

题目: [Balanced Multimodal Learning via On-the-fly Gradient Modulation](#)

作者: Xiaokang Peng^{1,†}, Yake Wei^{1,†}, Andong Deng², Dong Wang³, Di Hu^{1,*}

机构: Gaoling School of Artificial Intelligence, Renmin University of China

实验室主页: [GeWu-Lab](#)

动机

工作

Imbalance analysis: dominated modality

On-the-fly Gradient Modulation(OGM)

Generalization Enhancement(GE)

算法过程

实验

问题

限制

动机

- 为所有模态设定一个统一的优化目标会让单模态表征 **under-optimized**, 这可能是主导模态 (dominated modality) 导致的, 如刮风事件里的声音

工作

- 为了缓解第一个问题, 作者提出了 **on-the-fly gradient modulation**, 通过监控每个模态对优化目标的贡献差异, 自适应的控制每种模态的优化, 对于 under-optimized 的模态提供更多帮助
- 作者额外引入了一个动态改变的 **Gaussian noise** 来避免因 gradient modulation 导致的 generalization drop, 因为 gradient modulation 会导致 stochastic gradient noise 降低, 论文[1,2]中证明了 stochastic gradient noise 的强度与 generalization ability 成 positive correlation。

Imbalance analysis: dominated modality

$$\frac{\partial L}{\partial f(x_i)_c} = \frac{e^{(W^a \cdot \varphi_i^a + W^v \cdot \varphi_i^v + b)_c}}{\sum_{k=1}^M e^{(W^a \cdot \varphi_i^a + W^v \cdot \varphi_i^v + b)_k}} - 1_{c=y_i}$$

On-the-fly Gradient Modulation(OGM)

首先对梯度下降GD进行定义：

$$\theta_{t+1}^u = \theta_t^u - \eta \nabla_{\theta^u} L(\theta_t^u) \quad ()$$

实际中，我们使用SGD：

$$\theta_{t+1}^u = \theta_t^u - \eta \tilde{g}(\theta_t^u) \quad ()$$

其中 $\tilde{g}(\theta_t^u) = \frac{1}{m} \sum_{x \in B_t} \nabla_{\theta^u} \ell(x; \theta_t^u)$ ，其是 $\nabla_{\theta^u} L(\theta_t^u)$ 的unbiased estimation，有：

$$\begin{aligned} \tilde{g}(\theta_t^u) &\sim \mathcal{N}(\nabla_{\theta^u} L(\theta_t^u), \Sigma^{sgd}(\theta_t^u)) \\ \Sigma^{sgd}(\theta_t^u) &\approx \frac{1}{m} \left[\frac{1}{N} \sum_{i=1}^N \nabla_{\theta^u} \ell(x_i; \theta_t^u) \nabla_{\theta^u} \ell(x_i; \theta_t^u)^\top \right. \\ &\quad \left. - \nabla_{\theta^u} L(\theta_t^u) \nabla_{\theta^u} L(\theta_t^u)^\top \right]. \end{aligned}$$

接下来为了自适应的控制每种模态的优化，我们要得到每个模态梯度的加权系数，从而调整梯度大小，同步各个模态，

首先得到discrepancy ratio ρ_t^v ：

$$\begin{aligned} s_i^a &= \sum_{k=1}^M 1_{k=y_i} \cdot \text{softmax} \left(W_t^a \cdot \varphi_t^a(\theta^a, x_i^a) + \frac{b}{2} \right)_k, \\ s_i^v &= \sum_{k=1}^M 1_{k=y_i} \cdot \text{softmax} \left(W_t^v \cdot \varphi_t^v(\theta^v, x_i^v) + \frac{b}{2} \right)_k, \\ \rho_t^v &= \frac{\sum_{i \in B_t} s_i^v}{\sum_{i \in B_t} s_i^a}. \end{aligned}$$

然后计算系数 k_t^u ：

$$k_t^u = \begin{cases} 1 - \tanh(\alpha \cdot \rho_t^v) & \rho_t^v > 1 \\ 1 & \text{others} \end{cases}$$

带入SGD更新公式() $：$

$$\theta_{t+1}^u = \theta_t^u - \eta \cdot k_t^u \tilde{g}(\theta_t^u) \quad ()$$

Generalization Enhancement(GE)

- 噪声确实减少

考虑噪声的SGD过程:

$$\theta_{t+1}^u = \theta_t^u - \eta \nabla_{\theta^u} L(\theta_t^u) + \eta \xi_t, \xi_t \sim \mathcal{N}(0, \Sigma^{sgd}(\theta_t^u))$$

添加系数后:

$$\begin{aligned}\theta_{t+1}^u &= \theta_t^u - \eta \nabla_{\theta^u} L'(\theta_t^u) + \eta \xi_t' \\ \xi_t' &\sim \mathcal{N}(0, (k_t^u)^2 \cdot \Sigma^{sgd}(\theta_t^u))\end{aligned}$$

其中 $\eta \nabla_{\theta^u} L'(\theta_t^u) = k_t^u \cdot \eta \nabla_{\theta^u} L(\theta_t^u)$, 而 $k_t^u \in (0, 1]$, 所以 $\xi_t' < \xi_t$, 所以需要增强

- 增加高斯噪声 $h(\theta_t^u) \sim \mathcal{N}(0, \Sigma^{sgd}(\theta_t^u))$

$$\begin{aligned}\theta_{t+1}^u &= \theta_t^u - \eta (k_t^u \tilde{g}(\theta_t^u) + h(\theta_t^u)) \\ &= \theta_t^u - \eta \nabla_{\theta^u} L'(\theta_t^u) + \eta \xi_t' + \eta \epsilon_t \\ &= \theta_t^u - \eta \nabla_{\theta^u} L'(\theta_t^u) + \eta \xi_t''\end{aligned}\tag{0}$$

其中 $\epsilon_t \sim \mathcal{N}(0, \Sigma^{sgd}(\theta_t^u))$, $\xi_t'' \sim \mathcal{N}(0, (k_t^u)^2 + 1) \Sigma^{sgd}(\theta_t^u)$, SGD噪声得以加强

算法过程

```
1 \begin{algorithm}
2 \caption{Multimodal learning with OGM-GE strategy}
3 \label{alg:1}
4 \begin{algorithmic}
5 \Require Training dataset  $\mathcal{D} = \{(\mathbf{x}^a_i, \mathbf{x}^v_i), \mathbf{y}_i\}_{i=1,2,\dots,N}$ ,
iteration number  $T$ , hyper-parameter  $\alpha$ , initialized modal-specific
parameters  $\theta^u$ ,  $\mathbf{u} \in \{a, v\}$ .
6 \For{ $t=0, \dots, T-1$ }
7   \State Sample a fresh mini-batch  $B_t$  from  $\mathcal{D}$ ;
8   \State Feed-forward the batched data  $B_t$  to the model;
9   \State Calculate  $\rho^u$  using Equation-\ref{8} and-\ref{calcu_ratio};
10  \State Calculate  $k_t^u$  using Equation-\ref{k};
11  \State Calculate gradient  $\tilde{g}(\theta_t^u)$  using back-propagation;
12  \State Sample  $h(\theta_t^u)$  based on covariance of gradient  $\tilde{g}(\theta_t^u)$ ;
13  \State Update using  $\theta_{t+1}^u = \theta_t^u - \eta (k_t^u \tilde{g}(\theta_t^u) + h(\theta_t^u))$ .
14 \EndFor
15 \end{algorithmic}
16 \end{algorithm}
```

实验

- 与传统方法对比

作者提出的OGM-GE确实能够提高传统方法的性能

- 与其他modulation策略对比

Dataset	CREMA-D	KS
Method	Acc	Acc
Concatenation	51.7	59.8
Modality-Drop [9] (audio)	54.4	60.3
Modality-Drop [9] (visual)	53.3	61.3
Grad-Blending [39]	56.8	62.2
OGM	59.0	61.1
OGM-GE	61.9	62.3

- 将OGM-GE插入现有方法

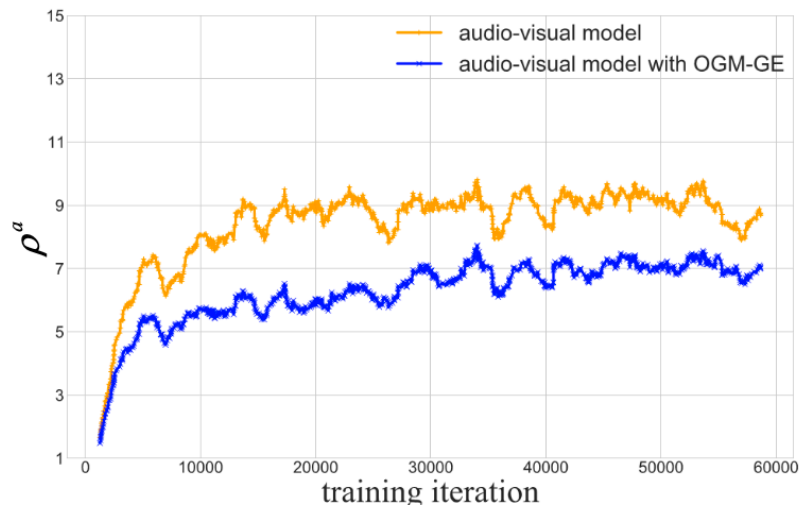
- 分类任务以外的应用

作者在audio-visual event localization任务上进行尝试:

Audio-visual Event Localization		
w/ or w/o OGM-GE	w/o	w/
AVGA [36]	72.0	72.8
PSP [46]	76.2	76.9

- 消融

- more balanced



- optimizer

Dataset	CREMA-D	KS	VGGSound
Method	Acc	Acc	Acc
SGD	51.7	59.8	49.1
SGD†	61.9	63.1	50.6
Adam	49.7	57.4	47.3
Adam†	54.6	58.9	48.2

- different noise intensities

Settings	CREMA-D	VGGSound
(b=64, lr=1e-4)	50.4	48.3
(b=64, lr=5e-4)	51.0	48.7
(b=64, lr=1e-3)	51.8	49.1
(b= 64, lr=1e-3)	51.8	49.1
(b=128, lr=1e-3)	50.2	48.8
(b=256, lr=1e-3)	48.6	47.7
(b= 64, lr=1e-3) w/ GE	60.2	50.3

问题

- CREMA-D数据集里val集合的划分比例在论文是0.1，但所提供的代码仓库里只划分了train/test
- CREMA-D: 0.716 (bs=32训练) > 0.641 (Repo提供的ckpt) > 0.619 (论文)，论文中解释了这个现象，更小的bs会导致更大的高斯噪声，但这个差距实在是太大了

限制

- the uni-modal performance in multimodal model still do not surpass the best uni-modal model.
- solely leveraging optimization-oriented method could not thoroughly solve the imbalance problem, more advanced fusion strategy or network architectures are needed
- 作者只尝试了两个模态, depth, optical flow, language也会有用

[1] Jastrzębski, Stanisław, et al. "Three factors influencing minima in sgd." arXiv preprint arXiv:1711.04623 (2017).

[2] Zhou, Mo, et al. "Toward understanding the importance of noise in training neural networks." International Conference on Machine Learning. PMLR, 2019.