

Some Papers about Reweighting

Narsil Zhang

Oct. 2019

Outline

Foreword

Focal Loss

GHM

Class-balanced loss

Robust Learning via Reweight

Meta-Weight-Net

MentorNet



Reweight is a method, not our goal.
Goal: class imbalance, noisy data(label)

Problems

- ▶ imbalance of number of data across each class
- ▶ the gradient information is dominant by numerous "easy" examples

Focal Loss [4]

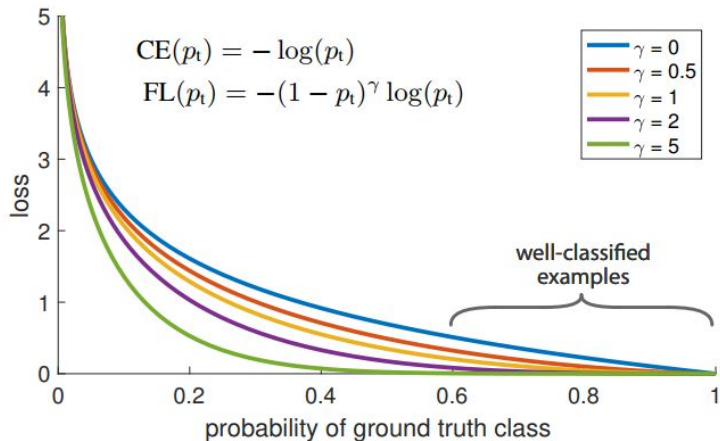


Figure: Focal Loss, t denotes class

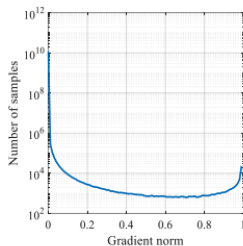
Gradient Harmonized Mechanism [3]

$$L_{GHM} = \sum_{i=1}^N \frac{L_{CE}(p_i)}{GD(g_i)}$$

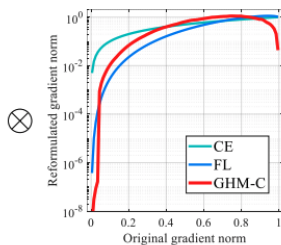
where g_i is the gradient norm relevant with i -th data and GD denotes some "gradient density" estimation

Effects

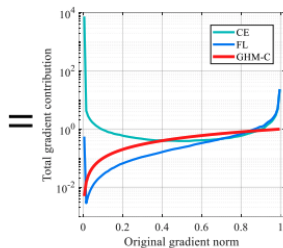
Gradient statistics



Gradient regularization



Overall gradient



GHM can suppress those extremely hard examples.

Effective number of samples

Denote the effective number of samples E_n is the expected volume of n samples.

Proposition

$$E_n = \frac{1 - \beta^n}{1 - \beta}, \quad \beta = \frac{n-1}{n}$$

Effective number of samples

Denote the effective number of samples E_n is the expected volume of n samples.

Proposition

$$E_n = \frac{1 - \beta^n}{1 - \beta}, \quad \beta = \frac{n-1}{n}$$

Proof.

Induction. Denote $p = \frac{E_{n-1}}{N}$, then

$$E_n = pE_{n-1} + (1-p)(E_{n-1} + 1) = 1 + \frac{N-1}{N}E_{n-1}$$



Class-balanced loss[1]

Suppose class y has n_y training samples, the class-balanced (CB) softmax cross-entropy loss is:

$$\text{CB}_{\text{softmax}}(\mathbf{z}, y) = -\frac{1 - \beta}{1 - \beta n_y} \log \left(\frac{\exp(z_y)}{\sum_{j=1}^C \exp(z_j)} \right)$$

where class y has n_y samples; In practice β is hyperparameter.

Class-balanced loss[1]

Suppose class y has n_y training samples, the class-balanced (CB) softmax cross-entropy loss is:

$$\text{CB}_{\text{softmax}}(\mathbf{z}, y) = -\frac{1 - \beta}{1 - \beta n_y} \log \left(\frac{\exp(z_y)}{\sum_{j=1}^C \exp(z_j)} \right)$$

where class y has n_y samples; In practice β is hyperparameter.

Class-balanced loss can be combined with focal loss.

Bilevel Optimization

$$\theta^*(w) = \arg \min_{\theta} \sum_{i=1}^N w_i f_i(\theta)$$

$$w^* = \arg \min_{w, w \geq 0} \frac{1}{M} \sum_{i=1}^M f_i^v(\theta^*(w))$$

Influence analysis

$$f_{i,\epsilon}(\theta) = \epsilon_i f_i(\theta)$$

$$\hat{\theta}_{t+1}(\epsilon) = \theta_t - \alpha \nabla \sum_{i=1}^n f_{i,\epsilon}(\theta) \Big|_{\theta=\theta_t}$$

$$\epsilon_t^* = \arg \min_{\epsilon} \frac{1}{M} \sum_{i=1}^M f_i^v(\theta_{t+1}(\epsilon))$$

v denotes validation loss.

Learning to Reweight Examples for Robust Deep Learning[5]

$$u_{i,t} = -\eta \frac{\partial}{\partial \epsilon_{i,t}} \frac{1}{m} \sum_{j=1}^m f_j^v(\theta_{t+1}(\epsilon)) \Bigg|_{\epsilon_{i,t}=0}$$
$$w_{i,t} = \frac{|u_{i,t}|}{\sum_j |u_{j,t}|}$$

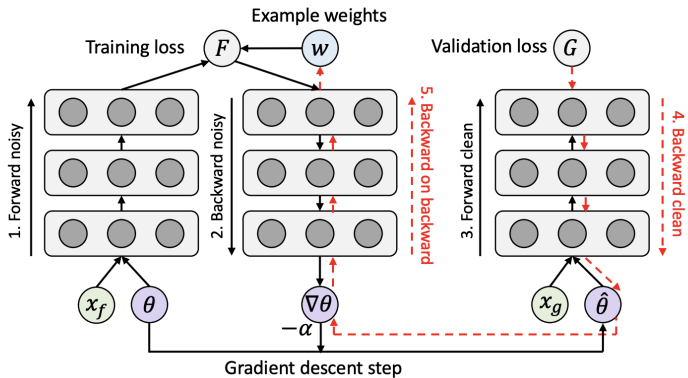


Figure 1. Computation graph of our algorithm in a deep neural network, which can be efficiently implemented using second order automatic differentiation.

Meta-NN

Use a meta-NN $\mathcal{V}(L; \Theta)$ to give the reweight coefficient for loss term L .

$$\mathbf{w}^*(\Theta) = \arg \min_{\mathbf{w}} \mathcal{L}^{\text{train}}(\mathbf{w}; \Theta) \triangleq \frac{1}{N} \sum_{i=1}^N \mathcal{V}(L_i^{\text{train}}(\mathbf{w}); \Theta) \cdot L_i^{\text{train}}(\mathbf{w})$$

Meta-NN

Use a meta-NN $\mathcal{V}(L; \Theta)$ to give the reweight coefficient for loss term L .

$$\mathbf{w}^*(\Theta) = \arg \min_{\mathbf{w}} \mathcal{L}^{\text{train}}(\mathbf{w}; \Theta) \triangleq \frac{1}{N} \sum_{i=1}^N \mathcal{V}(L_i^{\text{train}}(\mathbf{w}); \Theta) \cdot L_i^{\text{train}}(\mathbf{w})$$

Alternatively update \mathbf{w} and Θ to minimize loss

$$\Theta^* = \arg \min_{\Theta} \mathcal{L}^{\text{meta}}(\mathbf{w}^*(\Theta)) \triangleq \frac{1}{M} \sum_{i=1}^M L_i^{\text{meta}}(\mathbf{w}^*(\Theta))$$

Meta-weight-net[6]

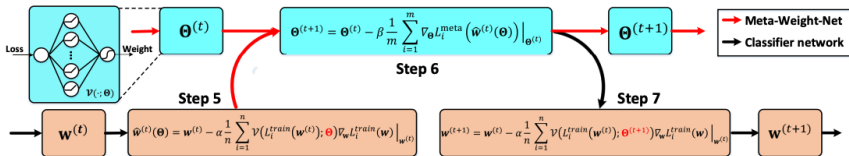


Figure 2: Main flowchart of the proposed MW-Net Learning algorithm (steps 5-7 in Algorithm 1).

Curriculum

$$\min_{\mathbf{w}, \mathbf{v}} \mathbb{F}(\mathbf{w}, \mathbf{v}) = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i^T \mathbf{L}(\mathbf{y}_i, g_s(\mathbf{x}_i, \mathbf{w})) + G(\mathbf{v}; \lambda) + \theta \|\mathbf{w}\|_2^2$$

Each G gives a curriculum.

MentorNet[2]





Use an nn (MentorNet $g(\cdot; \Theta)$) to learn data-driven curriculum :

$$g_m(\mathbf{z}_i; \Theta^*) = \arg \min_{\mathbf{v}_i \in [0,1]} \mathbb{F}(\mathbf{w}, \mathbf{v}), \forall i \in [1, n]$$



or

$$\Theta^* = \arg \min_{\Theta} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}} g_m(\mathbf{z}_i; \Theta) \ell_i + G(g_m(\mathbf{z}_i; \Theta); \lambda)$$

References I

-  Y. CUI, M. JIA, T.-Y. LIN, Y. SONG, AND S. BELONGIE, *Class-balanced loss based on effective number of samples*, in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 9268–9277.
-  L. JIANG, Z. ZHOU, T. LEUNG, L. LI, AND L. FEI-FEI, *Mentornet: Regularizing very deep neural networks on corrupted labels*, CoRR, abs/1712.05055 (2017).
-  B. LI, Y. LIU, AND X. WANG, *Gradient harmonized single-stage detector*, CoRR, abs/1811.05181 (2018).
-  T. LIN, P. GOYAL, R. B. GIRSHICK, K. HE, AND P. DOLLÁR, *Focal loss for dense object detection*, CoRR, abs/1708.02002 (2017).

References II

-  M. REN, W. ZENG, B. YANG, AND R. URTASUN, *Learning to reweight examples for robust deep learning*, arXiv preprint arXiv:1803.09050, (2018).
-  J. SHU, Q. XIE, L. YI, Q. ZHAO, S. ZHOU, Z. XU, AND D. MENG, *Meta-weight-net: Learning an explicit mapping for sample weighting*, arXiv preprint arXiv:1902.07379, (2019).