

《Make Sharpness-Aware Minimization Stronger: A Sparsified Perturbation Approach》-NIPS2022

[NIPS2022]<https://arxiv.org/pdf/2210.05177.pdf>

Improved SAM algorithm based on sparse perturbation

1. 背景介绍

SAM的方法有一个缺点，那就是它需要对所有参数进行计算和扰动，这导致了大量的计算开销。为了解决这个问题，论文提出了稀疏SAM (SSAM) 的方法。SSAM的核心思想是通过一个二进制掩码，只对部分参数进行扰动和计算，从而实现稀疏扰动。

为了得到这个稀疏掩码，论文给出了两种方法，一种是基于Fisher信息的方法，另一种是基于动态稀疏训练的方法。论文还从理论上证明了，即使只对部分参数进行扰动和计算，

SSAM也能像原始的SAM一样，以 $O(\log T / \sqrt{T})$ 的速率收敛。

通过这种方法，SSAM不仅可以大大减少计算开销，还能有效地提高模型的泛化能力。实验结果也证实了这一点，表明在CIFAR10、CIFAR100和ImageNet-1K上，SSAM比SAM更高效，性能在只有50%稀疏度的扰动下被保留或者甚至更好。

Fisher Information: measure the information that an observable random variable carries about an unknown parameter of a distribution; 作者使用它来衡量一个weight是否需要扰动以获得更flat的minima。

Dynamic sparse training: pruning unimportant weights。

Rethinking the Perturbation in SAM: SAM perturbs all parameters indiscriminately, but the fact is that merely about 5% parameter space is sharp while the rest is flat。作者发现，有些维度不进行perturbation反而能获得更好的generalization。并且作者还对比了SAM和SGD，发现两者的大部分梯度没有显著差异。由此，作者认为most parameters of the model require no perturbations for achieving the flat minima。

2. Methodology

2.1. SSAM

稀疏SAM (SSAM) 使用一个稀疏二进制掩码来决定哪些参数应该被扰动，也就是将 ϵ 乘稀疏的二进制掩码 m ，目标函数可被改写成：

$$\min_{\mathbf{w}} f_S \left(\mathbf{w} + \rho \cdot \frac{\nabla_{\mathbf{w}} f(\mathbf{w})}{\|\nabla_{\mathbf{w}} f(\mathbf{w})\|_2} \odot \mathbf{m} \right)$$

并且 m 在训练过程中会被定期更新。作者提供了两种获取稀疏掩码 m 的解决方案，分别是基于Fisher信息的稀疏SAM (SSAM-F)和基于动态稀疏训练的稀疏SAM (SSAM-D)。

Algorithm 1 Sparse SAM (SSAM)

Input: sparse ratio s , dense model w , binary mask m , update interval T_m , number of samples N_F , learning rate η , training set \mathcal{S} .

- 1: Initialize w and m randomly.
- 2: **for** epoch $t = 1, 2, \dots, T$ **do**
- 3: **for** each training iteration **do**
- 4: Sample a batch from \mathcal{S} : \mathcal{B}
- 5: Compute perturbation ϵ by Eq. (2)
- 6: **if** $t \bmod T_m = 0$ **then**
- 7: Generate mask m via Option I or II
- 8: **end if**
- 9: $\epsilon \leftarrow \epsilon \odot m$
- 10: **end for**
- 11: $w \leftarrow w - \eta \nabla f(w + \epsilon)$
- 12: **end for**
- 13: **return** Final weight of model w

Algorithm 2 Sparse Mask Generation

- 1: Option I:(Fisher Information Mask)
- 2: Sample N_F data from \mathcal{S} : \mathcal{B}_F
- 3: Compute Fisher \hat{F}_w by Eq. (5)
- 4: $m_1 = \{m_i = 1 | m_i \in m\} \leftarrow \text{ArgTopK}(\hat{F}_w, s \cdot |w|)$
- 5: $m_0 = \{m_i = 0 | m_i \in m\} \leftarrow \{m_i | m_i \notin m_1\}$
- 6: $m \leftarrow m_0 \cup m_1$
- 7: Option II:(Dynamic Sparse Mask)
- 8: $N_{drop} = f_{decay}(t; \alpha) \cdot (1 - s) \cdot |w|$
- 9: $N_{growth} = N_{drop}$
- 10: $m_1 = \{m_i = 1 | m_i \in m\} \leftarrow \{m_i = 1 | m_i \in m\} - \text{ArgTopK}_{m_i \in m_1}(-|\nabla f(w)|, N_{drop})$
- 11: $m_1 \leftarrow \{m_i = 1 | m_i \in m\} + \text{Random}_{m_i \notin m_1}(N_{growth})$
- 12: $m_0 = \{m_i = 0 | m_i \in m\} \leftarrow \{m_i | m_i \notin m_1\}$
- 13: $m \leftarrow m_0 \cup m_1$
- 14: **return** Sparse mask m

2.2. SSAM-F

The Fisher information is defined by:

$$F_w = \mathbb{E}_{x \sim p(x)} [\mathbb{E}_{y \sim p_w(y|x)} \nabla_w \log p_w(y|x) \nabla_w \log p_w(y|x)^T]$$

因为 $p(x)$ 不是available的，所以作者使用sampling来估计：

$$F_w = \frac{1}{N_F} \mathbb{E}_{y \sim p_w(y|x_i)} (\nabla_w \log p_w(y|x_i))^2$$

empirical Fisher:

$$\hat{F}_w = \frac{1}{N_F} (\nabla_w \log p_w(y_i|x_i))^2$$

获得经验费雪信息后按照元素降序排序，取前 k 个元素对应的参数添加扰动，掩码相应位置设置为1.

2.3. SSAM-D

由于费雪信息矩阵计算相对较高，本文也给出了基于动态稀疏训练的算法。动态稀疏训练包括扰动丢弃(Perturbation Dropping)步骤和扰动增长(Perturbation Growth)步骤。

扰动丢弃步骤：丢弃最平坦维度的权重。衡量指标即对应梯度的绝对值。

扰动增长步骤：核心目标是探索尽可能多的扰动组合，主要方法即随机选取一定数量的参数添加扰动。