mathcal{T}(w_3)$

		w_1	w_2	w_3	w_4	w_5
x_1						
x_2						
x_3						

y_i^*

AMBIGUITY IN CLASSICAL SUPERVISED SETTING

AREA UNDER THE MARGIN (AUM)



Goal: identify issues in classical datasets $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$

- ▶ AUM⁽⁹⁾: monitor margin during training
- ▶ Classifier: at epoch $t \in [T]$, $\mathcal{C}^{(t)}(x_i) \in \mathbb{R}^K$ a vector of **scores** (logits)

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

Diagram annotations:

- Average = Stability (red bracket above the fraction)
- Score of assigned label (blue bracket under $\mathcal{C}^{(t)}(x_i)_{y_i}$)
- Margin between scores: content of Hinge loss (pink bracket above the difference term)
- Other maximum score (pink bracket under $\max_{\ell \neq y_i} \mathcal{C}^{(t)}(x_i)_\ell$)

⁽⁹⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: NeurIPS.

AMBIGUITY IN CLASSICAL SUPERVISED SETTING

AREA UNDER THE MARGIN (AUM)



Goal: identify issues in classical datasets $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$

- ▶ AUM⁽⁹⁾: monitor margin during training
- ▶ 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

⁽⁹⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: NeurIPS.

AMBIGUITY IN CLASSICAL SUPERVISED SETTING

AREA UNDER THE MARGIN (AUM)



Goal: identify issues in classical datasets $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$

- ▶ AUM⁽⁹⁾: monitor margin during training
- ▶ 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

⁽⁹⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: NeurIPS.

AMBIGUITY IN CLASSICAL SUPERVISED SETTING

AREA UNDER THE MARGIN (AUM)



Goal: identify issues in classical datasets $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times [K]$

- ▶ AUM⁽⁹⁾: monitor margin during training
- ▶ Classifier: at epoch $t \in [T]$, $\mathcal{C}^{(t)}(x_i) \in \mathbb{R}^K$ a vector of **scores** (logits)

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

Diagram annotations:

- Average = Stability: Points to the term $\frac{1}{T} \sum_{t=1}^T$.
- Margin between scores: content of Hinge loss: Points to the term $\mathcal{C}^{(t)}(x_i)_{y_i} - \max_{\ell \neq y_i} \mathcal{C}^{(t)}(x_i)_\ell$.
- Score of assigned label: Points to the term $\mathcal{C}^{(t)}(x_i)_{y_i}$.
- Other maximum score: Points to the term $\max_{\ell \neq y_i} \mathcal{C}^{(t)}(x_i)_\ell$.

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 $\ell \neq y_i$.

⁽⁹⁾ G. Pleiss et al. (2020). "Identifying mislabeled data using the area under the margin ranking". In: NeurIPS.

GOING TO THE CROWDSOURCING SETTING

AUMC

Naive extension: Use the MV label \hat{y}_i^{MV} instead of the unknown y_i + use previous work of margins' properties⁽¹⁰⁾

$$\text{AUMC}(x_i, \hat{y}_i^{\text{MV}}) = \frac{1}{T} \sum_{t=1}^T \left[\mathcal{C}^{(t)}(x_i)_{\hat{y}_i^{\text{MV}}} - \mathcal{C}^{(t)}(x_i)_{[2]} \right]$$

Score of MV label

Average = Stability

Margin between scores:
margin for top-1 classification

Other maximum score

The diagram illustrates the formula for AUMC. It shows the equation $\text{AUMC}(x_i, \hat{y}_i^{\text{MV}}) = \frac{1}{T} \sum_{t=1}^T \left[\mathcal{C}^{(t)}(x_i)_{\hat{y}_i^{\text{MV}}} - \mathcal{C}^{(t)}(x_i)_{[2]} \right]$. A red bracket above the term $\mathcal{C}^{(t)}(x_i)_{\hat{y}_i^{\text{MV}}} - \mathcal{C}^{(t)}(x_i)_{[2]}$ is labeled "Margin between scores: margin for top-1 classification". A blue bracket below the term $\mathcal{C}^{(t)}(x_i)_{\hat{y}_i^{\text{MV}}}$ is labeled "Score of MV label". A red bracket above the entire sum is labeled "Average = Stability". A blue bracket below the entire sum is labeled "Other maximum score".

Issue:

- Lose all worker-related information
- Sensitive to poorly performing workers

⁽¹⁰⁾ M. Lapin, M. Hein, and B. Schiele (2016). "Loss functions for top-k error: Analysis and insights". In: CVPR, pp. 1468–1477.

Weighted Areas Under the Margins:

- ▶ Scale effects in the scores discarded, need normalization⁽¹¹⁾

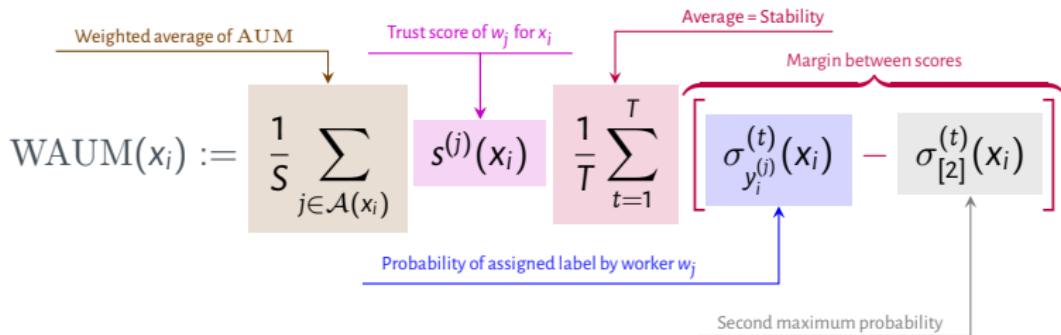
⁽¹¹⁾ C. Ju, A. Bibaut, and M. van der Laan (2018). "The relative performance of ensemble methods with deep convolutional neural networks for image classification". In: *J. Appl. Stat.* 45.15, pp. 2800–2818.

Weighted Areas Under the Margins:

- Scale effects in the scores discarded, need normalization⁽¹¹⁾

With:

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



⁽¹¹⁾ 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.

WEIGHTS IN THE WAUM

LEVERAGE BOTH TASKS AND LABELS



Our chosen worker/task score:

- Consider a score of the form: worker term \times task term (similar to GLAD⁽¹²⁾)
- Use multidimensionality of DS confusion matrices
- Use a neural network as a control agent to measure task difficulty⁽¹³⁾

⁽¹²⁾J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". In: *NeurIPS*. vol. 22.

⁽¹³⁾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.

WEIGHTS IN THE WAUM

LEVERAGE BOTH TASKS AND LABELS



15

Our chosen worker/task score:

- Consider a score of the form: worker term \times task term (similar to GLAD⁽¹²⁾)
- Use multidimensionality of DS confusion matrices
- Use a neural network as a control agent to measure task difficulty⁽¹³⁾

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

↑
Worker j overall ability ↑
Difficulty of task i

⁽¹²⁾ J. Whitehill et al. (2009). "Whose Vote Should Count More: Optimal Integration of Labels from Labelers of Unknown Expertise". In: *NeurIPS*. vol. 22.

⁽¹³⁾ 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.

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED

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

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED

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

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED

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

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED



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

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED



- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
 - Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
 - Compute all AUM($x_i, y_i^{(j)}$)
 - 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 $\pi^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
 - Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
 - Compute all $\text{AUM}(x_i, y_i^{(j)})$
 - 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 $\alpha = 0.01$)
- **Estimate confusion matrices** $\hat{\pi}^{(j)}$ on pruned training dataset

COMPUTING THE WAUM

THE PIPELINE SUMMARIZED

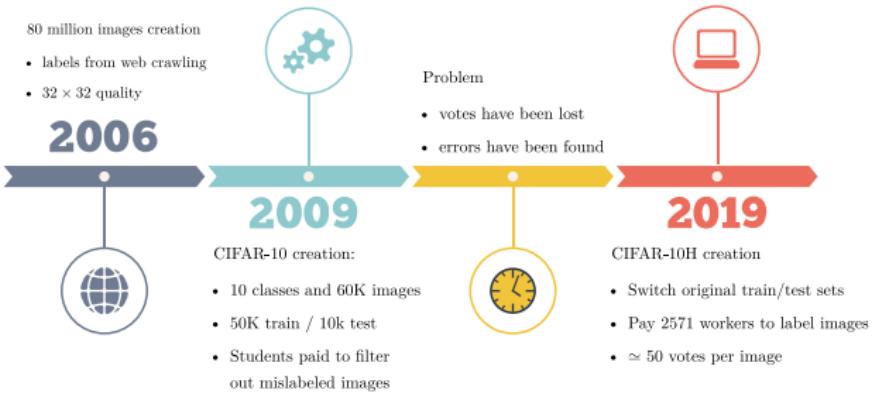


- Estimate confusion matrices $\pi^{(j)} \in \mathbb{R}^{K \times K}$, for all $j \in [n_{\text{worker}}]$
 - Train a network on all crowdsourced task/label pairs: $(x_i, y_i^{(j)})$
 - Compute all $\text{AUM}(x_i, y_i^{(j)})$
 - Compute trust scores $s^{(j)}(x_i)$
- $$\bullet \text{For each task compute } \text{WAUM}(x_i) = \frac{\sum_{j \in \mathcal{A}(x_i)} s^{(j)}(x_i) \text{AUM}(x_i, y_i^{(j)})}{\sum_{j' \in \mathcal{A}(x_i)} s^{(j')}(x_i)}$$

Usage (for learning):

- **Prune** x_i 's with $\text{WAUM}(x_i)$ below quantile q_α (say $\alpha = 0.01$)
- **Estimate confusion matrices** $\hat{\pi}^{(j)}$ on pruned training dataset
- **Aggregate labels** and **Train** a classifier on the newly pruned dataset

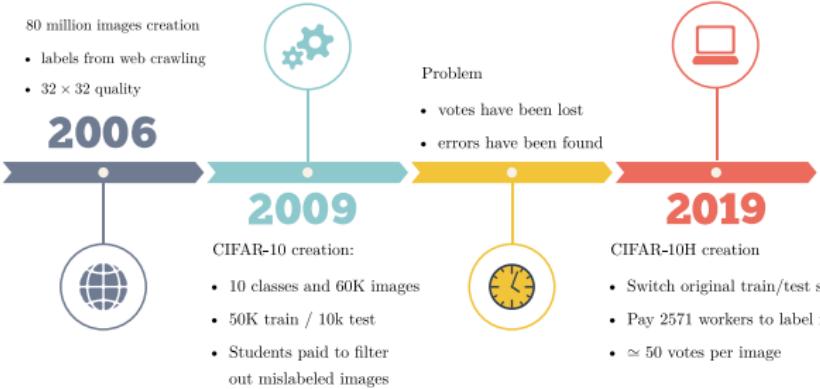
PRESENTING CIFAR-10H⁽¹⁴⁾ DATASET



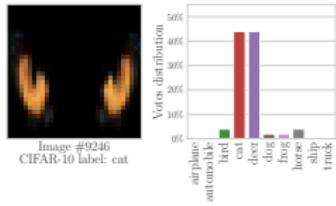
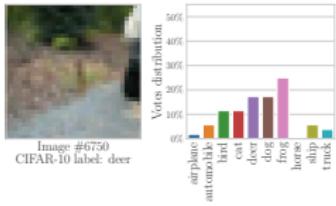
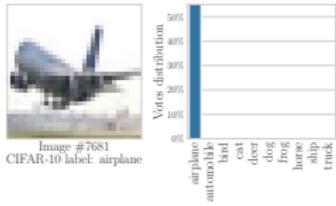
Labels: cat, dog, car, plane, bird, horse, frog, deer, ship, truck

⁽¹⁴⁾J. C. Peterson et al. (2019). "Human Uncertainty Makes Classification More Robust". In: ICCV, pp. 9617–9626.

PRESENTING CIFAR-10H⁽¹⁴⁾ DATASET



Labels: cat, dog, car, plane, bird, horse, frog, deer, ship, truck



⁽¹⁴⁾J. C. Peterson et al. (2019). "Human Uncertainty Makes Classification More Robust". In: ICCV, pp. 9617–9626.

PRESENTING LABELME DATASET⁽¹⁵⁾



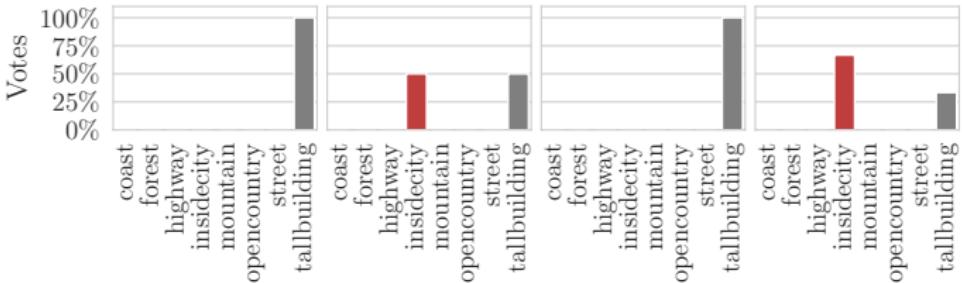
- ▶ 1000 training / 500 validation / 1188 test images
- ▶ 59 workers: each task has up to 3 votes
- ▶ 8 classes: highway, insidecity, tallbuilding, street, forest, coast, mountain, opencountry

⁽¹⁵⁾ F. Rodrigues, F. Pereira, and B. Ribeiro (2014). "Gaussian process classification and active learning with multiple annotators". In: ICML. PMLR, pp. 433–441.

PRESENTING LABELME DATASET⁽¹⁵⁾



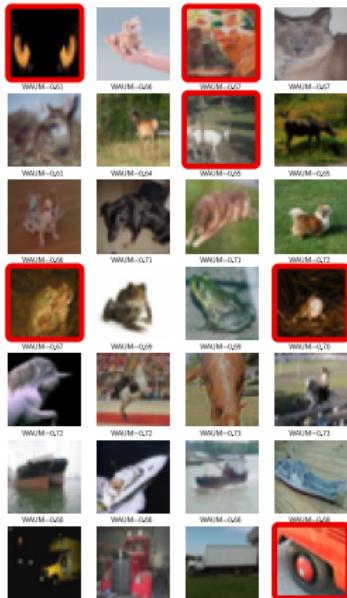
- ▶ 1000 training / 500 validation / 1188 test images
- ▶ 59 workers: each task has up to 3 votes
- ▶ 8 classes: highway, insidecity, tallbuilding, street, forest, coast, mountain, opencountry



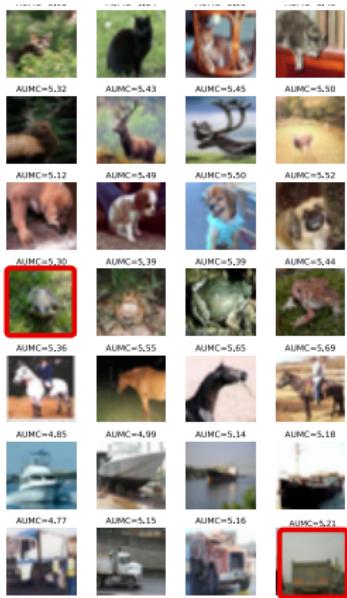
⁽¹⁵⁾ F. Rodrigues, F. Pereira, and B. Ribeiro (2014). "Gaussian process classification and active learning with multiple annotators". In: ICML. PMLR, pp. 433–441.

QUALITATIVE RESULTS

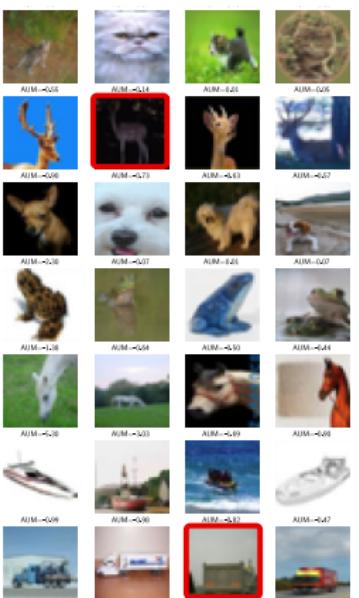
WAUM



AUMC



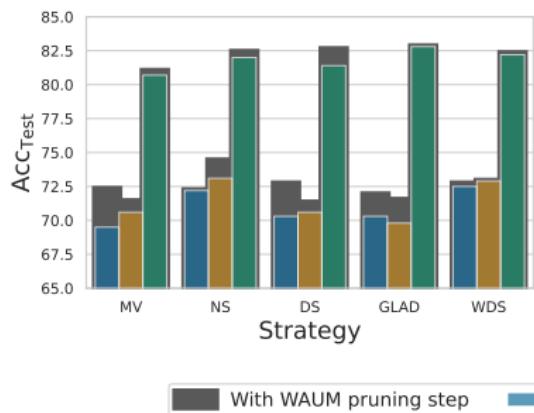
AUM



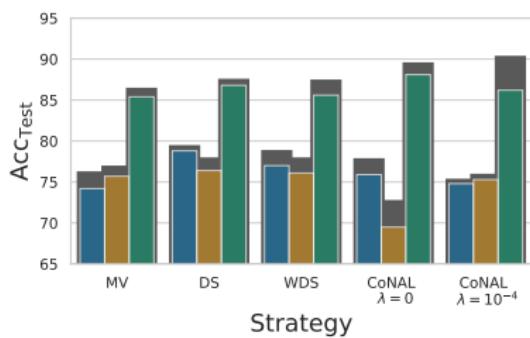
ABLATION STUDY



CIFAR-10H



LabelMe





In short

- ▶ Introduced the WAUM to find ambiguous images
- ▶ Better quality data can improve performance

In short

- ▶ Introduced the WAUM to find ambiguous images
- ▶ Better quality data can improve performance

Towards large-scale problems

- ▶ DS model and confusion matrices do not scale
- ▶ What is currently done in large-scale settings?
- ▶ Can we evaluate their performance?

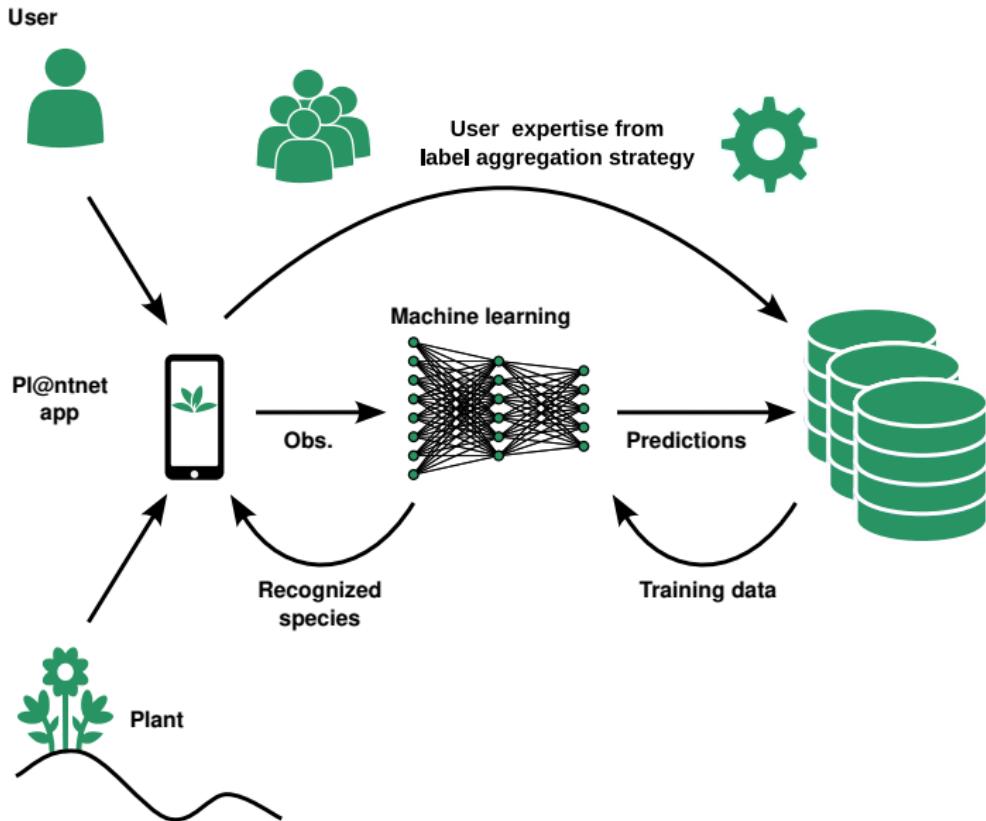
In short

- ▶ Introduced the WAUM to find ambiguous images
- ▶ Better quality data can improve performance

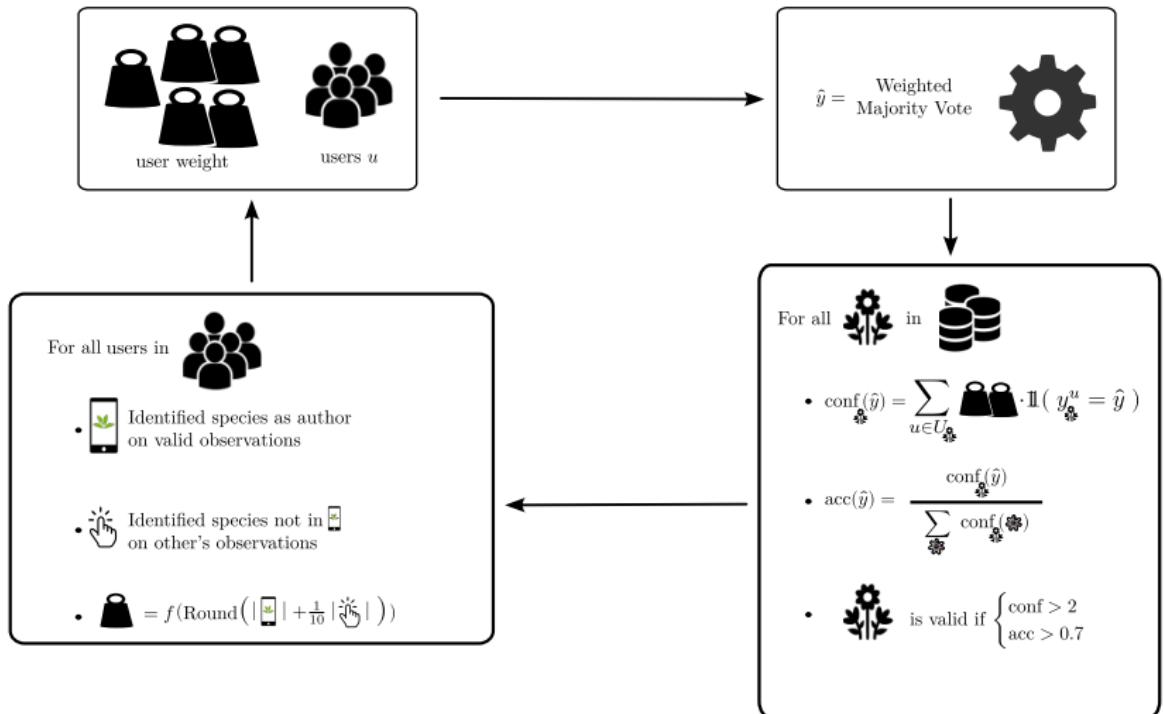
Towards large-scale problems

- ▶ DS model and confusion matrices do not scale
- ▶ What is currently done in large-scale settings?
- ▶ Can we evaluate their performance?
 - ▶ To evaluate we need the dataset!

PRESENTING PL@NTNET PIPELINE



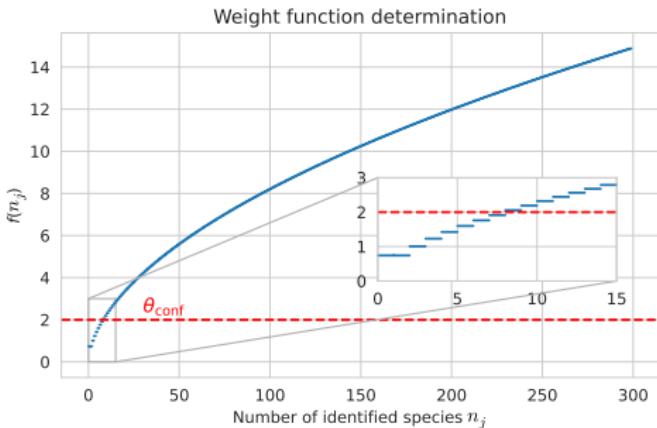
PL@NTNET AGGREGATION STRATEGY



PL@NTNET AGGREGATION STRATEGY

WEIGHT FUNCTION

$$f(n_j) = n_j^\alpha - n_j^\beta + \gamma \text{ with } \begin{cases} \alpha &= 0.5 \\ \beta &= 0.2 \\ \gamma &\simeq 0.74 \end{cases}$$

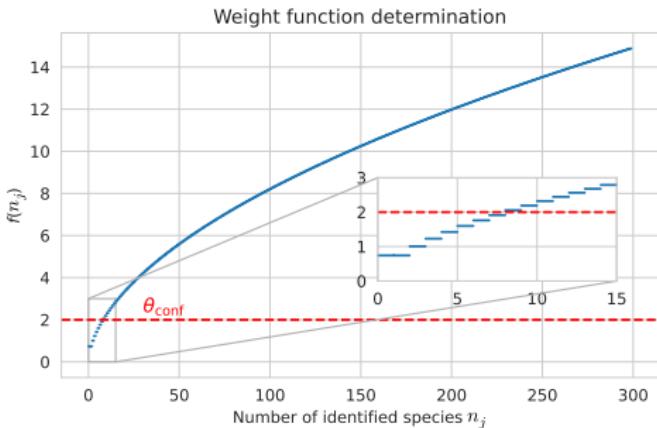


- With 8 identified species one becomes self-validating

PL@NTNET AGGREGATION STRATEGY

WEIGHT FUNCTION

$$f(n_j) = n_j^\alpha - n_j^\beta + \gamma \text{ with } \begin{cases} \alpha &= 0.5 \\ \beta &= 0.2 \\ \gamma &\simeq 0.74 \end{cases}$$



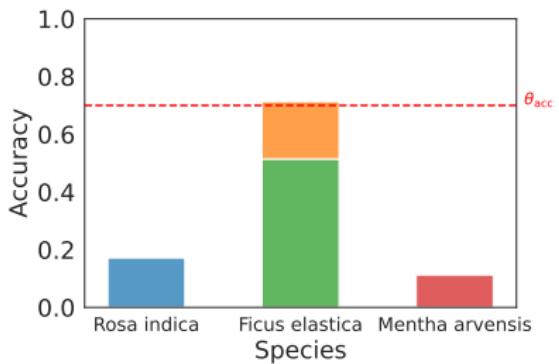
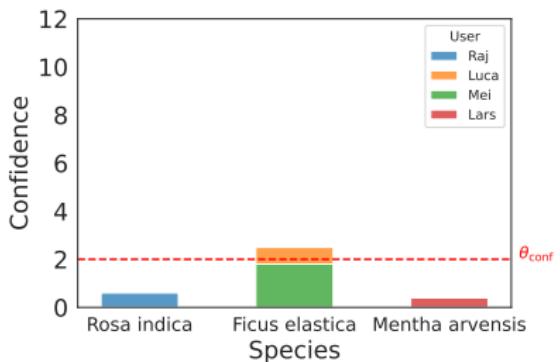
- With 8 identified species one becomes self-validating
- But observations can be invalidated at any time in the future

PL@NTNET AGGREGATION STRATEGY

EXAMPLES



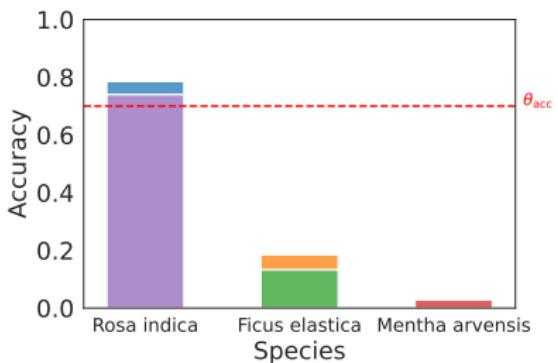
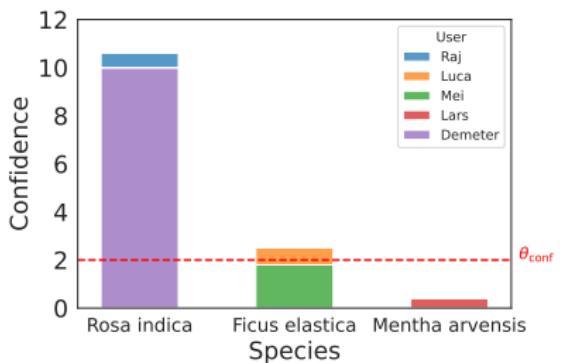
Initial setting



PL@NTNET AGGREGATION STRATEGY EXAMPLES



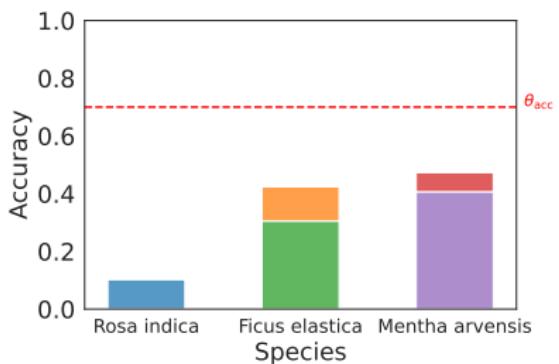
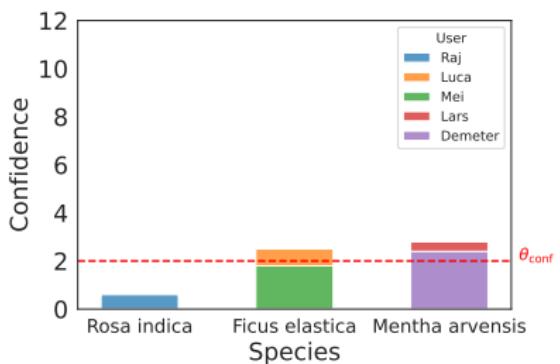
Label switch



PL@NTNET AGGREGATION STRATEGY EXAMPLES



Invalidate



REALEASING A NEW DATASET



- ▶ South Western European flora obs since 2017
- ▶ 823 000 users answered more than 11000 species
- ▶ 6 700 000 observations
- ▶ 9 000 000 votes casted
- ▶ **Imbalance:** 80% of observations are represented by 10% of total votes

REALEASING A NEW DATASET



- ▶ South Western European flora obs since 2017
 - ▶ 823 000 users answered more than 11000 species
 - ▶ 6 700 000 observations
 - ▶ 9 000 000 votes casted
 - ▶ **Imbalance:** 80% of observations are represented by 10% of total votes
-
- ▶ Extraction of 98 experts (TelaBotanic + prior knowledge – thanks to Pierre Bonnet)

COMPARED STRATEGIES

- ▶ **Majority Vote (MV)**

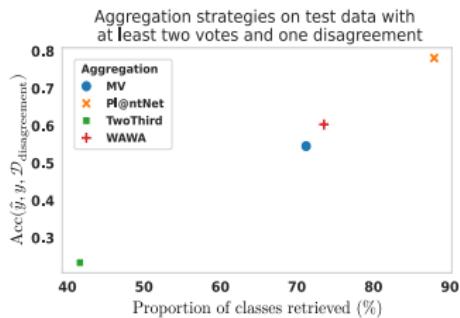
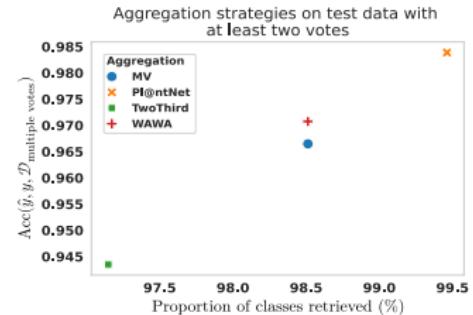
COMPARED STRATEGIES

- ▶ **Majority Vote (MV)**
- ▶ **Worker agreement with aggregate (WAWA) (appen)**
 - ▶ Majority vote
 - ▶ Weight user by how much they agree with the majority
 - ▶ Weighted majority vote

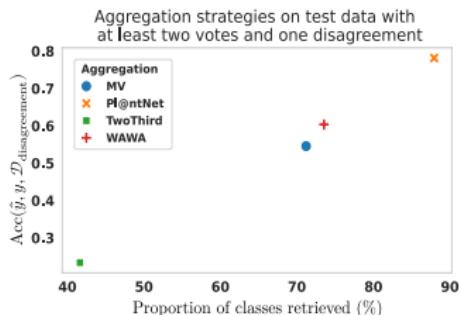
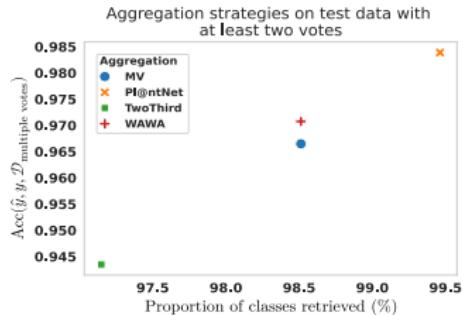
COMPARED STRATEGIES

- ▶ **Majority Vote (MV)**
- ▶ **Worker agreement with aggregate (WAWA) (appen)**
 - ▶ Majority vote
 - ▶ Weight user by how much they agree with the majority
 - ▶ Weighted majority vote
- ▶ **iNaturalist**
 - ▶ Need 2 votes
 - ▶ 2/3 of agreements

RESULTS



RESULTS



In short

- ▶ PI@ntNet aggregation performs better overall
- ▶ TwoThird is highly impacted by their reject threshold
- ▶ In ambiguous settings (right), strategies weighting users are better

MORE EXPERIMENTS



- ▶ AI vote's integration
- ▶ Penalizing user's mistakes

MORE EXPERIMENTS

- ▶ AI vote's integration
 - ▶ include AI's vote in the aggregation
 - ▶ Risks of model collapse
 - ▶ Results suggest using AI's vote if the score is above a threshold
 $\theta_{score} = 0.7$ can improve performance
- ▶ Penalizing user's mistakes

- ▶ AI vote's integration
 - ▶ include AI's vote in the aggregation
 - ▶ Risks of model collapse
 - ▶ Results suggest using AI's vote if the score is above a threshold
 $\theta_{score} = 0.7$ can improve performance
- ▶ Penalizing user's mistakes
 - ▶ users with large number of erroneous votes would be penalized in weight
 - ▶ Currently this does not improve performance
 - ▶ Taking the Accuracy of each user in the weight function shows however close performance from the current strategy

CONCLUSION AND PERSPECTIVES

KEY POINTS



- ▶ We released a new dataset for large-scale crowdsourced data:
<https://zenodo.org/records/10782465>
- ▶ We evaluated the current label aggregation in Pl@ntNet against other scaling ones and show better performance overall

CONCLUSION AND PERSPECTIVES

KEY POINTS



- ▶ We released a new dataset for large-scale crowdsourced data:
<https://zenodo.org/records/10782465>
 - ▶ We evaluated the current label aggregation in Pl@ntNet against other scaling ones and show better performance overall
-
-
-
- ▶ There is a need to *better collect votes* and not *waste* expert knowledge → recommendation system
 - ▶ We worked on the aggregation of species determination, users can enter other information – the observation is a flower/organ/etc.

Thank you

Peerannot library

<https://peerannot.github.io>

-  Chu, Z., J. Ma, and H. Wang (2021). "Learning from Crowds by Modeling Common Confusions.". In: *AAAI*, pp. 5832–5840.
-  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.
-  Hovy, D. et al. (2013). "Learning whom to trust with MACE". In: *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1120–1130.
-  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.
-  Lapin, M., M. Hein, and B. Schiele (2016). "Loss functions for top-k error: Analysis and insights". In: *CVPR*, pp. 1468–1477.
-  Lefort, T. et al. (2022). "Identify ambiguous tasks combining crowdsourced labels by weighting Areas Under the Margin". In: *arXiv preprint arXiv:2209.15380*.

-  Lefort, T. et al. (May 17, 2024a). "Peerannot: Classification for Crowdsourced Image Datasets with Python". In: *Computo*.
-  — (July 2024b). "Weighted majority vote using Shapley values in crowdsourcing". In: *CAp 2024 - Conférence sur l'Apprentissage Automatique*. Lille, France.
-  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: *AAAI*. Vol. 32.
-  Rodrigues, F., F. Pereira, and B. Ribeiro (2014). "Gaussian process classification and active learning with multiple annotators". In: *ICML*. PMLR, pp. 433–441.
-  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.

REFERENCES III



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