Decoupling "when to update" from "how to update"

Eran Malach and Shai Shalev-Shwartz School of Computer Science, The Hebrew University, Israel



Introduction

- Deep learning requires enormous amounts of labeled data
- Labeling is often done by crowdsourcing or data-mining a fast labeling process which does not require manual work
- This usually results in noisy labels
- Using data with noisy labels can crucially fail the training process:
- Updates are performed (mostly) on examples with wrong predictions
- As the network improves most of its updates will be on noisy examples, thus further improvement is impossible

Method

- Solution idea: decoupling the decision of "when to update" from the update rule itself
- "When to update": train two classifiers and update only when they disagree
- ightharpoonup Meta-algorithm for training a classifier from hypothesis class ${\cal H}$ on data sampled from noisy distribution $\hat{\mathcal{D}}$ using update rule \mathbf{U}

$$\begin{split} &\text{for } t=1,2,\ldots,N \ \, \text{do} \\ &\text{draw mini-batch } (x_1,y_1),\ldots,(x_b,y_b) \sim \tilde{\mathcal{D}}^b \\ &\text{let } S=\{(x_i,y_i):h_1(x_i)\neq h_2(x_i)\} \\ &h_1\leftarrow \text{U}(h_1,S) \\ &h_2\leftarrow \text{U}(h_2,S) \\ &\text{end for} \end{split}$$

Theoretical Analysis

- lacksquare Suppose $m{\mathcal{D}}$ is a linearly separable distribution, and $m{\mathcal{D}}$ is the same distribution with labels flipped with probability $\mu < 0.5$
- lacksquare Denote \mathbf{w}^* the linear separator: $\mathbf{y}\langle\mathbf{x},\mathbf{w}^*
 angle \geq \mathbf{1}$ for $(\mathbf{x},\mathbf{y}) \sim \mathcal{D}$
- ► Initialize two linear classifiers with norm < K
- Perceptron as the base classifier, update rule: $\mathbf{w}_{t}^{(i)} = \mathbf{w}_{t-1}^{(i)} + \tilde{\mathbf{y}}_{t}\mathbf{x}_{t}$
- The algorithm converges within a number of iterations that is similar up to a constant factor to that of a perceptron on noise-less data

Theorem 1. Let **T** be the number of updates performed by the algorithm. Then, $\mathbb{E}[T] \leq \frac{3(4K+1)}{(1-2\mu)^2} \|\mathbf{w}^*\|^2$

- Does not imply convergence to an optimum, as the algorithm can get "stuck" before reaching an optimum
- Empirical evaluation shows that this can be avoided

Proof Idea

- The proof is essentially similar to the classical perceptron convergence proof
- ightharpoonup Denote $\mathbf{w}_{t}^{(1)}, \mathbf{w}_{t}^{(2)}$ the weight vectors of the two classifiers after \mathbf{t} updates
- By the disagreement rule: $\operatorname{sign}\langle w_{t-1}^{(1)}, x_t \rangle \neq \operatorname{sign}\langle w_{t-1}^{(2)}, x_t \rangle$
- Upper bound $||\mathbf{w}_{t}^{(i)}||$:

$$\begin{aligned} \|\mathbf{w}_{t}^{(i)}\|^{2} &= \|\mathbf{w}_{t-1}^{(i)}\|^{2} + 2\mathbf{y}_{t}\langle\mathbf{w}_{t-1}^{(i)}, \mathbf{x}_{t}\rangle + \|\mathbf{x}_{t}\|^{2} \\ &\leq \|\mathbf{w}_{t-1}^{(i)}\|^{2} + \|\mathbf{w}_{t-1}^{(1)} - \mathbf{w}_{t-1}^{(2)}\| \|\mathbf{x}_{t}\| + \|\mathbf{x}_{t}\|^{2} \\ &\stackrel{\text{constant}}{\end{aligned}}$$

Lower bound $\mathbb{E}_{\tilde{\mathcal{D}}}[\langle \mathbf{w_t^{(i)}}, \mathbf{w}^* \rangle]$:

$$\mathbb{E}[\langle \mathbf{w}_{\mathsf{t}}^{(\mathsf{i})}, \mathbf{w}^* \rangle] = \mathbb{E}[\langle \mathbf{w}_{\mathsf{t}-1}^{(\mathsf{i})}, \mathbf{w}^* \rangle] + \mathbb{E}[\tilde{y}_{\mathsf{t}} \langle \mathbf{x}_{\mathsf{t}}, \mathbf{w}^* \rangle]$$

$$= \mathbb{E}[\langle \mathbf{w}_{\mathsf{t}-1}^{(\mathsf{i})}, \mathbf{w}^* \rangle] + (1 - 2\mu) \underbrace{\mathbb{E}[\mathbf{y}_{\mathsf{t}} \langle \mathbf{x}_{\mathsf{t}}, \mathbf{w}^* \rangle]}_{\geq 1}$$

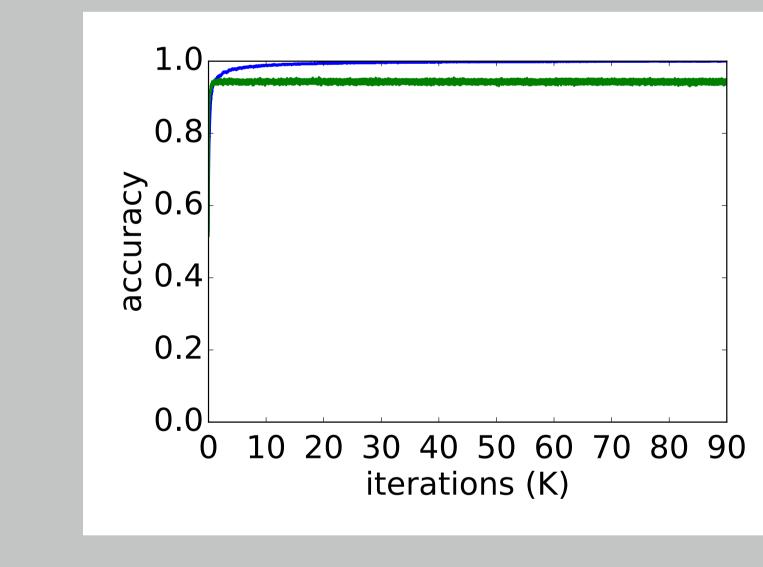
Combining these bounds:

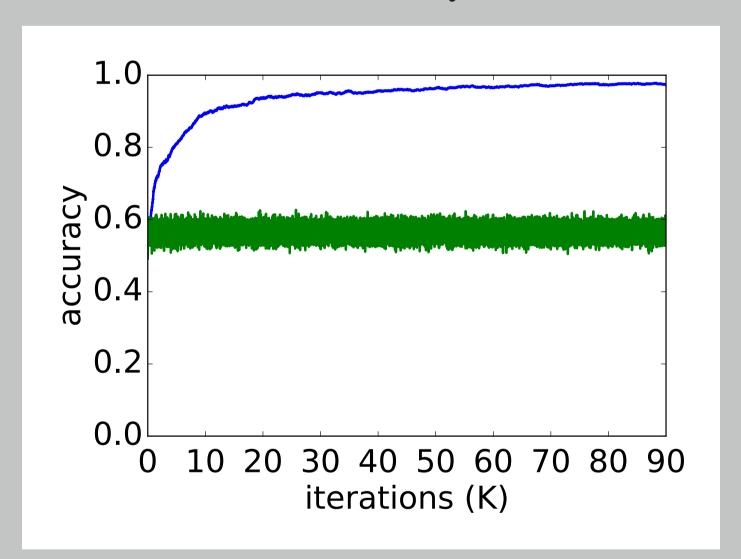
$$(1-2\mu)\,\mathbb{E}[\mathsf{T}] \tilde{\leq}\,\mathbb{E}[\langle\mathsf{w}_\mathsf{t}^\mathsf{(i)},\mathsf{w}^*
angle] \leq \mathbb{E}[\|\mathsf{w}_\mathsf{t}^\mathsf{(i)}\|\|\mathsf{w}^*\|] \ \tilde{\leq}\|\mathsf{w}^*\|\,\mathbb{E}[\sqrt{\mathsf{T}}] \leq \|\mathsf{w}^*\|\sqrt{\mathbb{E}[\mathsf{T}]}$$

Notation: we use \leq to hide constants

Experiments: Linear Classification

- Linear classification on Gaussian synthetic data with noise μ

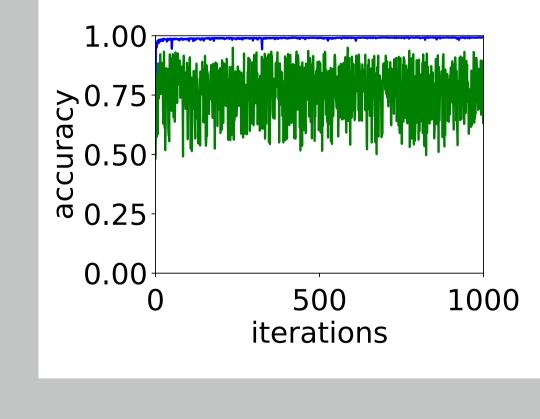


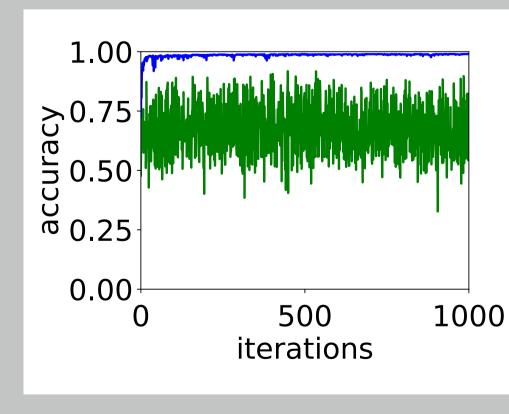


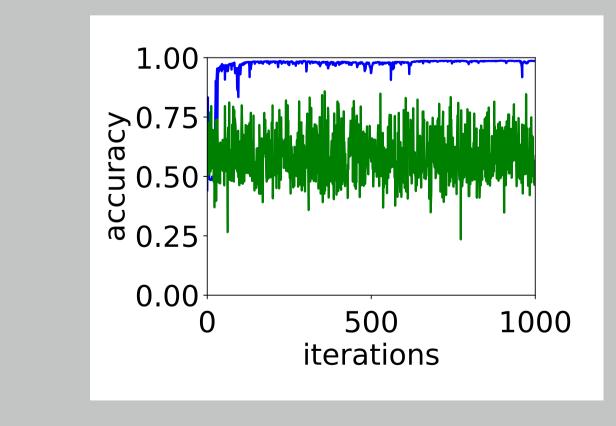
 $\mu = 0.01$

 $\mu = 0.4$ Figure 1: Our algorithm (blue) compared to a vanilla perceptron (green) on Gaussian data

- Linear classification on MNIST data with synthetic noise μ







 $\mu = 0.2$ $\mu = 0.3$

 $\mu = 0.4$ Figure 2: Mean accuracy of our algorithm vs. perceptron on MNIST classification with noise

Experiments: Deep Learning (Dataset)

- Testing on deep neural-networks trained with momentum SGD
- Construct a real-world noisy dataset for gender classification:
- ▶ Labeled Faces in the Wild benchmark celebrities images with names [1]
- ▶ To provide labeling, use an online service to match gender to name
- > Filter examples with small confidence to get a dataset with less noise

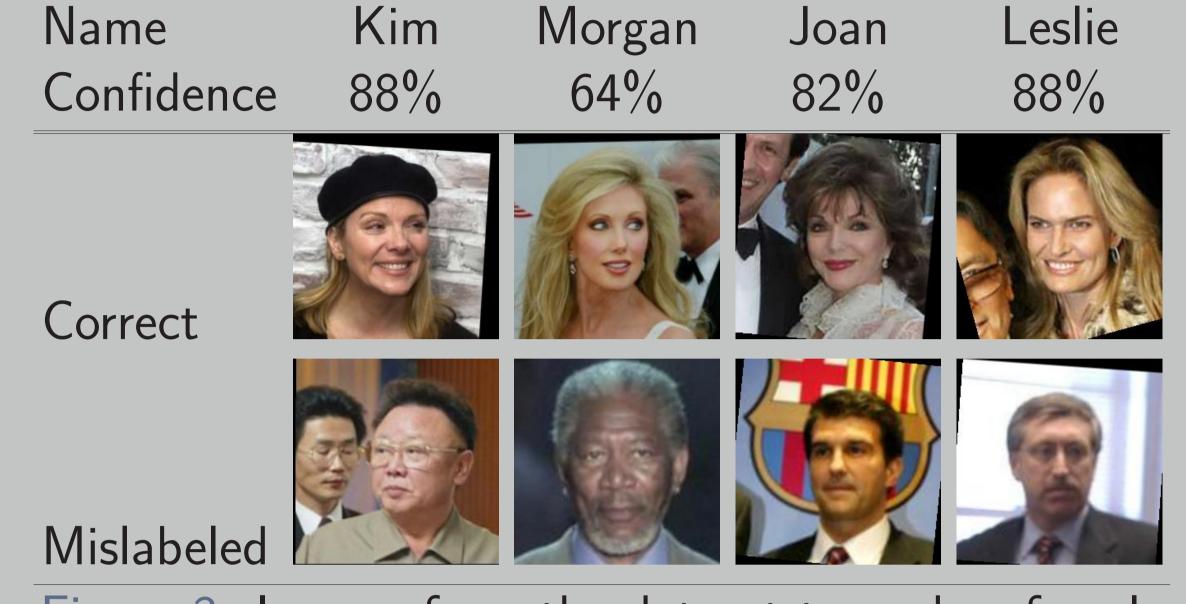
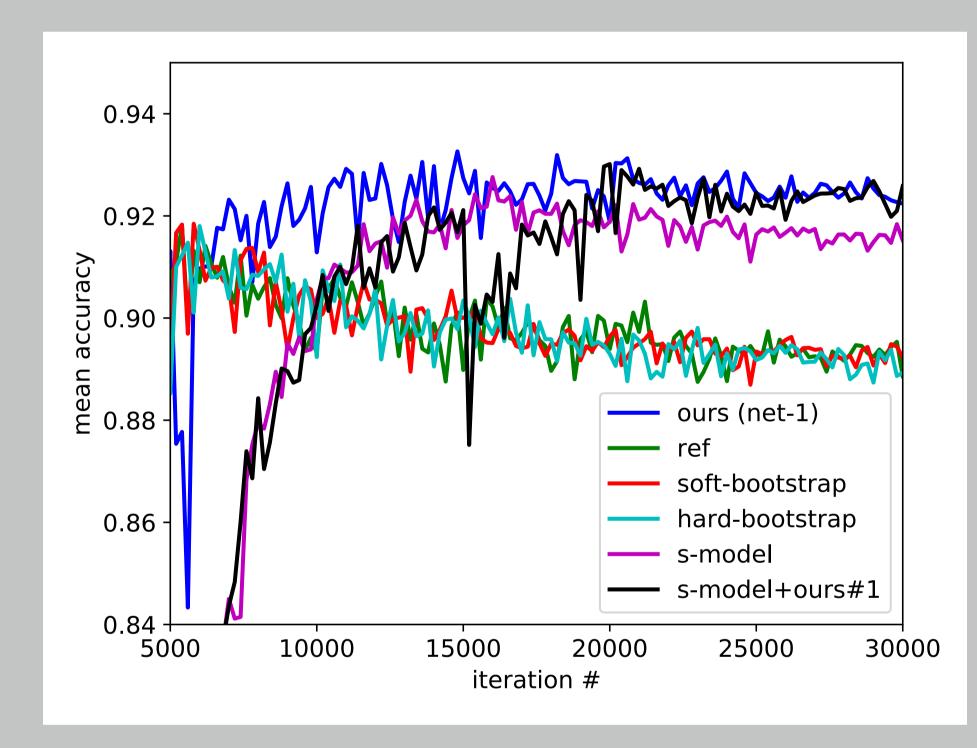
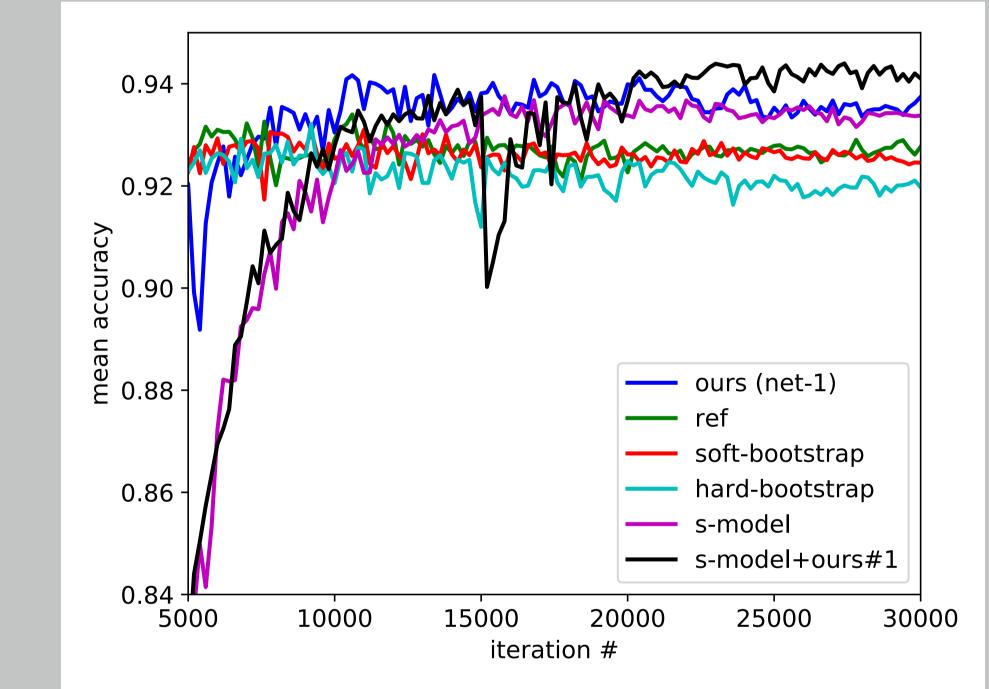


Figure 3: Images from the dataset tagged as female

Experiments: Deep Learning (Results)

- ► This method proves to be more effective than competing methods [2, 3]
- Combination with other methods can further improve results





Dataset #1 - more noise

Dataset #2 - less noise

Figure 4: Balanced accuracy of all methods on clean test data, trained on different datasets.

Conclusion

- A simple solution to supervised learning in the presence of noisy labels, which can be adapted to any update-based algorithm
- Theoretical analysis shows fast convergence for linear classification
- Empirical results show that the method improves the robustness of neural-networks to noisy labels in a real-world scenario

References

- 1] Gary B Huang, Manu Ramesh, Tamara Berg, and Erik Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical report, Technical Report 07-49, University of Massachusetts, Amherst, 2007.
- [2] Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. arXiv preprint arXiv:1412.6596, 2014.
- 3] Jacob Goldberger and Ehud Ben-Reuven. Training deep neural networks using a noise adaptation layer.