



# Revisiting Sharpness-Aware Minimization: A More Faithful and Effective Implementation

**Jianlong Chen, Zhiming Zhou\***

Shanghai University of Finance and Economics

April, ICLR 2026

# Sharpness-Aware Minimization (SAM)

SAM has attained significant attention for its potential to enhance the generalization of machine learning models, *in a direct optimization manner*.

## The Objective and Classical Approximation of SAM

- The formal objective:

$$\min_{\theta} \max_{\|\delta\| \leq \rho} L(\theta + \delta)$$

Exactly solving the inner maximization is computationally expensive!

**Main Approximation 1:** SAM approximates best perturbation  $\delta^*$  by gradient ascent(s). Formally,

$$\delta^* = \vartheta_k - \theta = \sum_{i=0}^{k-1} \rho_i \frac{\mathbf{g}_i}{\|\mathbf{g}_i\|}, \quad \text{where } k \geq 1.$$

- After such approximation of  $\delta^*$ , the SAM objective reduces to:

$$\min_{\theta} L(\theta + \delta^*)$$

Directly computing its gradient requires expensive higher-order derivatives!

**Main Approximation 2:** SAM approximates  $\nabla_{\theta} L(\theta + \delta^*)$  using  $\nabla_{\vartheta_k} L(\vartheta_k)$ . Formally,

$$\nabla_{\theta} L(\theta + \delta^*) = \nabla_{\theta} L(\vartheta_k) = \nabla_{\vartheta_k} L(\vartheta_k) \cdot \underbrace{\nabla_{\theta}(\vartheta_k)}_{\text{Approximated as identity matrix } I} \approx \nabla_{\vartheta_k} L(\vartheta_k).$$

- The resulting algorithm essentially applies the gradient at the final ascent point  $\vartheta_k$  to  $\theta$ :

$$\theta_{t+1} = \theta_t - \eta_t \cdot \nabla_{\vartheta_k} L(\vartheta_k).$$

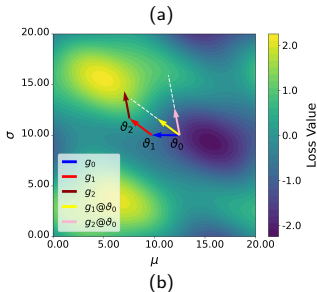
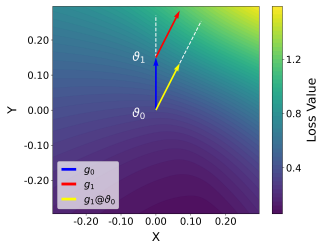
- After several approximations, SAM achieves empirical successes across various tasks using this implementation.
- A body of research has sought to demystify SAM after such approximations.
- However, a direct and intuitive understanding of why applying the nonlocal gradient at the ascent point to update the current parameter works superiorly is still lacking.

$$\theta_{t+1} = \theta_t - \eta_t \cdot \nabla_{\vartheta_k} L(\vartheta_k).$$

**This gap necessitates a deeper investigation into SAM's fundamental mechanisms, which motivates our work.**

# Novel Interpretation of SAM's Mechanism

## Empirical Observations



## Key Mechanism

The gradient  $g_1$  at the single-step ascent point  $\vartheta_1$ , when applied to the current parameters  $\vartheta_0$ , provides a better approximation of the direction from  $\vartheta_0$  toward the maximum within the neighborhood than gradient  $g_0$ .

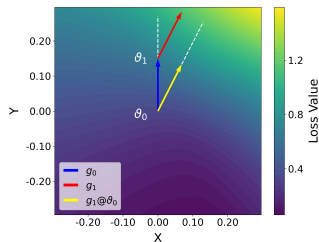
## Theoretical Confirmation

**Proposition.** Under some assumptions, there exists  $\rho_0 > 0$  such that for all  $\rho_m > \rho_0$ :

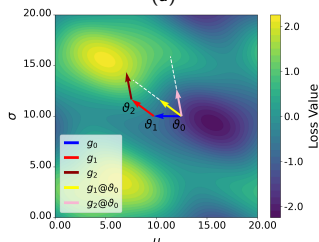
SAM better approximates the direction toward the maximum in the vicinity than SGD

$$L\left(\vartheta_0 + \rho_m \frac{g_1}{\|g_1\|}\right) > L\left(\vartheta_0 + \rho_m \frac{g_0}{\|g_0\|}\right);$$

## Empirical Observations



(a)



(b)

## Limitations

- L1: The approximation by the gradient at the single-step ascent point is rough, inaccurate and unstable.
- L2: The approximation quality may degrade as the number of ascent steps increases.

## Theoretical Confirmation

**Proposition.** Under some assumptions, there exists  $\rho_0 > 0$  such that for all  $\rho_m > \rho_0$ :

There exist better approximations than SAM there exists  $\alpha \in \mathbb{R}$  such that

$$L\left(\vartheta_0 + \rho_m \frac{g_\alpha}{\|g_\alpha\|}\right) > L\left(\vartheta_0 + \rho_m \frac{g_1}{\|g_1\|}\right),$$
$$g_\alpha = \alpha g_1 + (1 - \alpha) g_0.$$

So, how to improve?

# Explicit Sharpness-Aware Minimization (XSAM)

**Takeaway:** XSAM explicitly and dynamically estimates the maximum direction during training.

## Main steps of XSAM

- **Explicit direction probing:** Find the direction of the maximum in the local neighborhood; (line 11)
  - ▷ Constrain the probe to a 2D hyperplane spanned by  $\mathbf{g}_k$  at  $\vartheta_k$  and the vector from  $\vartheta_0$  to  $\vartheta_k$ ; (line 8)
- **Dynamic adaptation:** Update direction during training. (line 10)
- **Guided optimization:** Steer the optimization using the explicitly estimated direction. (line 16)

▶ Full algorithm on the right

## Algorithm 1: XSAM

**Input:**  $\theta_0, T, k \geq 1, \{\rho_i\}, \rho_m, T_\alpha, \{\eta_t\}$

**Output:**  $\theta_T$

```
1: for  $t = 0$  to  $T - 1$  do
2:    $\vartheta_0 = \theta_t$ 
3:   for  $i = 0$  to  $k - 1$  do           ▷ Single-step:  $k = 1$ 
4:      $\mathbf{g}_i = \nabla_{\vartheta_i} L(\vartheta_i)$ 
5:      $\vartheta_{i+1} = \vartheta_i + \rho_i \frac{\mathbf{g}_i}{\|\mathbf{g}_i\|}$ 
6:   end for
7:    $\mathbf{g}_k = \nabla_{\vartheta_k} L(\vartheta_k)$ 
8:    $\mathbf{v}_0 = \frac{\vartheta_k - \vartheta_0}{\|\vartheta_k - \vartheta_0\|}, \mathbf{v}_1 = \frac{\mathbf{g}_k}{\|\mathbf{g}_k\|}$ 
9:    $\psi = \arccos(\mathbf{v}_0 \cdot \mathbf{v}_1)$ 
10:  if  $t \bmod T_\alpha = 0$  then
11:     $\alpha_t^* = \arg \max_{\alpha} L(\vartheta_0 + \rho_m \cdot \mathbf{v}(\alpha)),$ 
12:    where  $\mathbf{v}(\alpha) = \frac{\sin((1-\alpha)\psi)}{\sin(\psi)} \mathbf{v}_0 + \frac{\sin(\alpha\psi)}{\sin(\psi)} \mathbf{v}_1$ 
13:  else
14:     $\alpha_t^* = \alpha_{t-1}^*$ 
15:  end if
16:   $\theta_{t+1} = \theta_t - \eta_t \cdot \mathbf{v}(\alpha_t^*) \cdot \|\mathbf{g}_k\|$ 
17: end for
```

# Experimental Results of Single-step Setting

XSAM consistently improves SAM across different datasets and architectures.

## ■ CIFAR-10

Dataset	CIFAR-10		
Model	VGG-11	ResNet-18	DenseNet-121
SGD	93.19 $\pm$ 0.11	96.15 $\pm$ 0.05	96.34 $\pm$ 0.11
SAM	93.83 $\pm$ 0.06	96.59 $\pm$ 0.06	96.97 $\pm$ 0.02
XSAM	<b>94.25</b> $\pm$ 0.14	<b>96.74</b> $\pm$ 0.04	<b>97.15</b> $\pm$ 0.03

## ■ CIFAR-100

Dataset	CIFAR-100		
Model	VGG-11	ResNet-18	DenseNet-121
SGD	71.46 $\pm$ 0.17	78.55 $\pm$ 0.20	81.78 $\pm$ 0.06
SAM	74.01 $\pm$ 0.05	80.93 $\pm$ 0.11	83.81 $\pm$ 0.02
XSAM	<b>74.21</b> $\pm$ 0.14	<b>81.24</b> $\pm$ 0.07	<b>83.96</b> $\pm$ 0.10

## ■ Tiny-ImageNet

Dataset	Tiny-ImageNet		
Model	VGG-11	ResNet-18	DenseNet-121
SGD	47.44 $\pm$ 0.33	57.02 $\pm$ 0.42	61.93 $\pm$ 0.10
SAM	51.96 $\pm$ 0.26	62.81 $\pm$ 0.09	66.31 $\pm$ 0.09
XSAM	<b>52.58</b> $\pm$ 0.38	<b>63.82</b> $\pm$ 0.23	<b>66.81</b> $\pm$ 0.08

## ■ ImageNet, Transformer, and ViT-Ti

	ImageNet ResNet-50	Transformer IWSLT2014	ViT-Ti CIFAR-100
SAM	77.04 $\pm$ 0.09	35.30 $\pm$ 0.04	67.80 $\pm$ 0.22
XSAM	<b>77.22</b> $\pm$ 0.07	<b>35.63</b> $\pm$ 0.13	<b>68.32</b> $\pm$ 0.18

# Experimental Results of Multi-step Setting

XSAM achieves the best performance among SAM variants.

## ■ Multi-step SAMs on CIFAR-100 with ResNet-18

Methods	$k = 1$	$k = 2$	$k = 4$
SAM	80.93 $\pm$ 0.11	80.91 $\pm$ 0.10	80.65 $\pm$ 0.26
LSAM	80.93 $\pm$ 0.11	80.94 $\pm$ 0.09	80.74 $\pm$ 0.18
LSAM+	80.61 $\pm$ 0.20	80.83 $\pm$ 0.11	80.41 $\pm$ 0.03
MSAM	80.93 $\pm$ 0.11	81.18 $\pm$ 0.06	81.01 $\pm$ 0.09
MSAM+	80.83 $\pm$ 0.05	80.86 $\pm$ 0.34	80.77 $\pm$ 0.08
XSAM	<b>81.27</b> $\pm$ 0.07	<b>81.44</b> $\pm$ 0.09	<b>81.37</b> $\pm$ 0.24

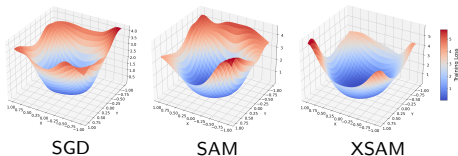
## ■ Multi-step ( $k = 3$ ) SAMs on CIFAR-100 using ResNet-18 across different $\rho$

Method	$\rho = 0.04$	$\rho = 0.08$	$\rho = 0.12$
SAM	80.79 $\pm$ 0.41	80.75 $\pm$ 0.27	79.72 $\pm$ 0.33
LSAM	81.00 $\pm$ 0.21	81.20 $\pm$ 0.24	81.16 $\pm$ 0.04
LSAM+	80.56 $\pm$ 0.20	80.77 $\pm$ 0.04	80.21 $\pm$ 0.27
MSAM	81.04 $\pm$ 0.06	81.12 $\pm$ 0.17	80.93 $\pm$ 0.11
MSAM+	80.72 $\pm$ 0.16	81.16 $\pm$ 0.05	81.16 $\pm$ 0.05
XSAM	<b>81.23</b> $\pm$ 0.06	<b>81.36</b> $\pm$ 0.08	<b>81.29</b> $\pm$ 0.06

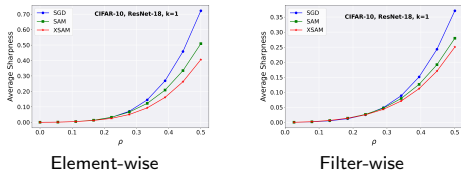
# Additional Experimental Results

## Flatness Analysis

### Visualization of loss landscape

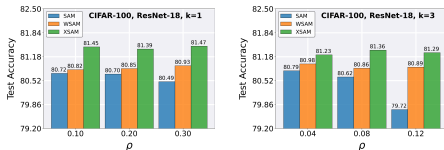


### Average Sharpness (perturbation)

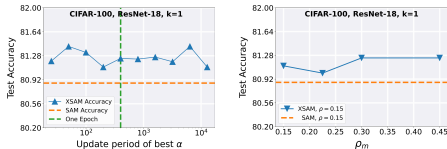


## Robustness Analysis

### Robustness to $\rho$



### Robustness to $T_\alpha$ and $\rho_m$



**Takeaway:** More accurately estimate the direction toward the maximum and then move away along it for improved sharpness-aware minimization.

- SAM better approximates the direction toward the maximum in the vicinity than SGD.
- The classic SAM gradient cannot accurately approximate the direction of the maximum, and may get worse as the number of ascent steps increases.
- XSAM explicitly and dynamically estimate the direction of the maximum, leading to a more faithful and effective sharpness-aware minimization.
- XSAM features a unified formulation that applies to both single-step and multi-step settings and only incurs negligible computational overhead.

# Thank you for your attention!

*If you have any questions, please feel free to contact me.*