

Summary of Read Papers

Sixu Li

- ”Using Trusted Data to Train Deep Networks on Labels Corrupted by Severe Noise” [1]

Preliminaries:

We are given an untrusted dataset $\tilde{\mathcal{D}}$ of u examples (x, \tilde{y}) , and we assume that these examples are *potentially* corrupted examples from the true data distribution $p(x, y)$ with K classes. Corruption is specified by a label noise distribution $p(\tilde{y}|y, x)$. We are also given a trusted dataset \mathcal{D} of t examples drawn from $p(x, y)$, where $t/(t + u) \ll 1$. Our method makes use of \mathcal{D} to estimate the $K \times K$ matrix of corruption probabilities $C_{ij} = p(\tilde{y} = j | y = i)$.

Key Results:

By Bayes’ theorem, we have,

$$p(\tilde{y} | y, x)p(x | y) = p(\tilde{y} | y)p(x | \tilde{y}, y) \quad (1)$$

Intergrating over all x on both sides gives us,

$$\int p(\tilde{y} | y, x)p(x | y)dx = p(\tilde{y} | y) \int p(x | \tilde{y}, y)dx = p(\tilde{y} | y) \quad (2)$$

We can approximate the intergral on the left with the expectation of $p(\tilde{y} | y, x)$ over the empirical distribution of x given y . Assuming conditional independence of \tilde{y} and y given x , we have $p(\tilde{y} | y, x) = p(\tilde{y} | x)$, which can be directly approximated by $\hat{p}(\tilde{y} | x)$, the classifier trained on $\tilde{\mathcal{D}}$. More explicitly, let A_i be the subset of x in \mathcal{D} with label i . Denote our estimate of C by \hat{C} . We have,

$$\hat{C}_{ij} = \frac{1}{|A_i|} \sum_{x \in A_i} \hat{p}(\tilde{y} = j | x) = \frac{1}{|A_i|} \sum_{x \in A_i} \hat{p}(\tilde{y} = j | y = i, x) \approx p(\tilde{y} = j | y = i) \quad (3)$$

The accuracy of this approximation relies on three effects:

- (i). $\hat{p}(\tilde{y} | x)$ being a good estimate of $p(\tilde{y} | x)$.
- (ii). The number of trusted examples of each class, which effects the second approximation of Eq. (3).

(iii). Conditional independence assumption is satisfied.

- ”Combating Label Noise in Deep Learning Using Abstention” [2]

Preliminaries:

We assume we are interested in training a k -class multi-class classifier with a deep neural network(DNN) where x is the input and y is the output. For a given x , we define $p_i = p_w(y = i|x)$ as the i^{th} output of the DNN that implements the probability model $p_w(y = i|x)$ where w is the parameters of the DNN.

Key Results:

We train the deep abstaining classifier(DAC) with following modified version of the k -class cross-entropy per-sample loss:

$$\mathcal{L}(x_j) = (1 - p_{k+1}) \left(- \sum_{i=1}^k q_i \log \frac{p_i}{1 - p_{k+1}} \right) + \alpha \log \frac{1}{1 - p_{k+1}} \quad (4)$$

References

- [1] HENDRYCKS, D., MAZEIKA, M., WILSON, D., AND GIMPEL, K. Using trusted data to train deep networks on labels corrupted by severe noise. In *Advances in neural information processing systems* (2018), pp. 10456–10465.
- [2] THULASIDASAN, S., BHATTACHARYA, T., BILMES, J., CHENNUPATI, G., AND MOHD-YUSOF, J. Combating label noise in deep learning using abstention. *arXiv preprint arXiv:1905.10964* (2019).