

Robust Losses

Sixu Li

1 Preliminary

Given a C-class dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, with $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$ denoting a sample and $y_i \in \mathcal{Y} = \{1, \dots, C\}$ its associated label. For each sample \mathbf{x}_i , a classifier $f(\mathbf{x}_i)$ (DNNs here) computes its probability of each label $j \in 1, \dots, C$: $p(j|\mathbf{x}_i) = \frac{\exp(\mathbf{z}_{ij})}{\sum_{k=1}^C \exp(\mathbf{z}_{ik})}$, where \mathbf{z}_i are the logits on sample \mathbf{x}_i . We denote the ground-truth distribution over labels for sample \mathbf{x}_i by $q(j|\mathbf{x}_i)$, and $\sum_{j=1}^C q(j|\mathbf{x}_i) = 1$. Consider the case of a single ground-truth label y_i , then $q(y_i|\mathbf{x}_i) = 1$ and $q(j|\mathbf{x}_i) = 0$ for all $j \neq y_i$. Then categorical cross entropy (CCE) loss and mean absolute error (MAE) loss for sample \mathbf{x}_i can be written as:

$$\mathcal{L}_{\text{CCE}}(\mathbf{x}_i) = - \sum_{j=1}^C q(j|\mathbf{x}_i) \log p(j|\mathbf{x}_i) = - \log p(y_i|\mathbf{x}_i) \quad (1)$$

$$\mathcal{L}_{\text{MAE}}(\mathbf{x}_i) = \sum_{j=1}^C |p(j|\mathbf{x}_i) - q(j|\mathbf{x}_i)| = 2(1 - p(y_i|\mathbf{x}_i)) \quad (2)$$

2 Summary of Papers

- "Robust Loss Functions under Label Noise for Deep Neural Networks" [1]

Key Results:

(i). Define "Symmetric Losses":

We call a loss function L symmetric if it satisfies, for some constant C ,

$$\sum_{i=1}^k \mathcal{L}(f(x), i) = C, \quad \forall x \in \mathcal{X}, \forall f \quad (3)$$

For these loss functions, we have,

$$\sum_{i=1}^k \mathcal{L}(f(x), e_i) = \begin{cases} \sum_{i=1}^k \log \frac{1}{f_i(x)}, & \text{CCE} \\ \sum_{i=1}^k (2 - 2f_i(x)) = 2k - 2, & \text{MAE} \\ k \|f_i(x)\|_2^2 + k - 2, & \text{MSE} \end{cases} \quad (4)$$

Thus, among these, only MAE satisfies symmetry condition given by Eq.(3).

(ii). Theoretically proved that Symmetric Losses are robust to label noise with a little additional assumptions and verified it by empirical experiments

(iii). Use empirical results showing that Bounded Losses are more robust to unbounded losses.

• ”Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels” [5]

Key Results:

(i). Give an explanation of robustness of MAE and noise-sensitiveness of CCE.

Notations: Let $\mathcal{X} \subset \mathbb{R}^d$ be the feature space and $\mathcal{Y} = 1, \dots, c$ be the label space, $(\mathbf{x}_i, y_i) \in (\mathcal{X}, \mathcal{Y})$. Denote classifier $f(\mathbf{x}; \boldsymbol{\theta}) : \mathcal{X} \rightarrow \mathcal{R}^c$ and f_j denotes the j 'th element of f .

Let's look at the gradient of the loss functions:

$$\sum_{i=1}^n \frac{\partial \mathcal{L}(f(\mathbf{x}_i; \boldsymbol{\theta}), y_i)}{\partial \boldsymbol{\theta}} = \begin{cases} \sum_{i=1}^n -\frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}), & \text{for CCE} \\ \sum_{i=1}^n -\nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}), & \text{for MAE} \end{cases} \quad (5)$$

In this paper, they claims that, in CCE, $\frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})}$ can be viewed as the weight put on each example when updating the gradient. Therefore, this means that, when training with CCE, more emphasis is put on difficult examples. This implicit weighting scheme is desirable for training with clean data, but can cause overfitting to noisy labels. Conversely, MAE treats every sample equally, which makes it more robust to noisy labels but meanwhile slower convergence rate. Moreover, without the implicit weighting scheme to focus on challenging samples, the stochasticity involved in the training process can make learning difficult. As a result, classification accuracy might suffer.

My Question: Is the perspective that viewing $\frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})}$ as the weight put on each example when updating the gradient correct?

In my opinion, large $\frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})}$ does not guarantee that sample \mathbf{x}_i has large contribution to the gradient updating because $\nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})$ also contains information that (\mathbf{x}_i, y_i) provides. It is possible that for large $\frac{1}{f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})}$, it corresponds to small $\nabla_{\boldsymbol{\theta}} f_{y_i}(\mathbf{x}_i; \boldsymbol{\theta})$. Similarly, the explanation for MAE also should be doubted.

(ii). They proposed a generalized version of MAE and CCE:

$$\mathcal{L}_q(f(\mathbf{x}, \mathbf{e}_j)) = \frac{(1 - f_j(\mathbf{x})^q)}{q} \quad (6)$$

Using L'Hopital's rule, it can be shown that the proposed loss function is equivalent to CCE for $\lim_{q \rightarrow 0} \mathcal{L}_q(f(\mathbf{x}, \mathbf{e}_j))$, and becomes MAE when $q = 1$. This means that \mathcal{L}_q utilizes

hyperparameter q to control the trade-off between fast convergence of CCE and noise-robustness of MAE.

Besides, in order to let \mathcal{L}_q to have a tight bound, they proposed Truncated \mathcal{L}_q loss:

$$\mathcal{L}_{\text{trunc}}(f(\mathbf{x}, \mathbf{e}_j)) = \begin{cases} \mathcal{L}_q(k), & \text{if } f_j(\mathbf{x}) \leq k \\ \mathcal{L}_q(f(\mathbf{x}, \mathbf{e}_j)), & \text{if } f_j(\mathbf{x}) > k \end{cases} \quad (7)$$

In this way, if the softmax output for the provided label is below a threshold, truncated \mathcal{L}_q loss becomes a constant. Thus, the loss gradient is zero for that sample, and it does not contribute to learning dynamics.

Drawback: this process is equivalent to drop out those "hard" samples, which, in some senses, filter out those noisy samples. However, it will also ignore those challenging but clean data. Therefore, how to distinguish "hard" but clean samples and noisy samples is an important problem to think about.

- "IMAE for Noise-Robust Learning: Mean Absolute Error Does Not Treat Examples Equally and Gradient Magnitude's Variance Matters"[2]

Preliminaries:

Training set $\mathbf{X} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, training sample $\mathbf{x}_i \in \mathbb{R}^D$ and its annotated class $y_i \in \{1, 2, \dots, C\}$. Let f_θ be a deep neural network, which transforms \mathbf{x}_i to a representation $\mathbf{f}_i = f_\theta(\mathbf{x}_i) \in \mathbb{R}^E$. The linear classifier follows the output embeddings f_θ and is composed of one fully connected(FC) layer, one softmax layer and one loss layer. The FC layer can be represented as $\mathbf{z}_i = \mathbf{W}^T \mathbf{f}_i \in \mathbb{R}^C$, where $\mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_C] \in \mathbb{R}^{E \times C}$. $\mathbf{z}_{ij} = \mathbf{w}_j^T \mathbf{f}_i$ is a logit that indicates the compatibility between sample \mathbf{x}_i and class j . A softmax function is given by,

$$p(j|\mathbf{x}_i) = \frac{\exp(\mathbf{z}_{ij})}{\sum_{m=1}^C \exp(\mathbf{z}_{im})} \quad (8)$$

where $p(j|\mathbf{x}_i)$ is the probability of sample \mathbf{x}_i being predicted by class j .

Key Results:

(i). They define gradient magnitude w.r.t logit vector \mathbf{z} and the gradients' variance(See derivation in the original paper)

$$\omega_{\text{CCE}}(\mathbf{x}_i) = \left\| \frac{\partial L_{\text{CCE}}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \right\|_1 = 2(1 - p(y_i|\mathbf{x}_i)) \quad (9)$$

$$\omega_{\text{MAE}}(\mathbf{x}_i) = \left\| \frac{\partial L_{\text{MAE}}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \right\|_1 = 4p(y_i|\mathbf{x}_i)(1 - p(y_i|\mathbf{x}_i)) \quad (10)$$

Assuming that samples' probabilities $p(y_i|\mathbf{x}_i)$ are uniformly distributed, then the gradi-

ents' variance of MAE over training data points can be written as,

$$\sigma_{\text{MAE}} = \int_0^1 \omega_{\text{MAE}}^2(p) dp - \left(\int_0^1 \omega_{\text{MAE}}(p) dp \right)^2 \quad (11)$$

They treat gradient magnitude as weight on each sample, i.e, if one sample's gradient is larger, its impact is larger during gradient back-propagation.

Based on this consideration, they provide two explanations of MAE's properties:

1. MAE emphasizes on uncertain examples($p(y_i|\mathbf{x}_i)$ is closer to 0.5) and those samples with high-probability and low-probability have smaller weight. While low-probability ones are highly likely to be noisy as a model improves during training, therefore, MAE is noise-robustness.

2. MAE's gradient magnitude's variance over data points is only 0.09(too small). As a consequence, the impact ratio of one example versus another is too small. Therefore, the majority contribute almost equally. Therefore, MAE generally underfits to training data.

My Question: Is the gradient magnitudes truly reflect the samples' impact on the gradient updating? Is the L_1 norm of gradients meaningful?

My suspicion on this problem: For \forall loss function \mathcal{L} , following the previous notation, when training on data point (\mathbf{x}_i, y_i) , now we consider the gradients on parameters W . By chain rule, we have,

$$\begin{aligned} \frac{\partial \mathcal{L}(\mathbf{x}_i)}{\partial \mathbf{W}} &= \frac{\partial \mathcal{L}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \frac{\partial \mathbf{z}_i}{\partial \mathbf{W}} \\ &= \frac{\partial \mathcal{L}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \mathbf{f}_i \in \mathbb{R}^{E \times C} \end{aligned} \quad (12)$$

Following original paper's definition, if we define L_1 norm of gradients to be gradients' magnitude, which measures the impact of each training sample on gradient updating. Then, let's consider two data point (\mathbf{x}_1, y_1) and (\mathbf{x}_2, y_2) . If we have $\|\frac{\partial \mathcal{L}(\mathbf{x}_1)}{\partial \mathbf{z}_1}\|_1 > \|\frac{\partial \mathcal{L}(\mathbf{x}_2)}{\partial \mathbf{z}_2}\|_1$, under this paper's conclusion, the impact of sample \mathbf{x}_1 is larger than \mathbf{x}_2 during gradient back-propagation. However, it doesn't guarantee that $\|\frac{\partial \mathcal{L}(\mathbf{x}_1)}{\partial \mathbf{W}}\|_1 > \|\frac{\partial \mathcal{L}(\mathbf{x}_2)}{\partial \mathbf{W}}\|_1$.

So, one important problem to consider, how can measure training samples' impact on gradient updating?

(ii). In order to solve the problem that MAE's gradient magnitude's variance over data points is too small. They define IMAE that transforms MAE's weighting scheme non-linearly:

$$\omega_{\text{IMAE}}(\mathbf{x}_i) = \exp(Tp(y_i|\mathbf{x}_i))(1 - p(y_i|\mathbf{x}_i)) \quad (13)$$

And in back-propagation, they simply scale the gradient w.r.t logits as follows:

$$\frac{\partial \mathcal{L}_{\text{IMAE}}(\mathbf{x}_i)}{\partial \mathbf{z}_i} = \frac{\partial \mathcal{L}_{\text{MAE}}(\mathbf{x}_i)}{\partial \mathbf{z}_i} \frac{\omega_{\text{IMAE}}(\mathbf{x}_i)}{\omega_{\text{MAE}}(\mathbf{x}_i)} \quad (14)$$

(Figure that compares the gradients' variance is shown in the original paper.)

Drawback: By looking at the Figure 4 in the original paper, we can find out that the convergence rate of IMAE is the same as MAE, which is still a problem.

My Question: When training on clean data point (\mathbf{x}_i, y_i) , let's look at the following two gradient on the logit \mathbf{z}_{ij} :

$$\frac{\partial \mathcal{L}_{\text{CCE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} = \begin{cases} p(y_i|\mathbf{x}_i) - 1, & j = y_i \\ p(j|\mathbf{x}_i), & j \neq y_i \end{cases} \quad (15)$$

$$\frac{\partial \mathcal{L}_{\text{MAE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} = \begin{cases} 2p(y_i|\mathbf{x}_i)(p(y_i|\mathbf{x}_i) - 1), & j = y_i \\ 2p(y_i|\mathbf{x}_i)p(j|\mathbf{x}_i), & j \neq y_i \end{cases} \quad (16)$$

We can observe that when $j = y_i$, both $\frac{\partial \mathcal{L}_{\text{CCE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}}$ and $\frac{\partial \mathcal{L}_{\text{MAE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}}$ are less than 0 and when $j \neq y_i$, they are all greater than zero. Therefore, is the sign of gradients meaningful?

The expressions of these two gradients when training on noise label $(\mathbf{x}_i, \hat{y}_i)$ can be extended easily, but need to take into consideration.

- **"Symmetric Cross Entropy for Robust Learning with Noisy Labels" [3]**

Key Results:

(i). They did a class-based experiment, which tracks the test accuracy of each class along the training procedure(see the experimental results in the original paper). In this experiment, they found out that DNN learning with CE can be **class-biased**: some classes("easy" classes) are easy to learn and converge faster than other classes("hard" classes). And this phenomenon is amplified when training labels are noisy: while easy classes already overfit to noisy labels, hard classes still suffer from significant **under learning**(class accuracy significantly lower than clean label setting). They claim that it appears that the under learning of hard classes is a major cause for the overall performance degradation.(The experiment phenomenons to prove their claim are shown in the original paper).

My Question: They only do this experiment on data set CIFAR-10 with 40% symmetric noise, hence, is this conclusion true on other data set with different kinds of noise?

(ii). Inspired by the symmetric KL-divergence, they proposed a Symmetric Cross Entropy(SCE) as,

$$SCE = CE + RCE = H(q, p) + H(p, q) \quad (17)$$

And we can write loss function Reverse Cross Entropy(RCE) as,

$$\mathcal{L}_{\text{RCE}}(\mathbf{x}_i) = - \sum_{j=1}^C p(j|\mathbf{x}_i) \log q(j|\mathbf{x}_i) \quad (18)$$

And define,

$$\mathcal{L}_{\text{SL}} = \alpha \mathcal{L}_{\text{CE}} + \beta \mathcal{L}_{\text{RCE}} \quad (19)$$

Because $q(j|\mathbf{x}_i)$ is an one-hot vector, in order to let RCE valid, they define $\log 0$ to be a negative number A . And similar to [1], they prove that RCE is noise-robustness under a little assumption.

And let's look at the gradients of SL on the logits \mathbf{z}_{ij} when training on clean data point (\mathbf{x}_i, y_i) .

$$\frac{\partial \mathcal{L}_{\text{SL}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} = \begin{cases} \frac{\partial \mathcal{L}_{\text{CE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} + Ap(y_i|\mathbf{x}_i)(1 - p(y_i|\mathbf{x}_i)), & j = y_i \\ \frac{\partial \mathcal{L}_{\text{CE}}(\mathbf{x}_i)}{\partial \mathbf{z}_{ij}} - Ap(y_i|\mathbf{x}_i)p(j|\mathbf{x}_i), & j \neq y_i \end{cases} \quad (20)$$

We can take a look at the experimental results of SL in the original paper.

My Question:

- (1). Look at their explanation of this gradient, I think it does not make sense.
- (2). All of these paper only explain the behavior of these gradients when training on clean data, however, they do not consider the form of these gradients when training on noisy label.(Write out those gradients when training on noisy data and take a look). And this question relates to the meaning of gradients' sign on each component.

- " \mathcal{L}_{DMI} : A Novel Information-theoretic Loss Function for Training Deep Nets Robust to Label Noise" [4]

Preliminaries:

Denote the set of classes by \mathcal{C} and the size of \mathcal{C} is C . Also, denote the domain of datapoints by \mathcal{X} . A classifier is denoted by $h : \mathcal{X} \rightarrow \Delta_{\mathcal{C}}$, where $\Delta_{\mathcal{C}}$ is the set of all possible distributions over \mathcal{C} . h represents a randomized classifier such that given $x \in \mathcal{X}$, $h(x)_c$ is the probability that h maps x into class c . Note that fixing the input x , the randomness of a classifier is independent of everything else.

There are N datapoints $\{x_i\}_{i=1}^N$. For each datapoint x_i , there is an unknown ground truth $y_i \in \mathcal{C}$. We assume that there is an unknown prior distribution $Q_{X,Y}$ over $\mathcal{X} \times \mathcal{C}$ such that $\{(x_i, y_i)\}_{i=1}^N$ are i.i.d. samples drawn from $Q_{X,Y}$ and

$$Q_{X,Y}(x, y) = Pr[X = x, Y = y] \quad (21)$$

And we Denote matrix format of joint distribution over X and Y as,

$$\mathbf{Q}_{X,Y} = [Pr[X = x_i, Y = y_i]] \quad (22)$$

Besides, define noise transition matrix as,

$$\mathbf{T}_{Y \rightarrow \tilde{Y}} = [Pr[\tilde{Y} = \tilde{y}|Y = y]]_{C \times C} \quad (23)$$

Key Results:

(1) Problem of Shannon mutual information and new proposed determinant based mutual information.

For every two random variables W_1 and W_2 , Shannon mutual information defined as

$$\mathbf{MI}(W_1, W_2) \triangleq \sum_{w_1, w_2} Pr[W_1 = w_1, W_2 = w_2] \log \frac{Pr[W_1 = w_1, W_2 = w_2]}{Pr[W_1 = w_1]Pr[W_2 = w_2]} \quad (24)$$

Properties of $\mathbf{MI}(W_1, W_2)$:

- (i). Non-negative.
- (ii). Symmetric, i.e., $\mathbf{MI}(W_1, W_2) = \mathbf{MI}(W_2, W_1)$.
- (iii). Information-monotonicity:

For all random variables W_1, W_2, W_3 , when W_3 is less informative for W_2 than W_1 , i.e., W_3 is independent of W_2 conditioning on W_1 , then,

$$\mathbf{MI}(W_3, W_2) \leq \mathbf{MI}(W_1, W_2) \quad (25)$$

Based on Shannon mutual information, a performance measure for a classifier h can be naturally defined. High quality classifier's output $h(X)$ should have high mutual information with the ground truth category Y . Thus, a classifier h 's performance can be measured by $\mathbf{MI}(h(X), Y)$.

Problem for $\mathbf{MI}(W_1, W_2)$ when training on noisy labels:

$$\forall h, h', \mathbf{MI}(h(X), Y) > \mathbf{MI}(h'(X), Y) \not\Leftrightarrow \mathbf{MI}(h(X), \tilde{Y}) > \mathbf{MI}(h'(X), \tilde{Y}) \quad (26)$$

In order to solve this problem, we proposed a new measure called Determinant based Mutual Information(**DMI**).

Definition:

Given two discrete random variables W_1, W_2 , we define the Determinant based Mutual Information between W_1 and W_2 as

$$\mathbf{DMI}(W_1, W_2) = |\det(\mathbf{Q}_{W_1, W_2})| \quad (27)$$

where \mathbf{Q}_{W_1, W_2} is the matrix format of the joint distribution over W_1 and W_2 .

Properties of **DMI**:

- (i). Non-negative.
- (ii). Symmetric.
- (iii). Information-monotonicity.
- (iv). Relative Invariance (See the Proof in original paper):

For random variables W_1, W_2, W_3 , when W_3 is less informative for W_2 than W_1 , i.e., W_3 is independent of W_2 conditioning W_1 ,

$$\mathbf{DMI}(W_2, W_3) = \mathbf{DMI}(W_2, W_1) |\det(\mathbf{T}_{W_1 \rightarrow W_3})| \quad (28)$$

where $\mathbf{T}_{W_1 \rightarrow W_3}$ is the matrix format of

$$T_{W_1 \rightarrow W_3}(w_1, w_3) = \Pr[W_3 = w_3 | W_1 = w_1] \quad (29)$$

Based on property (iv), we have,

$$\forall h, h', \mathbf{DMI}(h(X), Y) > \mathbf{DMI}(h'(X), Y) \Leftrightarrow \mathbf{DMI}(h(X), \tilde{Y}) > \mathbf{DMI}(h'(X), \tilde{Y}) \quad (30)$$

(2) Proposed new loss function $\mathcal{L}_{\mathbf{DMI}}$:

Definition:

$$\mathcal{L}_{\mathbf{DMI}}(Q_{h(X), \tilde{Y}}) \triangleq -\log(\mathbf{DMI}(h(X), \tilde{Y})) = -\log(|\det(\mathbf{Q}_{h(X), \tilde{Y}})|) \quad (31)$$

(See how to calculate $\mathcal{L}_{\mathbf{DMI}}$ and Main Theorem in original paper. Look through the proof.)

My Question:

In the proof of Main Theorem, I think,

$$\mathbf{DMI}(h(X), \tilde{Y}) = \mathbf{DMI}(h(X), Y) |\det(\mathbf{T}_{Y \rightarrow \tilde{Y}})| \quad (32)$$

doesn't hold. Because, when using property (iv) to derive this result, the assumption $h(X) \perp \tilde{Y} | Y$ may not hold. When training on classifier h on \tilde{Y} , Y should be less informative for $h(X)$ than \tilde{Y} .

References

- [1] GHOSH, A., KUMAR, H., AND SASTRY, P. Robust loss functions under label noise for deep neural networks, 2017.

- [2] WANG, X., KODIROV, E., HUA, Y., AND ROBERTSON, N. M. Imae for noise-robust learning: Mean absolute error does not treat examples equally and gradient magnitude’s variance matters.
- [3] WANG, Y., MA, X., CHEN, Z., LUO, Y., YI, J., AND BAILEY, J. Symmetric cross entropy for robust learning with noisy labels. In *Proceedings of the IEEE International Conference on Computer Vision* (2019), pp. 322–330.
- [4] XU, Y., CAO, P., KONG, Y., AND WANG, Y. L_dmi: A novel information-theoretic loss function for training deep nets robust to label noise. In *Advances in Neural Information Processing Systems* (2019), pp. 6222–6233.
- [5] ZHANG, Z., AND SABUNCU, M. Generalized cross entropy loss for training deep neural networks with noisy labels, 2018.