



# Semi-Supervised Aggregation of Dependent Weak Supervision Sources With Performance Guarantees

Alessio Mazzetto\*, Dylan Sam, Andrew Park, Eli Upfal, Stephen H. Bach

\*alessio\_mazzetto@brown.edu



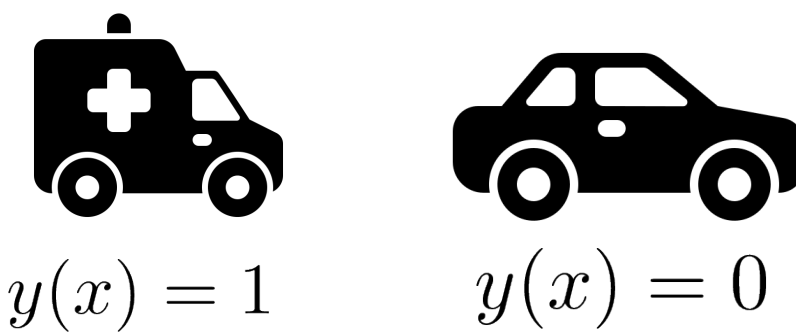
## Introduction

Binary classification

$$y : \mathcal{X} \rightarrow \{0, 1\}$$

Classification domain with distribution  $\mathcal{D}$

Example: classify ambulance.



Given hypothesis class  $\mathcal{H}$ , we want to find the hypothesis  $h$  s.t.

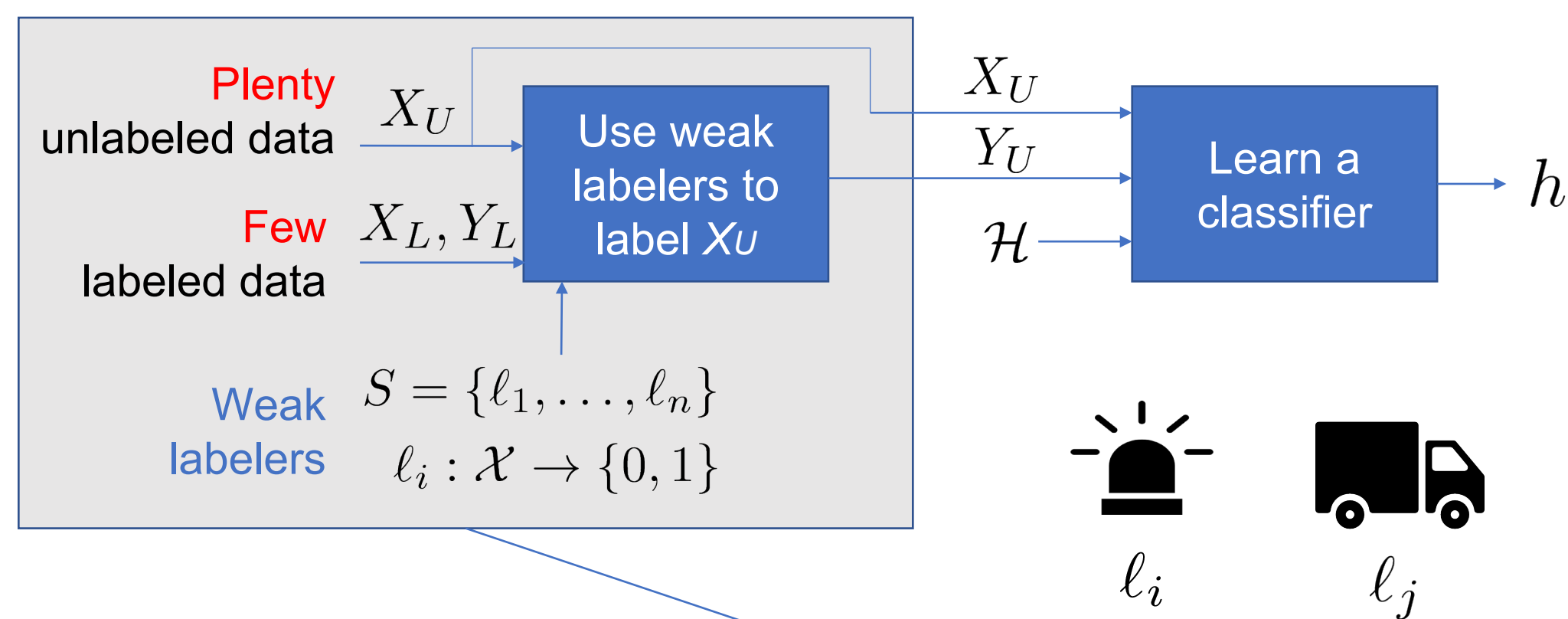
$$\min_{h \in \mathcal{H}} \varepsilon(h) := \min_{h \in \mathcal{H}} \Pr_{x \sim \mathcal{D}} (y(x) \neq h(x))$$

error of  $h$

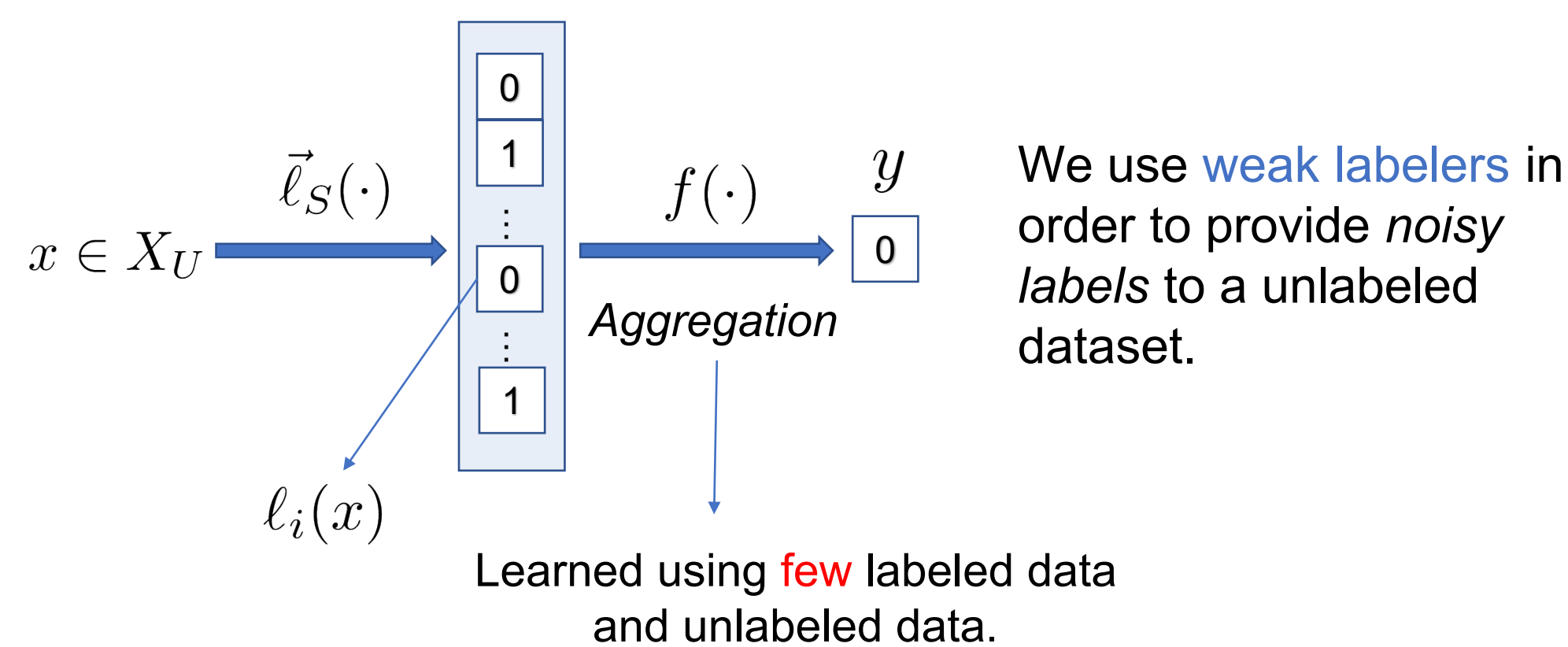
Supervised learning of a binary classification task requires a lot of labeled data for high-dimensional hypothesis classes (e.g., DNN).

Labeled data is **costly** and **scarce** for a lot of binary classification task of interest.

Weak Supervision Framework [1]



Our focus



## Contribution

Previous work unrealistically usually assumes *independence* or a *distribution family* between weak labelers' errors to do aggregation.

**Our contribution:**

- **First theoretical bound** to the worst-case error of the **majority vote** of a set of weak labelers without those assumptions.
- **Novel algorithm** that uses the bound above to provide **the first theoretical guarantees** in learning an aggregation of an arbitrary set of weak labelers.

## Intuition

Preliminary definitions:

- Let error rate of  $i$ -th labeler be  $\epsilon_i = \varepsilon(\ell_i) \rightarrow$  easy to estimate with few labeled data
- Let  $S(\vec{\epsilon})$  be the **set of all set of labelers** that have error rates equal to  $\vec{\epsilon} = (\epsilon_1, \dots, \epsilon_n)$
- Given a vector  $\vec{a} \in \{0, 1\}^n$ , let  $\lambda(\vec{a})$  be **its majority vote**.

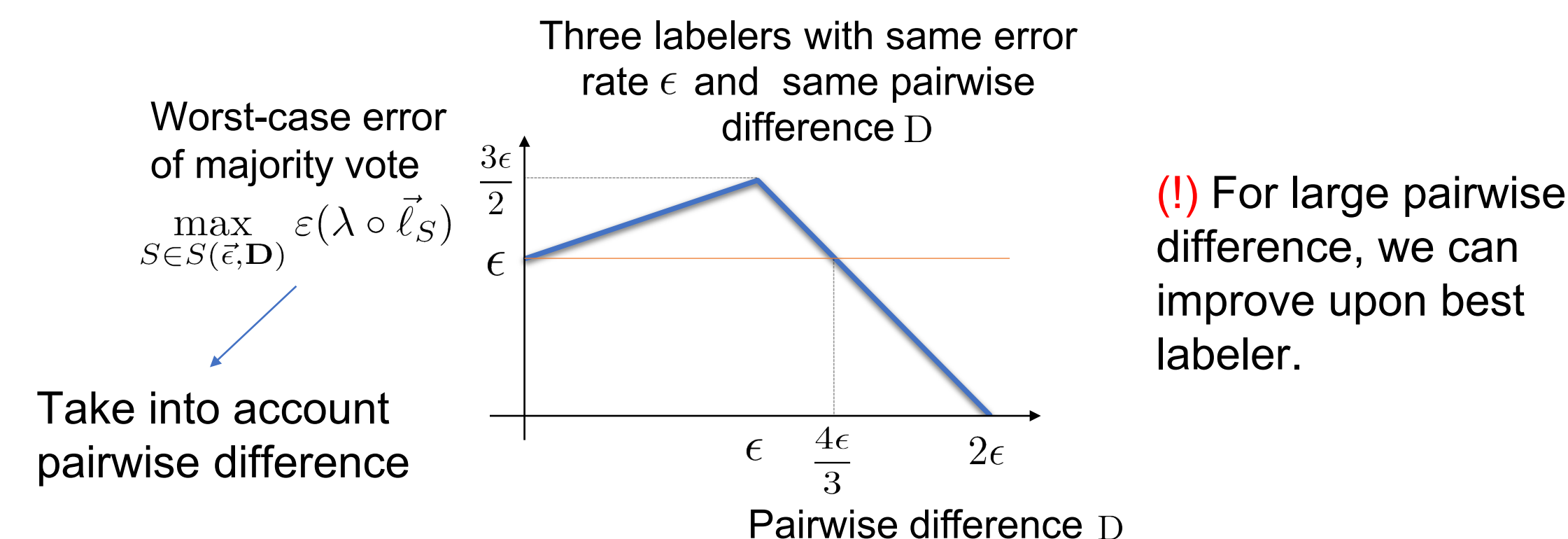
Our result (knowledge of error rates is not enough):

$$\max_{S \in S(\vec{\epsilon})} \varepsilon(\lambda \circ \vec{\ell}_S) \geq \text{median}\{\epsilon_1, \dots, \epsilon_n\}$$

Worst-case, we **cannot improve** upon the best labeler, if we only know their error rates (!) With independence assumption, error of majority vote would go to zero.

We need auxiliary information: **pairwise difference** between labelers.

$$D_{ij} = \Pr_{x \sim \mathcal{D}} (\ell_i(x) \neq \ell_j(x)) \rightarrow \text{only need unlabeled data to estimate it}$$



## Method

**Goal:** find subset of weak labelers with lowest worst-case error on their majority vote.

- **Closed formula** for set of three weak labelers.
- **Heuristic:** iteratively add the two labelers that yield the lowest worst-case error on their majority vote.

For a set of labelers  $S = \{\ell_1, \dots, \ell_n\}$  with error rate  $\vec{\epsilon}$  and pairwise difference  $\mathbf{D}$ , we have that:

$$\max_{S \in S(\vec{\epsilon}, \mathbf{D})} \varepsilon(\lambda \circ \vec{\ell}_S) = \max_{\vec{a} \in \{0, 1\}^n: |\vec{a}|_1 < n/2} \sum p_{\vec{a}}$$

(a)  $\sum_{\vec{a} \in \{0, 1\}^n: a_i = 0} p_{\vec{a}} = \epsilon_i$  for  $i = 1, \dots, n$

(b)  $\sum_{\vec{a} \in \{0, 1\}^n: a_i \neq a_j} p_{\vec{a}} = D_{ij}$  for  $i \neq j$

(c)  $\sum_{\vec{a}} p_{\vec{a}} = 1$  **Linear program** with  $O(2^n)$  variables and  $O(n^2)$  constraints

(d)  $p_{\vec{a}} \geq 0 \quad \forall \vec{a}$

## Experiments

Animals With Attribute (AwA2 [2]) dataset. Each class has 85 attributes, used to create weak classifiers.

Dataset	Baselines		State-of-the-art [3]		Our methods	
	Majority Vote	Dawid-Skene	ALL	PGMV	PGMV-P	PGMV-D
AwA2 (1)	79.1 ± 1.1	80.0 ± 1.8	84.2 ± 0.9	82.0 ± 1.1	85.5 ± 0.9	84.3 ± 1.3
AwA2 (2)	90.0 ± 0.7	94.7 ± 0.4	93.5 ± 0.5	93.7 ± 0.4	93.7 ± 0.5	94.1 ± 0.4
AwA2 (3)	92.3 ± 1.0	96.7 ± 0.3	95.5 ± 0.5	95.4 ± 0.3	95.9 ± 0.3	96.3 ± 0.2
AwA2 (4)	94.2 ± 0.6	96.8 ± 0.2	93.8 ± 0.8	96.8 ± 0.2	97.0 ± 0.3	96.8 ± 0.2
AwA2 (5)	97.6 ± 0.6	99.0 ± 0.2	96.3 ± 0.7	97.5 ± 0.3	98.3 ± 0.3	98.8 ± 0.2

Accuracy over different tasks, grouped by quality of the weak labelers, using ~800 unlabeled and labeled data.

References:

- [1]: Ratner, A., Bach, S. H., Ehrenberg, H., Fries, J., Wu, S., and Ré, C. (2017). Snorkel: Rapid training data creation with weak supervision. PVLDB.
- [2]: Xian, Y., Lampert, C. H., Schiele, B., and Akata, Z. (2018). Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. PAMI
- [3]: Arachic, C. and Huang, B. (2019). Adversarial label learning. AAAI.

Images from flaticon.com