# Confident Learning applied to MNIST Label Error

David Szczecina

# Confident Learning and Cleanlab

Noisy Data, $X$
$(x, \tilde{y})^n \in (\mathbb{R}^d, \mathbb{Z}_{\geq 0})^n$

Model, $\boldsymbol{\theta}$

Noisy Predicted Probs, $\hat{p}(\tilde{y}; x, \boldsymbol{\theta})$

Noisy inputs

Confident Joint, $C_{\tilde{y},y^*}$
Estimate of Joint, $\widehat{Q}_{\tilde{y},y^*}$

Prune

cleanlab

Clean Data

Count

| $C_{\tilde{y},y^*}$ | $y^*=dog$ | $y^*=fox$ | $y^*=cow$ |
|---|---|---|---|
| $\tilde{y}=dog$ | 100 | 40 | 20 |
| $\tilde{y}=fox$ | 56 | 60 | 0 |
| $\tilde{y}=cow$ | 32 | 12 | 80 |

Normalize rows to match prior & divide by total

| $\widehat{Q}_{\tilde{y},y^*}$ | $y^*=dog$ | $y^*=fox$ | $y^*=cow$ |
|---|---|---|---|
| $\tilde{y}=dog$ | 0.25 | 0.1 | 0.05 |
| $\tilde{y}=fox$ | 0.14 | 0.15 | 0 |
| $\tilde{y}=cow$ | 0.08 | 0.03 | 0.2 |

Dirty Data
$\begin{pmatrix} \text{Examples with} \\ \text{Label Issues} \end{pmatrix}$

1. Estimate the joint distribution of given, noisy labels and latent (unknown) uncorrupted labels to fully characterize class-conditional label noise.
2. Find and prune noisy examples with label issues.
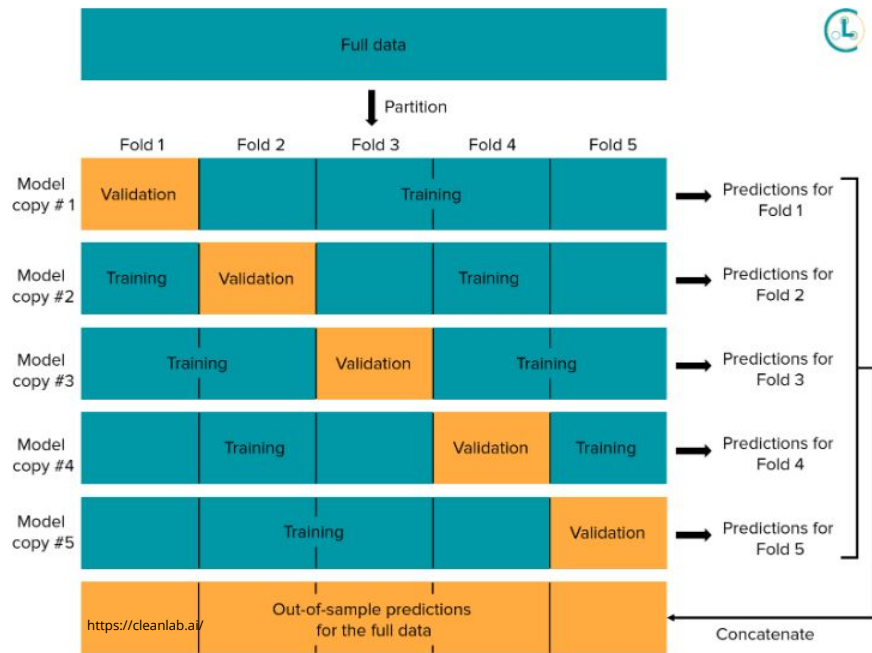3. Train with errors removed, re-weighting examples by the estimated latent prior.

https://arxiv.org/pdf/1911.00068.pdf

# Confident Learning and Cleanlab

Out of sample Predictions                    +                    Labels (can contain label errors)



Full data

↓ Partition

Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5

Model copy # 1 — Validation | Training → Predictions for Fold 1

Model copy #2 — Training | Validation | Training → Predictions for Fold 2

Model copy #3 — Training | Validation | Training → Predictions for Fold 3

Model copy #4 — Training | Validation | Training → Predictions for Fold 4

Model copy #5 — Training | Validation → Predictions for Fold 5

Out-of-sample predictions for the full data → Concatenate

https://cleanlab.ai/

0 1 2 3 4
5 6 7 8 9

# Predicted Probabilities

The central idea is that when the predicted probability of an example is greater than a per-class-threshold, we *confidently count* that example as actually belonging to that threshold's class.

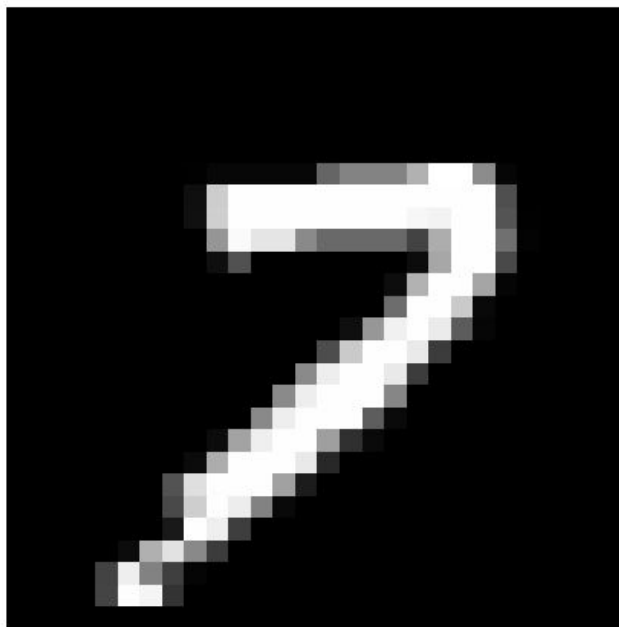The thresholds for each class are the average predicted probability of examples in that class.
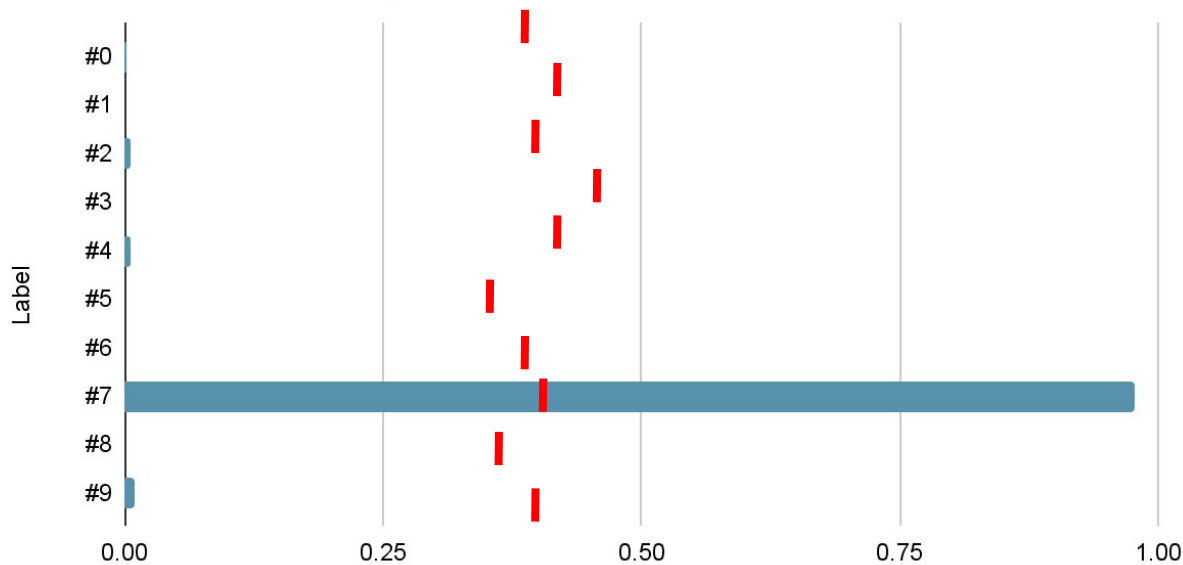
predicted probability vs. Number

# MNSIT Example

$$C_{\tilde{y}, y^*}[i][j] := |\hat{X}_{\tilde{y}=i, y^*=j}| \quad \text{where}$$

$$\hat{X}_{\tilde{y}=i, y^*=j} := \left\{ x \in X_{\tilde{y}=i} : \hat{p}(\tilde{y}=j; x, \theta) \geq t_j, \; j = \operatorname*{arg\,max}_{l \in [m]: \hat{p}(\tilde{y}=l; x, \theta) \geq t_l} \hat{p}(\tilde{y}=l; x, \theta) \right\}$$

$$t_j = \frac{1}{|X_{\tilde{y}=j}|} \sum_{x \in X_{\tilde{y}=j}} \hat{p}(\tilde{y}=j; x, \theta)$$

```
plot_examples([59915])
```

id: 59915
label: 4



Normalized Probability Score

Class Thresholds *example    Probability

# Results

148 out of 70000 labels identified as bad.

Many are mislabeled or questionable, some are correct and false positives.

# Ranking Calculation

'normalized_margin': normalized margin (p(label = k) - max(p(label != k)))

'self_confidence': [pred_probs[i][labels[i]] for i in label_issues_idx]

'confidence_weighted_entropy': entropy(pred_probs) / self_confidence

# Confusion Matrix

Predicted

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6891 | 0 | 1 | 0 | 1 | 0 | 7 | 1 | 1 | 1 |
| 1 | 0 | 7843 | 17 | 1 | 5 | 0 | 0 | 8 | 3 | 0 |
| 2 | 3 | 12 | 6925 | 3 | 8 | 0 | 2 | 27 | 9 | 1 |
| 3 | 0 | 7 | 33 | 7051 | 0 | 14 | 1 | 17 | 8 | 10 |
| 4 | 1 | 2 | 2 | 0 | 6794 | 0 | 8 | 1 | 3 | 13 |
| 5 | 1 | 2 | 2 | 9 | 1 | 6266 | 14 | 1 | 12 | 5 |
| 6 | 7 | 8 | 0 | 0 | 5 | 4 | 6847 | 0 | 5 | 0 |
| 7 | 4 | 5 | 23 | 3 | 10 | 0 | 0 | 7221 | 5 | 22 |
| 8 | 3 | 18 | 3 | 4 | 9 | 8 | 11 | 3 | 6745 | 21 |
| 9 | 10 | 4 | 0 | 7 | 17 | 7 | 0 | 24 | 9 | 6880 |

Given Label

Total Errors: 148

# 3 Folds vs 10 Folds

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6891 | 0 | 1 | 0 | 1 | 0 | 7 | 1 | 1 | 1 |
| 1 | 0 | 7843 | 17 | 1 | 5 | 0 | 0 | 8 | 3 | 0 |
| 2 | 3 | 12 | 6925 | 3 | 8 | 0 | 2 | 27 | 9 | 1 |
| 3 | 0 | 7 | 33 | 7051 | 0 | 14 | 1 | 17 | 8 | 10 |
| 4 | 1 | 2 | 2 | 0 | 6794 | 0 | 8 | 1 | 3 | 13 |
| 5 | 1 | 2 | 2 | 9 | 1 | 6266 | 14 | 1 | 12 | 5 |
| 6 | 7 | 8 | 0 | 0 | 5 | 4 | 6847 | 0 | 5 | 0 |
| 7 | 4 | 5 | 23 | 3 | 10 | 0 | 0 | 7221 | 5 | 22 |
| 8 | 3 | 18 | 3 | 4 | 9 | 8 | 11 | 3 | 6745 | 21 |
| 9 | 10 | 4 | 0 | 7 | 17 | 7 | 0 | 24 | 9 | 6880 |

10 fold cross validation: 98 errors

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6903 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 7874 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 3 | 6979 | 0 | 4 | 0 | 0 | 2 | 1 | 1 |
| 3 | 0 | 0 | 4 | 7122 | 0 | 2 | 0 | 5 | 2 | 6 |
| 4 | 0 | 0 | 1 | 0 | 6816 | 0 | 0 | 1 | 1 | 5 |
| 5 | 1 | 0 | 1 | 4 | 0 | 6305 | 2 | 0 | 0 | 0 |
| 6 | 1 | 0 | 0 | 0 | 2 | 4 | 6869 | 0 | 0 | 0 |
| 7 | 0 | 2 | 7 | 0 | 3 | 0 | 0 | 7276 | 1 | 4 |
| 8 | 2 | 0 | 4 | 0 | 1 | 1 | 0 | 0 | 6815 | 2 |
| 9 | 3 | 0 | 0 | 0 | 3 | 4 | 0 | 4 | 1 | 6943 |

3 fold cross validation: 148 errors

Highlighted in red if more than 10 errors.

Less false positives with more cross validation folds

# Next Steps

- See difference in model accuracy without bad labeled data
- See how well it works with additional augmented data
- Test on BIOSCAN-small dataset