

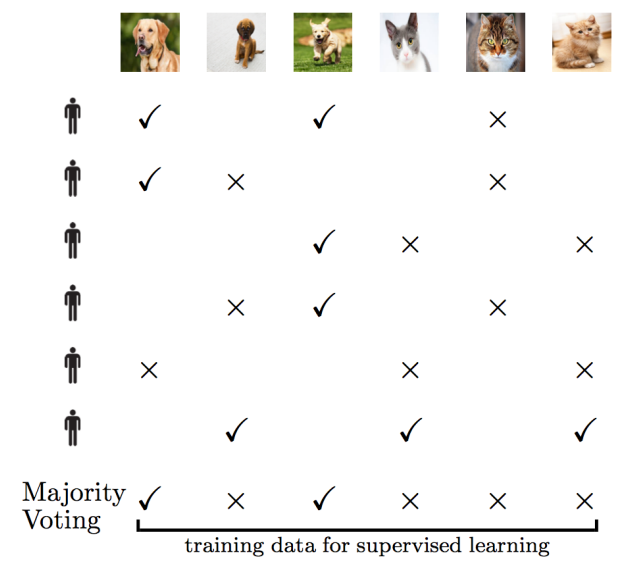
LEARNING FROM NOISY SINGLY-LABELED DATA

ASHISH KHETAN

ZACHARY C. LIPTON

ANIMASHREE ANANDKUMAR

CROWDSOURCING: NOISY ANNOTATIONS



PROBLEM FORMULATION

- n i.i.d. samples $(X, Y) \in (\mathcal{X} \times \mathcal{K}) \sim \mathcal{D}$
- r noisy labels $\{Z_{ij}\}_{j \in [r]}$ on each i -th example X_i
- given by workers $\{w_{ij}\}_{j \in [r]}$, $w_{ij} \in [m]$
- want to learn $\hat{f} \in \mathcal{F}$ such that $\hat{f}(X) = Y$ w.h.p.

DAWID SKENE MODEL (DS)

- Each a -th worker is characterized by its confusion matrix $\pi^{(a)}$
- $\pi^{(a)} \in [0, 1]^{K \times K}$: $\sum_{s \in \mathcal{K}} \pi_{ks} = 1$
- $\mathbb{P}[Z_{ij} = s | Y_i = k, w_{ij} = a] = \pi_{ks}^{(a)}$

LEARNING WITH NOISY LABELS

Posterior probability weighted loss:

$$\ell_{\hat{\pi}}(f(X), Z^{(r)}, w^{(r)}) \equiv \sum_{k \in \mathcal{K}} \mathbb{P}_{\hat{\pi}}[Y = k | Z^{(r)}; w^{(r)}] \ell(f(X), Y = k)$$

MODEL BOOTSTRAPPED EM (MBEM)

Input: data $\{(X_i, Z_i^{(r)}, w_i^{(r)})\}_{i \in [n]}$

Output: deep learning model \hat{f}

Initialize posterior distribution using weighted majority vote

$$\mathbb{P}_{\hat{\pi}}[Y_i = k | Z_i^{(r)}; w_i^{(r)}] \leftarrow (1/r) \sum_{j=1}^r \mathbb{I}[Z_{ij} = k], \text{ for } k \in \mathcal{K}$$

Repeat T times:

learn predictor function by minimizing probability weighted loss $\hat{f} \leftarrow \arg \min_{f \in \mathcal{F}} \frac{1}{n} \ell_{\hat{\pi}}(f(X_i), Z_i^{(r)}, w_i^{(r)})$

predict on the training examples

$$t_i \leftarrow \arg \max_{k \in \mathcal{K}} \hat{f}(X_i)_k, \text{ for } i \in [n]$$

estimate confusion matrices $\hat{\pi}$ given model predictions $\{t_i\}_{i \in [n]}$

$\hat{\pi}^{(a)} \leftarrow$ MLE under the DS model assuming $\{t_i\}$ are true labels, $a \in [m]$

estimate label posterior distribution given $\hat{\pi}$

$\mathbb{P}_{\hat{\pi}}[Y_i = k | Z_i^{(r)}; w_i^{(r)}] \leftarrow$ MLE under the DS model assuming $\hat{\pi}$ are true confusion matrices, for $k \in \mathcal{K}, i \in [n]$

Return \hat{f}

ℓ -RISK UNDER \mathcal{D}

Let $\ell(f(X), Y)$ denote a loss function.

$$R_{\ell, \mathcal{D}}(f) \triangleq \mathbb{E}_{(X, Y) \sim \mathcal{D}} [\ell(f(X), Y)]$$

MAIN THEOREM

- $N \triangleq nr$. For any hypothesis class \mathcal{F} with a finite VC dimension V , and binary classification with 0-1 loss ℓ .

- There exists a universal constant C such that for any $\delta < 1$, if N is large enough (characterized in the paper) then \hat{f} returned by the MBEM algorithm after $T = 2$ iterations satisfies

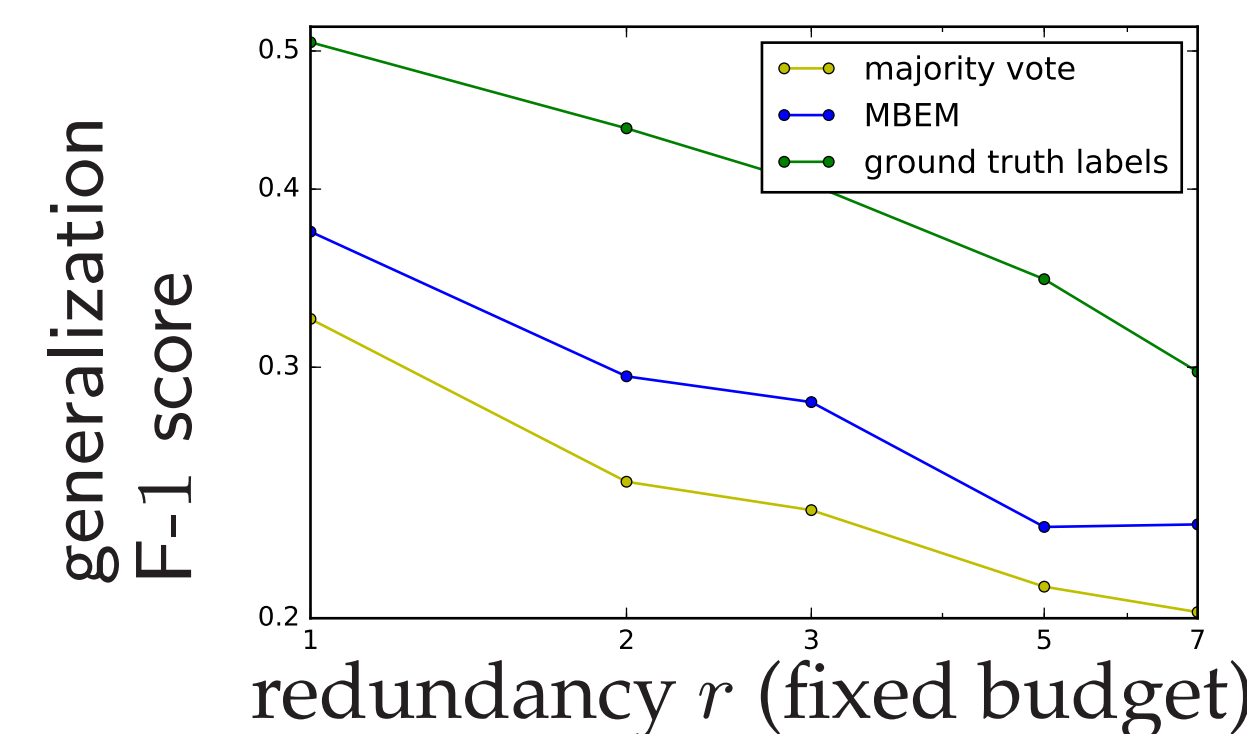
$$R_{\ell, \mathcal{D}}(\hat{f}) - \min_{f \in \mathcal{F}} R_{\ell, \mathcal{D}}(f) \leq \frac{C\sqrt{r}}{1 - 2g(\rho, r)} \left(\sqrt{\frac{V}{N}} + \sqrt{\frac{\log(1/\delta)}{N}} \right).$$

- $g(\rho, r)$ is an analytical function of worker quality ρ and redundancy r .
- If ρ is above a threshold then $\arg \min_{r \in \mathbb{N}} \frac{\sqrt{r}}{1 - 2g(\rho, r)} = 1$.

Labeling once is optimal. It is also seen in all the experiments.

MS-COCO: REAL ANNOTATIONS

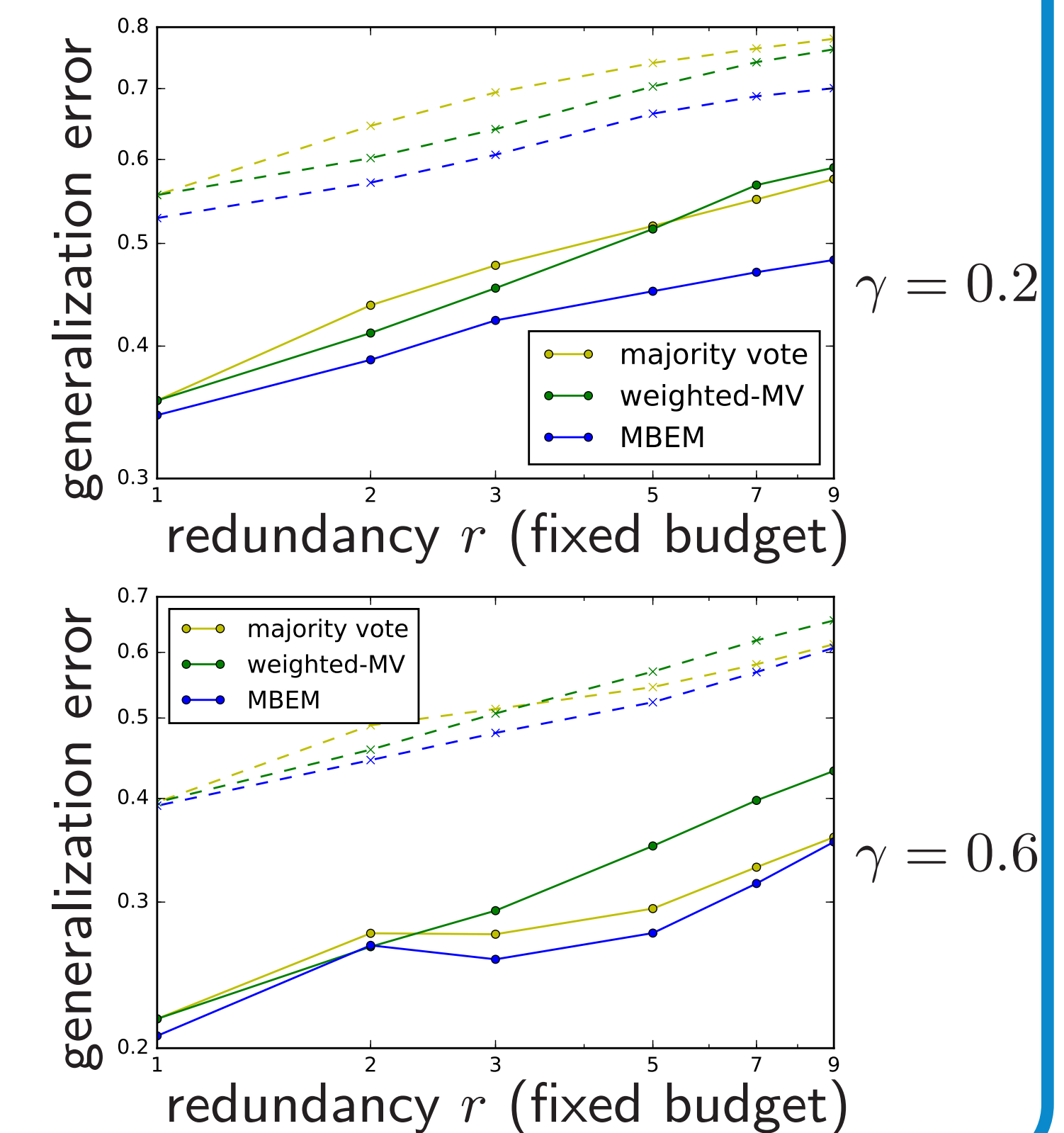
Labeling Once is Optimal.



IMAGENET 1K: SIMULATED WORKERS

- class-wise hammer-spammer workers: Always correct with probability γ for each class independently.

Labeling Once is Optimal.



CIFAR10: SIMULATED WORKERS

class-wise hammer spammer workers

Labeling Once is Optimal.

