

SUPERVISED LEARNING BY CROWDSOURCING

Joseph Salmon

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



UNIVERSITÉ DE
MONTPELLIER



Inria



- ▶ **Tanguy Lefort** (IMAG, Inria, LIRMM, Univ Montpellier, CNRS) Ph.D. student, looking for a post-doc next year!
- ▶ Benjamin Charlier (IMAG, Univ Montpellier, CNRS)
- ▶ Alexis Joly (Inria, LIRMM, Univ Montpellier CNRS)
- ▶ Maximilien Servajean (Paul Valery University, LIRMM, Univ Montpellier CNRS)
- ▶ Axel Dubar (IMAG, Univ Montpellier, CNRS)

PROBLEM: CAN WE TRUST OUR DATA?



(Deep) Learning pipeline with huge labeled dataset (of images):



PROBLEM: CAN WE TRUST OUR DATA?



(Deep) Learning pipeline with huge labeled dataset (of images):



... but labeling errors are common

PROBLEM: CAN WE TRUST OUR DATA?



(Deep) Learning pipeline with huge labeled dataset (of images):



... but labeling errors are common

CIFAR10⁽¹⁾



Given label: cat

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

PROBLEM: CAN WE TRUST OUR DATA?



(Deep) Learning pipeline with huge labeled dataset (of images):



... but labeling errors are common

CIFAR10⁽¹⁾



Given label: cat

Quickdraw⁽²⁾



Given label: T-shirt

⁽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>)

PROBLEM: CAN WE TRUST OUR DATA?



(Deep) Learning pipeline with huge labeled dataset (of images):



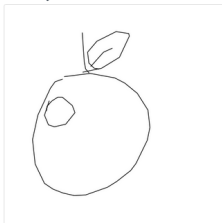
... but labeling errors are common

CIFAR10⁽¹⁾



Given label: cat

Quickdraw⁽²⁾



Given label: T-shirt

MNIST⁽³⁾



Given label: 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". *Proceedings of the IEEE* 86.11, pp. 2278–2324)



► Notation and setting

- Dataset : $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n_{\text{train}} + n_{\text{val}} + n_{\text{test}}} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}} \cup \mathcal{D}_{\text{test}}$
- Splitting : $|\mathcal{D}| = n_{\text{train}} + n_{\text{val}} + n_{\text{test}}$
- Tasks : $x_i \in \mathcal{X}$ (images here)
- Labels : $y_i \in [K] = \{1, \dots, K\}$



► Notation and setting

- Dataset : $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n_{\text{train}} + n_{\text{val}} + n_{\text{test}}} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}} \cup \mathcal{D}_{\text{test}}$
- Splitting : $|\mathcal{D}| = n_{\text{train}} + n_{\text{val}} + n_{\text{test}}$
- Tasks : $x_i \in \mathcal{X}$ (images here)
- Labels : $y_i \in [K] = \{1, \dots, K\}$

► Popular datasets for classification

- CIFAR10⁽⁴⁾
- CIFAR100⁽⁴⁾
- ImageNet⁽⁵⁾
- MNIST⁽⁶⁾
- Quickdraw⁽⁷⁾
- LabelMe⁽⁸⁾

⁽⁴⁾ (A. Krizhevsky and G. Hinton [2009]. *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto)

⁽⁵⁾ (J. Deng et al. [2009]. "ImageNet: A Large-Scale Hierarchical Image Database". *CVPR*)

⁽⁶⁾ (Y. LeCun et al. [1998]. "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* 86.11, pp. 2278–2324)

⁽⁷⁾ ([n.d.]. <https://github.com/googlecreativelab/quickdraw-dataset>)

⁽⁸⁾ (F. Rodrigues and F. Pereira [2018]. "Deep learning from crowds". *AAAI*. vol. 32)



► Notation and setting

- Dataset : $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n_{\text{train}}+n_{\text{val}}+n_{\text{test}}} = \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{val}} \cup \mathcal{D}_{\text{test}}$
- Splitting : $|\mathcal{D}| = n_{\text{train}} + n_{\text{val}} + n_{\text{test}}$
- Tasks : $x_i \in \mathcal{X}$ (images here)
- Labels : $y_i \in [K] = \{1, \dots, K\}$

► Popular datasets for classification

- CIFAR10⁽⁴⁾
- CIFAR100⁽⁴⁾
- ImageNet⁽⁵⁾
- MNIST⁽⁶⁾
- Quickdraw⁽⁷⁾
- LabelMe⁽⁸⁾

⁽⁴⁾ (A. Krizhevsky and G. Hinton [2009]. *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto)

⁽⁵⁾ (J. Deng et al. [2009]. "ImageNet: A Large-Scale Hierarchical Image Database". *CVPR*)

⁽⁶⁾ (Y. LeCun et al. [1998]. "Gradient-based learning applied to document recognition". *Proceedings of the IEEE* 86.11, pp. 2278–2324)

⁽⁷⁾ ([n.d.]. <https://github.com/googlecreativelab/quickdraw-dataset>)

⁽⁸⁾ (F. Rodrigues and F. Pereira [2018]. "Deep learning from crowds". *AAAI*. vol. 32)

CIFAR10

A SIMPLE DATASET EXAMPLE FOR MODERN DEEP LEARNING⁽⁹⁾



⁽⁹⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

CIFAR10

A SIMPLE DATASET EXAMPLE FOR MODERN DEEP LEARNING⁽⁹⁾



- ▶ $K = 10$ classes
- ▶ $n_{\text{train}} + n_{\text{val}} = 50\,000$
- ▶ x_i : 32×32 RGB images
- ▶ $n_{\text{test}} = 10\,000$

⁽⁹⁾ A. Krizhevsky and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.

DATASET CONSTRUCTION

HOW DO WE CREATE SUCH A DATASET?



Questions:

- ▶ Where do the tasks come from?

DATASET CONSTRUCTION

HOW DO WE CREATE SUCH A DATASET?



Questions:

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

DATASET CONSTRUCTION

HOW DO WE CREATE SUCH A DATASET?



Questions:

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

DATASET CONSTRUCTION

HOW DO WE CREATE SUCH A DATASET?



Questions:

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

DATASET CONSTRUCTION

HOW DO WE CREATE SUCH A DATASET?



Questions:

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

Notation:

- ▶ Tasks: $\mathcal{X}_{\text{train}} = \{x_1, \dots, x_{n_{\text{task}}}\}$
- ▶ True labels: $(y_i^*)_{i \in [n_{\text{task}}]}$ **unobserved**
- ▶ Workers: $(w_j)_{j \in [n_{\text{worker}}]}$, label some images
- ▶ Label answered by worker w_j for a task x_i : $y_i^{(j)} \in [K]$
- ▶ Annotators set: $\mathcal{A}(x_i) = \{j \in [n_{\text{worker}}] : \text{worker } w_j \text{ labeled task } x_i\}$

$$\mathcal{D}_{\text{train}} = \bigcup_{i=1}^{n_{\text{task}}} \left\{ (x_i, (y_i^{(j)})) \text{ for } j \in \mathcal{A}(x_i) \right\}$$

CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)

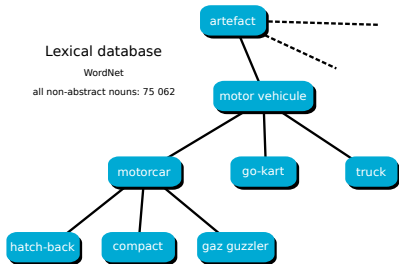


CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)

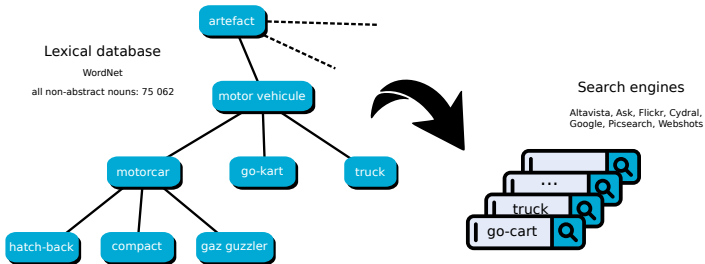


Lexical database
WordNet
all non-abstract nouns: 75 062



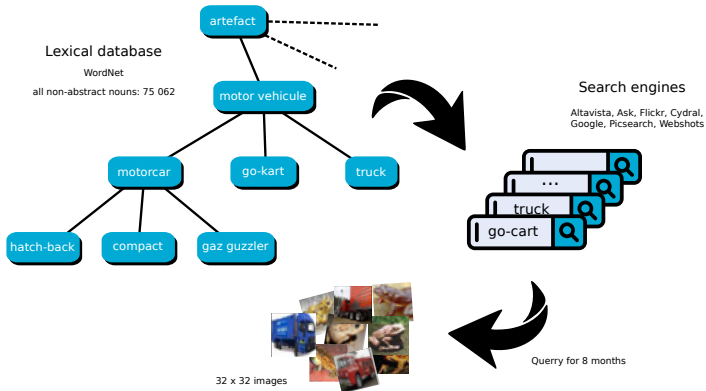
CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)



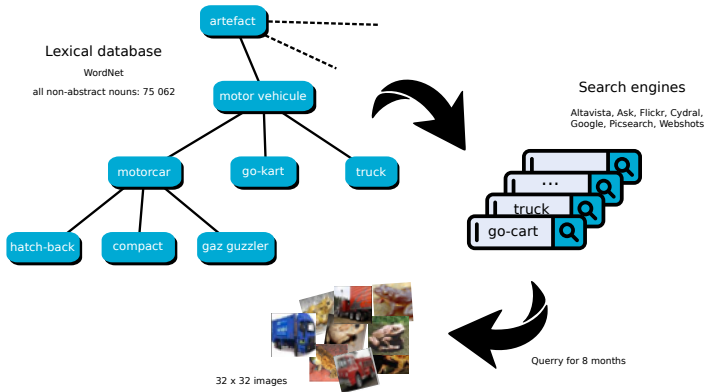
CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)



CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 1: DATA COLLECTION (80 MILLION TINY IMAGES)

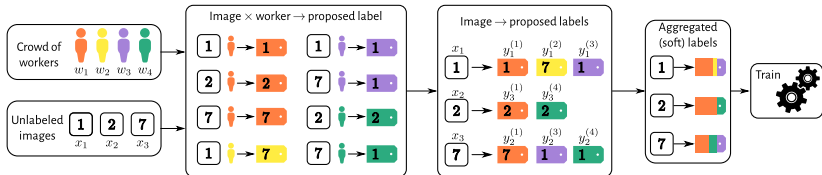


Many issues raised⁽¹⁰⁾: opacity, anonymity (face search/reverse image search), perpetuate stereotypes, etc.

⁽¹⁰⁾V. Uday Prabhu and A. Birhane (June 2020). "Large image datasets: A pyrrhic win for computer vision?" *arXiv e-prints*, arXiv:2006.16923.

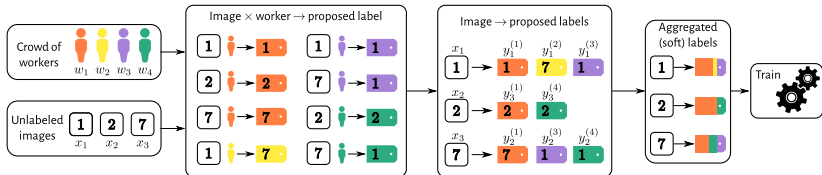
CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING



CIFAR10, AN ARCHETYPAL EXAMPLE

STEP 2: LABEL COLLECTION AND CROWDSOURCING

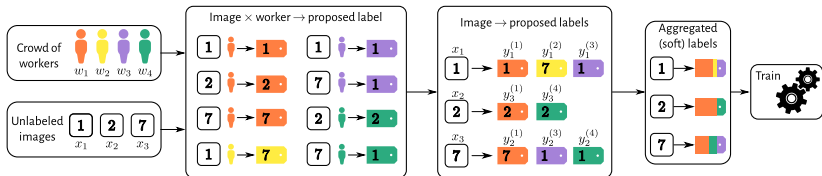


Krizhevsky and Hinton⁽¹¹⁾ :

⁽¹¹⁾ 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



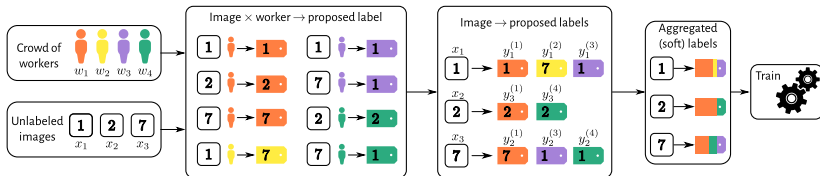
Krizhevsky and Hinton⁽¹¹⁾ :

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

⁽¹¹⁾ 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



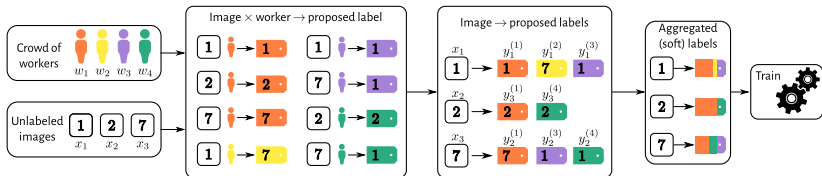
Krizhevsky and Hinton⁽¹¹⁾ :

- ▶ "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."

⁽¹¹⁾ 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



Krizhevsky and Hinton⁽¹¹⁾ :

- ▶ "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*

⁽¹¹⁾ 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 CIFAR-10H behavioral dataset consists of **511 400** human categorization decisions over the **10 000**-image testing subset of CIFAR10 (approx. 50 judgments per image)."

- ▶ Total number of workers: $n_{\text{worker}} = 2571$ (via Amazon Mechanical Turk)
- ▶ **Processing:** (After an initial training phase) 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, according to the authors...we could not reproduce that).

Note: workers were paid \$1.50 (average completion time ≈ 5 mn); poor worker conditions⁽¹²⁾

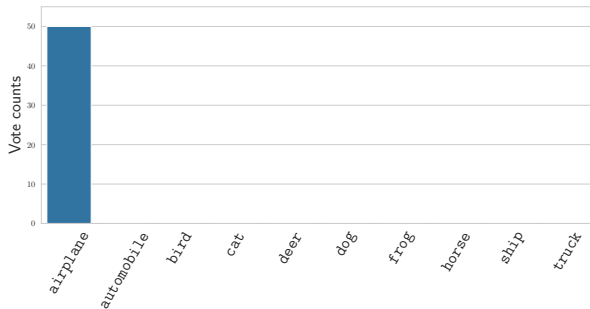
For learning: we will consider $n_{\text{train}} = 9500$ and $n_{\text{val}} = 500$

⁽¹²⁾ <https://time.com/6247678/openai-chatgpt-kenya-workers/>

⁽¹³⁾ J. C. Peterson et al. (2019). "Human Uncertainty Makes Classification More Robust". ICCV, pp. 9617–9626.



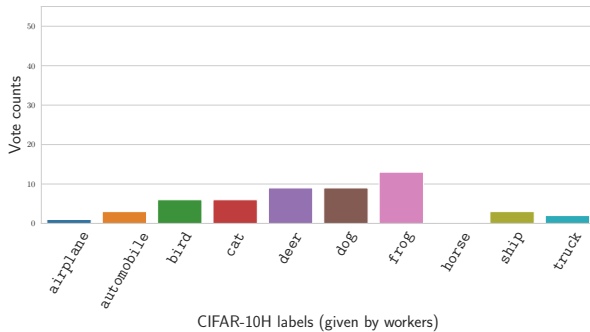
Image # 7681
CIFAR-10 label: airplane



CIFAR-10H labels (given by workers)



Image # 6750
CIFAR-10 label: deer



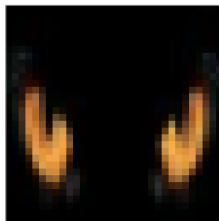
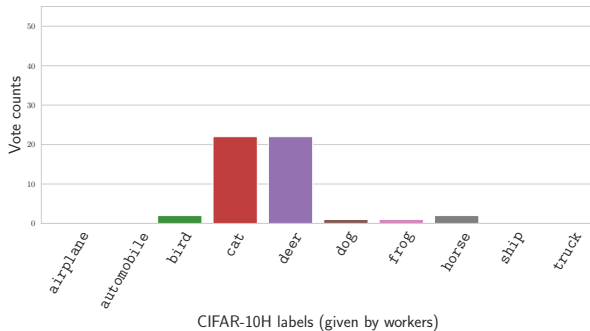


Image # 9246
CIFAR-10 label: cat



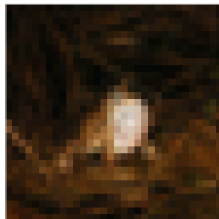


Image # 3724
CIFAR-10 label: frog

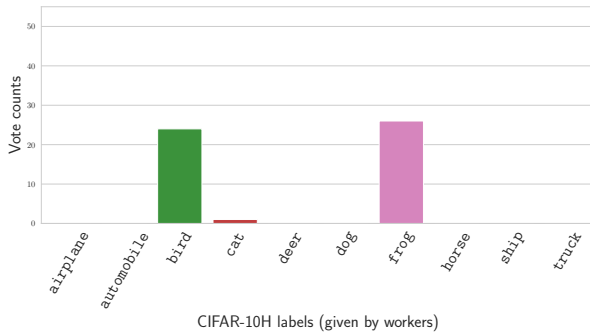




Image # 1353
CIFAR-10 label: cat

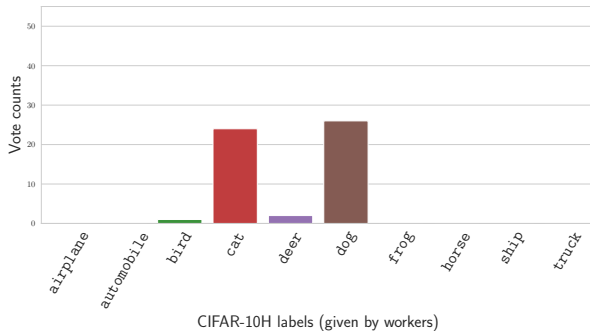
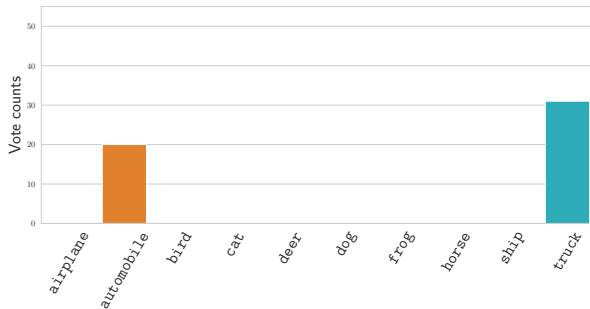




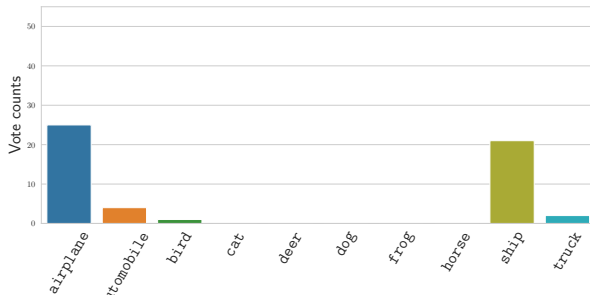
Image # 7455
CIFAR-10 label: automobile



CIFAR-10H labels (given by workers)



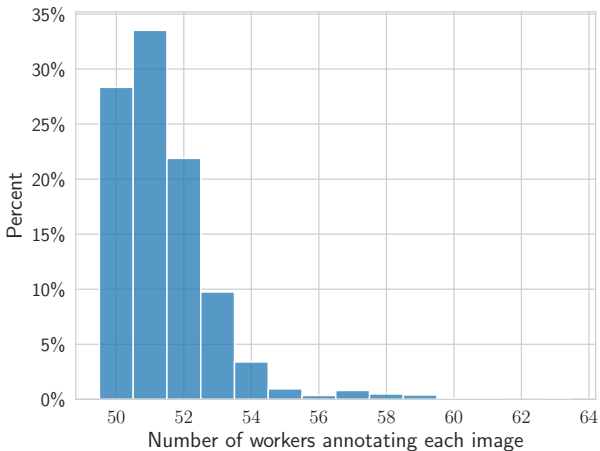
Image # 8872
CIFAR-10 label: ship



CIFAR-10H labels (given by workers)

CIFAR-10H: DATASET VISUALIZATION

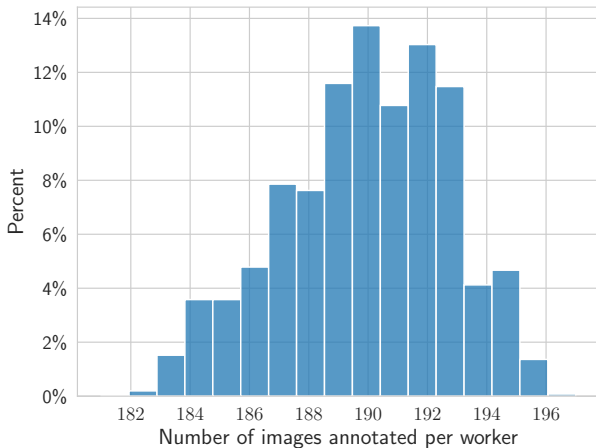
STATISTICS ON OUR TRAINING SET ($n_{\text{train}} = 9\,500$)



Feedback effort per task distribution

CIFAR-10H: DATASET VISUALIZATION

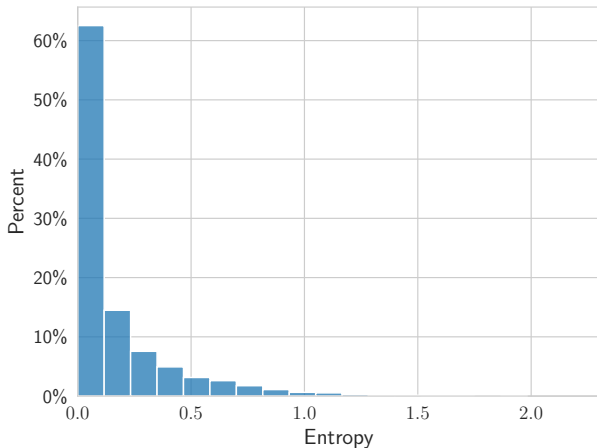
STATISTICS ON OUR TRAINING SET ($n_{\text{train}} = 9500$)



Load per worker distribution

CIFAR-10H: DATASET VISUALIZATION

STATISTICS ON OUR TRAINING SET ($n_{\text{train}} = 9500$)



Naive soft labels, entropy distribution



Definition: Majority Voting (MV)

Majority Voting outputs the most answered label:

$$\forall x_i \in \mathcal{X}_{\text{train}}, \quad \hat{y}_i^{\text{MV}} = \arg \max_{k \in [K]} \left(\sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_i^{(j)}=k\}} \right)$$

Properties:

- ✓ simple
- ✓ adapted for any number of workers
- ✓ usually efficient, often few labelers sufficient (say⁽¹⁴⁾ <5)
- ✗ ineffective for borderline cases
- ✗ suffer from spammers / adversarial workers

⁽¹⁴⁾ R. Snow et al. (2008). "Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks". *Conference on Empirical Methods in Natural Language Processing*. EMNLP 2008. Association for Computational Linguistics, pp. 254–263.

**Definition: Weighted Majority Voting (WMV)**

Majority voting but weighted by a confidence score per worker w_j :

$$\forall x_i \in \mathcal{X}_{\text{train}}, \quad \hat{y}_i^{\text{WMV}} = \arg \max_{k \in [K]} \left(\sum_{j \in \mathcal{A}(x_i)} \alpha_j \mathbb{1}_{\{y_i^{(j)} = k\}} \right)$$

$\alpha_j > 0$: reflects the confidence in worker w_j

- ✓ simple
- ✓ adapted for any number of workers
- ✓ usually efficient
- ✓ can leverage expert workers
- ✗ ineffective for borderline cases
- ✗ suffer from spammers / adversarial workers
- ✗ requires prior knowledge of the workers



Notation : for $z \in \mathbb{R}_+^d, \forall i \in [d], \text{Norm}(z)_i = \frac{z_i}{\sum_{i'=1}^d z_{i'}}$

Definition: Naive Soft (NS) labels

Naive soft outputs the empirical distribution of the answered votes:

$$\forall x_i \in \mathcal{X}_{\text{train}}, \hat{y}_i^{\text{NS}} = \text{Norm}(\tilde{y}_i), \quad \text{where } \tilde{y}_i = \left(\sum_{j \in \mathcal{A}(x_i)} \mathbb{1}_{\{y_i^{(j)}=k\}} \right)_{k \in [K]}$$

- ✓ simple
- ✓ adapted for any number of workers
- ✓ can reflect workers variability & task ambiguity
- ✗ suffer from spammers/adversarial workers

Dawid and Skene⁽¹⁵⁾ (DS)

Assumption: each worker answers independently

The j -th worker has his own **confusion matrix**: $\pi^{(j)} \in \mathbb{R}^{K \times K}$

$$\pi_{\ell,k}^{(j)} = \mathbb{P}(y_i^{(j)} = k | y_i^* = \ell)$$

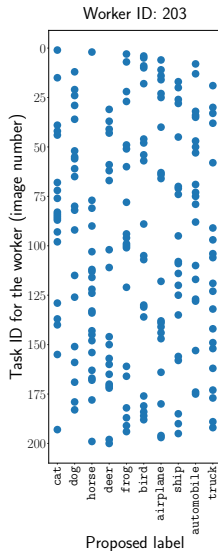
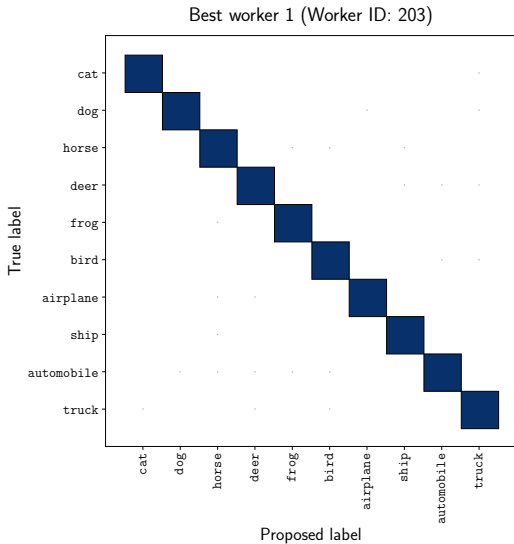
Conditionally on the true label, the j -th worker answers as follows:

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

⁽¹⁵⁾ 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.

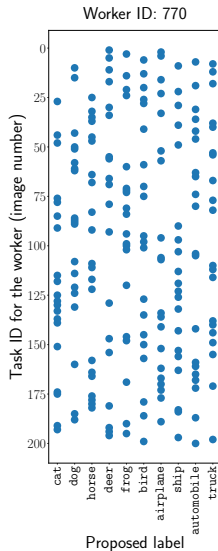
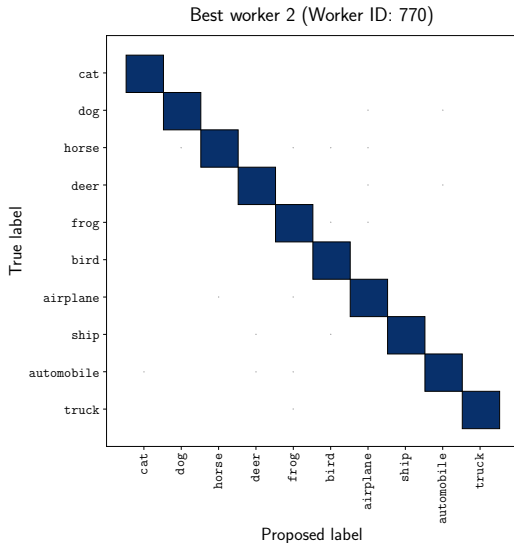
(ESTIMATED) CONFUSION MATRICES

ILLUSTRATION AND INTERPRETATION



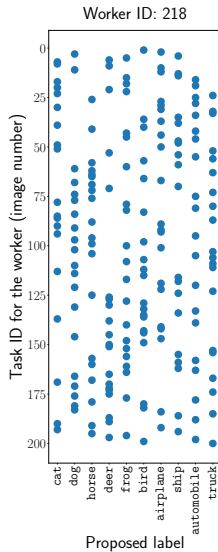
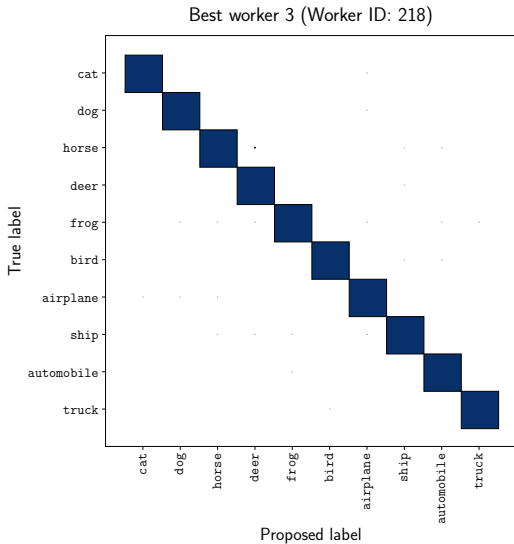
(ESTIMATED) CONFUSION MATRICES

ILLUSTRATION AND INTERPRETATION



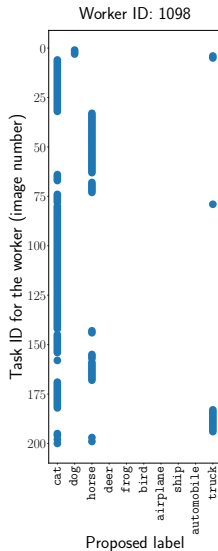
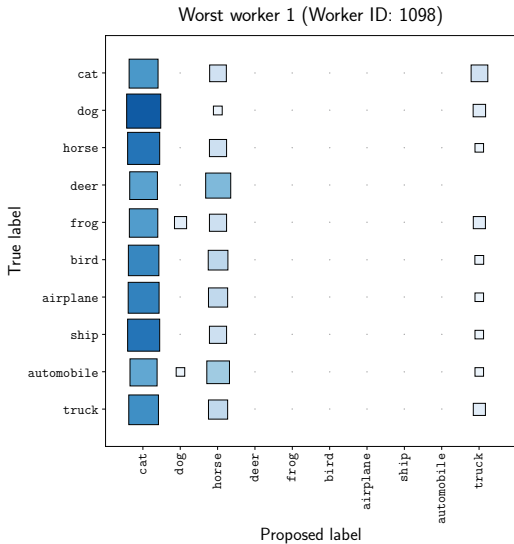
(ESTIMATED) CONFUSION MATRICES

ILLUSTRATION AND INTERPRETATION



(ESTIMATED) CONFUSION MATRICES

ILLUSTRATION AND INTERPRETATION

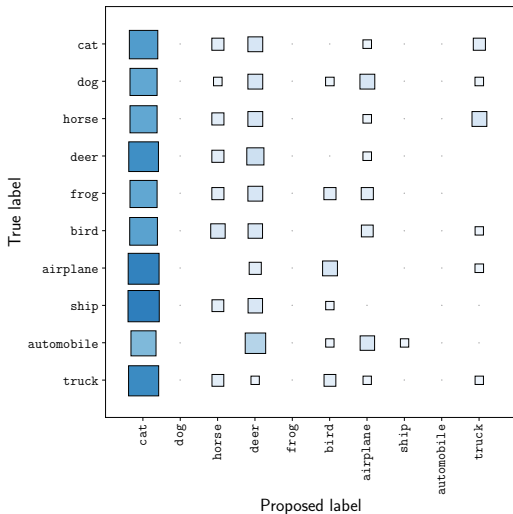


(ESTIMATED) CONFUSION MATRICES

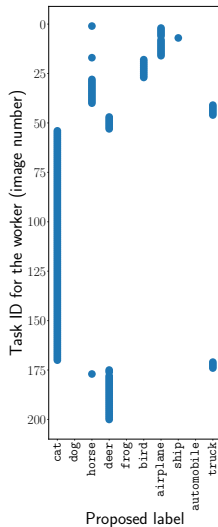
ILLUSTRATION AND INTERPRETATION



Worst worker 2 (Worker ID: 2160)

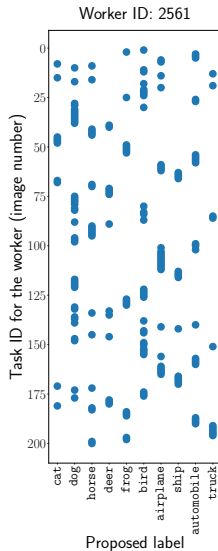
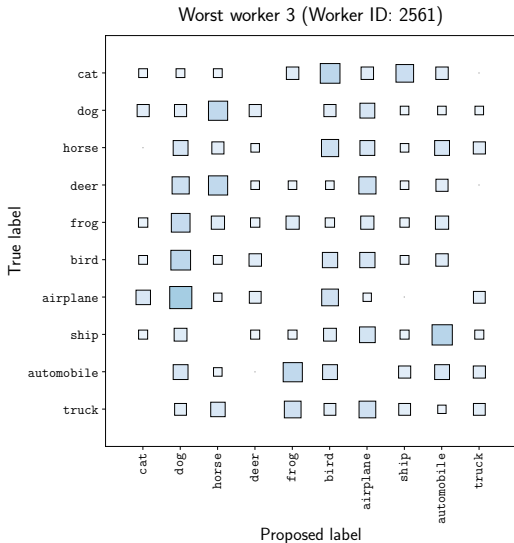


Worker ID: 2160



(ESTIMATED) CONFUSION MATRICES

ILLUSTRATION AND INTERPRETATION



Likelihood:

$$\prod_{k \in [K]} (\pi_{\ell, k}^{(j)})^{\mathbb{1}_{\{y_i^{(j)} = k\}}}$$

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

Likelihood:

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

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

Likelihood:

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

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



Likelihood:

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

- Multinomial with 1 task, 1 worker and 1 answer conditioned on $y_i^* = \ell$
- Multiple workers answer independently
- Remove conditioning assumption on y_i^* : $\mathbb{P}(y_i^* = \ell) = \rho_{\ell}$ (**prevalence**)
- Tasks independence and $T_{i, \ell} = \mathbb{1}_{\{y_i^* = \ell\}}$ (1 if task i has true label ℓ , 0 otherwise)



Likelihood:

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



Likelihood:

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

Annotations:

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

Likelihood:

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

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 $\pi_{\ell, k}^{(j)}$)

a) Estimate $\rho \in \Delta_{K-1} := \{p \in \mathbb{R}^K, \sum_{k=1}^K p_k = 1, p_k \geq 0\}$ assuming **known** $T_{i, \ell}$ s and the constraints $\sum_{\ell \in [K]} T_{i, \ell} = 1$ for all i

$$\hat{\rho} \in \arg \max_{\rho \in \Delta_{K-1}} \left(\log \prod_{i \in [n_{\text{task}}]} \prod_{\ell \in [K]} \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}} \right]^{T_{i, \ell}} \right)$$

$$\iff \hat{\rho} \in \arg \max_{\rho \in \Delta_{K-1}} \sum_{i \in [n_{\text{task}}]} \sum_{\ell \in [K]} T_{i, \ell} \log \left[\rho_{\ell} \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}} \right]$$

$$\iff \hat{\rho} \in \arg \max_{\rho \in \Delta_{K-1}} \sum_{i \in [n_{\text{task}}]} \sum_{\ell \in [K]} T_{i, \ell} \log(\rho_{\ell})$$

$$\iff \hat{\rho} = \frac{1}{n_{\text{task}}} \sum_{i \in [n_{\text{task}}]} \hat{T}_{i, :} \quad (\text{use Lagrange multipliers to get the solution})$$

Likelihood:

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

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 $\pi_{\ell, k}^{(j)}$)

b) Estimate $\pi_{\ell, :}^{(j)} \in \Delta_{K-1}$ assuming **known** $T_{i, \ell}$ s and the constraints $\sum_{\ell \in [K]} T_{i, \ell} = 1$ for all i

$$\hat{\pi}_{\ell, :}^{(j)} \in \arg \max_{\pi^{(j)} \in \Delta_{K-1}} \left(\log \prod_{i \in [n_{\text{task}}]} \prod_{\ell' \in [K]} \left[\rho_{\ell'} \prod_{j' \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell', k}^{(j')} \right)^{\mathbb{1}_{\{y_i^{(j')} = k\}}} \right]^{T_{i, \ell'}} \right)$$

$$\iff \hat{\pi}_{\ell, :}^{(j)} \in \arg \max_{\pi^{(j)} \in \Delta_{K-1}} \sum_{i \in [n_{\text{task}}]} T_{i, \ell} \log \left[\rho_{\ell} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}} \right]$$

$$\iff \hat{\pi}_{\ell, :}^{(j)} \in \arg \max_{\pi^{(j)} \in \Delta_{K-1}} \sum_{i \in [n_{\text{task}}]} \sum_{k \in [K]} T_{i, \ell} \cdot \mathbb{1}_{\{y_i^{(j)} = k\}} \log(\pi_{\ell, k}^{(j)})$$

$$\iff \hat{\pi}_{\ell, :}^{(j)} = \sum_{i \in [n_{\text{task}}]} T_{i, \ell} \cdot \mathbb{1}_{\{y_i^{(j)} = :\}} / \sum_{i \in [n_{\text{task}}]} \sum_{k' \in [K]} T_{i, \ell} \cdot \mathbb{1}_{\{y_i^{(j)} = k'\}}$$

Likelihood:

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

Annotations:

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

c) Estimate $T_{i, \ell}$ s as probabilities with ρ and $\pi^{(j)}$ s known, with the constraints $\sum_{\ell \in [K]} T_{i, \ell} = 1$ for all i

$$\begin{aligned} \hat{T}_{i, \ell} &= \mathbb{P}(T_{i, \ell} = 1 | \mathcal{D}_{\text{train}}) \\ &\propto \mathbb{P}(\mathcal{D}_{\text{train}} | T_{i, \ell} = 1) \mathbb{P}(T_{i, \ell} = 1) \\ &\propto \prod_{j \in [n_{\text{worker}}]} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}} \cdot \rho_{\ell} \\ &\propto \prod_{j \in \mathcal{A}(x_i)} \prod_{k \in [K]} \left(\pi_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}} \cdot \rho_{\ell} \end{aligned}$$



Likelihood:

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

Annotations:

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

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:

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

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

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**

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

Likelihood:

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

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

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 $\hat{\rho}$ and $\hat{\pi}$ assuming \hat{T} s are known

$$3 \quad \forall \ell \in [K], \quad \hat{\rho}_{\ell} \leftarrow \frac{1}{n_{\text{task}}} \sum_{i \in [n_{\text{task}}]} \hat{T}_{i, \ell}$$

$$4 \quad \forall (\ell, k) \in [K]^2, \quad \hat{\pi}_{\ell, k}^{(j)} \leftarrow \frac{\sum_{i \in [n_{\text{task}}]} \hat{T}_{i, \ell} \cdot \mathbb{1}_{\{y_i^{(j)} = k\}}}{\sum_{k' \in [K]} \sum_{i' \in [n_{\text{task}}]} \hat{T}_{i', \ell} \cdot \mathbb{1}_{\{y_{i'}^{(j)} = k'\}}}$$

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

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

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

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 $\hat{\rho}$ and $\hat{\pi}$ assuming \hat{T} s are known

$$3 \quad \forall \ell \in [K], \quad \hat{\rho}_{\ell} \leftarrow \frac{1}{n_{\text{task}}} \sum_{i \in [n_{\text{task}}]} \hat{T}_{i, \ell}$$

$$4 \quad \forall (\ell, k) \in [K]^2, \quad \hat{\pi}_{\ell, k}^{(j)} \leftarrow \frac{\sum_{i \in [n_{\text{task}}]} \hat{T}_{i, \ell} \cdot \mathbb{1}_{\{y_i^{(j)} = k\}}}{\sum_{k' \in [K]} \sum_{i' \in [n_{\text{task}}]} \hat{T}_{i', \ell} \cdot \mathbb{1}_{\{y_{i'}^{(j)} = k'\}}}$$

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

$$5 \quad \forall (i, \ell) \in [n_{\text{task}}] \times [K], \quad \hat{T}_{i, \ell} \leftarrow \frac{\prod_{j \in \mathcal{A}(x_i)} \prod_{k \in [K]} \hat{\rho}_{\ell} \cdot \left(\hat{\pi}_{\ell, k}^{(j)} \right)^{\mathbb{1}_{\{y_i^{(j)} = k\}}}}{\sum_{\ell' \in [K]} \prod_{j' \in \mathcal{A}(x_i)} \prod_{k' \in [K]} \hat{\rho}_{\ell'} \cdot \left(\hat{\pi}_{\ell', k'}^{(j')} \right)^{\mathbb{1}_{\{y_{i'}^{(j')} = k'\}}}}$$

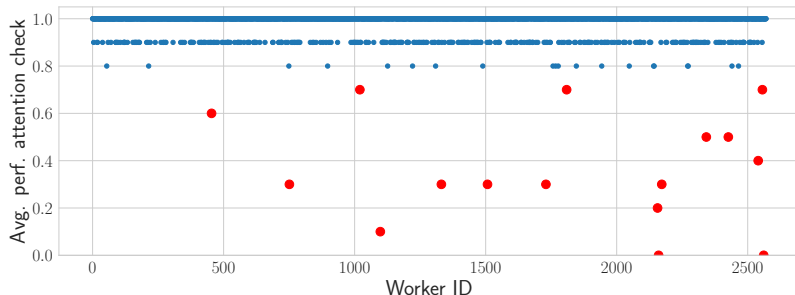
6 Labels: $\forall i \in [n_{\text{task}}], \hat{y}_i = \hat{T}_{i, :}$ (soft label)

SORTING WORKERS BY QUALITY

USE CASE ON CIFAR10H



- ▶ Use attention check / Trapping sets: 10 images per worker (out of 200) whose true label is known \implies get an average score for each worker (red: 16 workers < 0.8)



SORTING WORKERS BY QUALITY

USE CASE ON CIFAR10H



- Use spammer score⁽¹⁶⁾: measure the distance between $\hat{\pi}^j$ and rank 1 matrices (since a spammer has a distribution of answers independent of the true label)

$$\min_{v_j \in \mathbb{R}^K} \left\| \hat{\pi}^{(j)} - \mathbf{1}_K v_j \right\|_F^2$$



⁽¹⁶⁾V. C. Raykar and S. Yu (2011). "Ranking annotators for crowdsourced labeling tasks". *NeurIPS*, pp. 1809–1817.

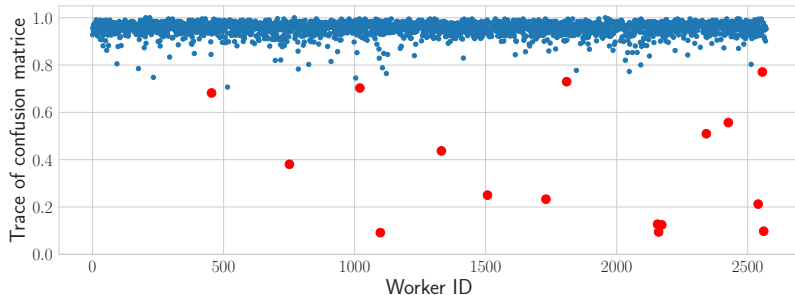
SORTING WORKERS BY QUALITY

USE CASE ON CIFAR10H



- ▶ Use DS: diagonal elements of $\hat{\pi}^{(j)}$ represents worker ability to be correct, get the average success across all labels with

$$\frac{1}{K} \text{trace}(\hat{\pi}^{(j)})$$

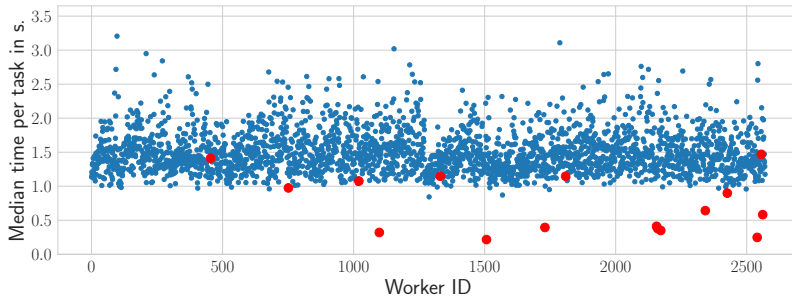


SORTING WORKERS BY QUALITY

USE CASE ON CIFAR10H



- ▶ Use time spent: get the median time spent per task





More to come after a short break



Contact:

Joseph Salmon

✉ `joseph.salmon@umontpellier.fr`









🌐 `https://josephsalmon.eu`



Github: @josephsalmon



Mastodon: @josephsalmon@sigmoid.social



-  (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”. *J. R. Stat. Soc. Ser. C. Appl. Stat.* 28.1, pp. 20–28.
-  Deng, J. et al. (2009). “ImageNet: A Large-Scale Hierarchical Image Database”. *CVPR*.
-  Krizhevsky, A. and G. Hinton (2009). *Learning multiple layers of features from tiny images*. Tech. rep. University of Toronto.
-  LeCun, Y. et al. (1998). “Gradient-based learning applied to document recognition”. *Proceedings of the IEEE* 86.11, pp. 2278–2324.
-  Peterson, J. C. et al. (2019). “Human Uncertainty Makes Classification More Robust”. *ICCV*, pp. 9617–9626.
-  Raykar, V. C. and S. Yu (2011). “Ranking annotators for crowdsourced labeling tasks”. *NeurIPS*, pp. 1809–1817.
-  Rodrigues, F. and F. Pereira (2018). “Deep learning from crowds”. *AAAI*. Vol. 32.

-  Snow, R. et al. (2008). “Cheap and Fast - But is it Good? Evaluating Non-Expert Annotations for Natural Language Tasks”. *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?” *arXiv e-prints*, arXiv:2006.16923.