



Improving Sharpness-Aware Minimization by Lookahead

Runsheng Yu¹, Youzhi Zhang², James T. Kwok¹

¹ Department of Computer Science and Engineering, Hong Kong University of Science and Technology

² Centre for Artificial Intelligence and Robotics, Hong Kong Institute of Science & Innovation, CAS.

Background

Sharpness-Aware Minimization (SAM)

Flat minima often imply better generalization(Chatterji et al., 2020; Jiang et al., 2020).

SAM (Foret et al., 2021) is designed to locate flat minima:

$$\min_{\boldsymbol{w} \in \mathbb{R}^n} \max_{\boldsymbol{\epsilon}: \|\boldsymbol{\epsilon}\| \leq \rho} L(\boldsymbol{w} + \boldsymbol{\epsilon}),$$

where \boldsymbol{w} is the model parameter, $\boldsymbol{\epsilon}$ is the perturbation whose magnitude is bounded ρ .

Background

Update Scheme for Sharpness-Aware Minimization (SAM)

To solve:

First-order approximation on the objective, the optimal ϵ for the maximization sub-problem: $\epsilon^*(\mathbf{w}) = \frac{\rho \nabla_{\mathbf{w}} L(\mathbf{w})}{\|\nabla_{\mathbf{w}} L(\mathbf{w})\|}$.

Then, the update scheme is:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta \nabla_{\mathbf{w}_{t-1}} L(\mathbf{w}_{t-1} + \epsilon_{t-1}), \epsilon_{t-1} = \epsilon^*(\mathbf{w}_{t-1})$$

However, SAM is easy to be trapped into saddle points (Kim et al., 2023; Compagnoni et al., 2023).

Background: Extra-Gradient and Optimistic Gradient

Consider the minimax problem: $\min_{x \in \mathbb{R}^m} \max_{y \in \mathbb{R}^n} f(x, y)$.

Let $z := [x, y]^\top$ and $F(z) := [\nabla_x f(x, y), -\nabla_y f(x, y)]^\top$

Gradient Descent-Ascent (GDA): $z_{t+1} = z_t - \eta_t F(z_t)$ Unstability (Gidel et al., 2019).

To ensure stability:

Extra-Gradient (Korpelevich, 1976):
$$\begin{aligned} \bar{z}_t &= z_t - \eta_t F(z_t) \\ z_{t+1} &= z_t - \eta_t F(\bar{z}_t) \end{aligned}$$
 Extra extrapolation step, looking one step ahead

Optimistic-Gradient (Popov, 1980):
$$\begin{aligned} \bar{z}_t &= \bar{z}_{t-1} - \eta_t F(\bar{z}_{t-1}) \\ z_{t+1} &= z_t - \eta_t F(\bar{z}_t) \end{aligned}$$
 Reuse gradient at t-1, Reduce computation

Method

To address the issue that SAM is easy to be trapped saddle points,

Direct EG and OG to SAM objective: $\min_{\mathbf{w} \in \mathbb{R}^n} \max_{\boldsymbol{\epsilon}: \|\boldsymbol{\epsilon}\| \leq \rho} L(\mathbf{w} + \boldsymbol{\epsilon}),$

EG

$$\begin{aligned}\hat{\boldsymbol{\epsilon}}_t &= \Pi(\boldsymbol{\epsilon}_{t-1} + \nabla_{\boldsymbol{\epsilon}_{t-1}} L(\mathbf{w}_{t-1} + \boldsymbol{\epsilon}_{t-1})) \\ \hat{\mathbf{w}}_t &= \mathbf{w}_{t-1} - \eta_t \nabla_{\mathbf{w}_{t-1}} L(\mathbf{w}_{t-1} + \boldsymbol{\epsilon}_{t-1}), \\ \boldsymbol{\epsilon}_t &= \Pi(\boldsymbol{\epsilon}_{t-1} + \nabla_{\hat{\boldsymbol{\epsilon}}_t} L(\hat{\mathbf{w}}_t + \hat{\boldsymbol{\epsilon}}_t)) \\ \mathbf{w}_t &= \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\boldsymbol{\epsilon}}_t).\end{aligned}$$

OG

$$\begin{aligned}\hat{\boldsymbol{\epsilon}}_t &= \Pi(\boldsymbol{\epsilon}_{t-1} + \eta'_t \nabla_{\hat{\boldsymbol{\epsilon}}_{t-1}} L(\hat{\mathbf{w}}_{t-1} + \hat{\boldsymbol{\epsilon}}_{t-1})), \\ \hat{\mathbf{w}}_t &= \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_{t-1}} L(\hat{\mathbf{w}}_{t-1} + \hat{\boldsymbol{\epsilon}}_{t-1}), \\ \boldsymbol{\epsilon}_t &= \Pi(\boldsymbol{\epsilon}_{t-1} + \nabla_{\hat{\boldsymbol{\epsilon}}_t} L(\hat{\mathbf{w}}_t + \hat{\boldsymbol{\epsilon}}_t)) \\ \mathbf{w}_t &= \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\boldsymbol{\epsilon}}_t).\end{aligned}$$

$\Pi(\cdot)$ is projection function.

However, directly using both EG and OG which converges at a $O(T^{-\frac{1}{4}})$, which is much slower than the $O(1/\sqrt{T})$ rate of SAM.

Method

A faster convergence by using approximated closed form solution for the max problem.

Lookahead-SAM:

$$\hat{\epsilon}_t = \epsilon^*(\mathbf{w}_{t-1}) \leftarrow \text{closed-form of } \epsilon \text{ w.r.t. } \mathbf{w}.$$

$$\hat{\mathbf{w}}_t = \mathbf{w}_{t-1} - \eta'_t \nabla_{\mathbf{w}_{t-1}} L(\mathbf{w}_{t-1} + \hat{\epsilon}_t) \leftarrow \text{Lookahead step.}$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\epsilon}_t) \leftarrow \text{Update } \mathbf{w} \text{ using the } \hat{\mathbf{w}}_t \text{ with perturbation.}$$

Intuition:

Reduces the perturbation.

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\epsilon}_t)$$

$$\hat{\mathbf{w}}_t = \mathbf{w}_{t-1} - \eta'_t \nabla_{\mathbf{w}_{t-1}} L(\mathbf{w}_{t-1} + \hat{\epsilon}_t)$$

As larger perturbations is prone to being trapped in saddle points (Kim et al., 2023; Compagnoni et al., 2023).

Method

Lookahead-SAM

Optimistic Lookahead-SAM (Opt-SAM):
to further save computation time by reuse gradient

$\hat{\epsilon}_t = \epsilon^* (\mathbf{w}_{t-1})$		$\hat{\epsilon}_t = \epsilon^* (\mathbf{w}_{t-1})$	Reduces the perturbation.
$\hat{\mathbf{w}}_t = \mathbf{w}_{t-1} - \eta'_t \nabla_{\mathbf{w}_{t-1}} L(\mathbf{w}_{t-1} + \hat{\epsilon}_t)$		$\hat{\mathbf{w}}_t = \mathbf{w}_{t-1} - \eta'_t \nabla_{\hat{\mathbf{w}}_{t-1}} L(\hat{\mathbf{w}}_{t-1} + \hat{\epsilon}_{t-1})$	↓
$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\epsilon}_t)$		$\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \nabla_{\hat{\mathbf{w}}_t} L(\hat{\mathbf{w}}_t + \hat{\epsilon}_t)$	

However, Opt-SAM still has to compute the gradient in each iteration, and can be expensive.

Background: AE-SAM (Jiang et al., 2023)

A SAM variant to reduce computation.

Empirically, $\|\nabla L(\mathbf{w}_t)\|^2 \sim \mathcal{N}(\mu_t, \sigma_t)$

$$\begin{aligned}\mu_t &= \delta\mu_{t-1} + (1 - \delta) \|\nabla L(\mathbf{w}_t)\|^2 \\ \sigma_t^2 &= \delta\sigma_{t-1}^2 + (1 - \delta) \left(\|\nabla L(\mathbf{w}_t)\|^2 - \mu_t \right)^2\end{aligned} \quad \leftarrow \text{Exponential Moving Average}$$

$\|\nabla L(\mathbf{w}_t)\|^2 \geq \mu_t + c_t\sigma_t$ SAM is chosen; Otherwise ERM.

$$c_t = \frac{t}{T}\kappa_1 + \left(1 - \frac{t}{T}\right)\kappa_2, \quad \kappa_1 \text{ and } \kappa_2 \text{ are two constants.}$$

Method: Adaptive Lookahead-SAM

Use AE-SAM technique to help reduce computation.

If $\|\frac{1}{b} \sum_{i \in I_t} \nabla_{\mathbf{w}_t} \ell_i(\mathbf{w}_t)\|^2 \geq \mu_t + c_t \sigma_t$ we use Opt-SAM.
Otherwise, SGD (i.e., ERM) is used instead.

Algorithm 2: Adaptive Optimistic SAM (AO-SAM).

Input: Training set S , number of epochs T , batch size b , \mathbf{w}_0 , $\epsilon_0 = 0$, $\mu_0 = 0$, and $\sigma_0 = e^{-10}$.

```
1 for  $t = 1, 2, \dots, T$  do
2   sample a minibatch  $I_t$  from  $S$  with size  $b$ ;
3    $\mathbf{g}_t = \frac{1}{b} \sum_{i \in I_t} \nabla_{\mathbf{w}_t} \ell_i(\mathbf{w}_t)$ ;
4   update  $\mu_t$  and  $\sigma_t$  as in AE-SAM (5);
5   if  $\|\frac{1}{b} \sum_{i \in I_t} \nabla_{\mathbf{w}_t} \ell_i(\mathbf{w}_t)\|^2 \geq \mu_t + c_t \sigma_t$  then
6      $\hat{\epsilon}_t = \frac{\rho_t \nabla_{\mathbf{w}_{t-1}} \frac{1}{b} \sum_{i \in I_t} \ell_i(\mathbf{w}_{t-1})}{\|\nabla_{\mathbf{w}_{t-1}} \frac{1}{b} \sum_{i \in I_t} \ell_i(\mathbf{w}_{t-1})\|}$ ;
7      $\hat{\mathbf{w}}_t = \mathbf{w}_{t-1} - \eta'_t \mathbf{g}_{t-1}$ ;
8      $\mathbf{g}_t = \nabla_{\hat{\mathbf{w}}_t} [\frac{1}{b} \sum_{i \in I_t} \ell_i(\hat{\mathbf{w}}_t + \hat{\epsilon}_t)]$ ;
9    $\mathbf{w}_t = \mathbf{w}_{t-1} - \eta_t \mathbf{g}_t$ ;
10 return  $\mathbf{w}_T$ .
```

Analysis: Region of attraction (ROA)

ROA (informal): The Region of Attraction is all the starting points from which the system will eventually settle into this stable state.

Objective: $\min_{\mathbf{w}} \mathbf{w}^\top \mathbf{H} \mathbf{w}$

ODE SAM:

$$d\mathbf{w}_\tau = -\mathbf{H} \left(\mathbf{w}_\tau + \frac{\rho \mathbf{H} \mathbf{w}_\tau}{\|\mathbf{H} \mathbf{w}_\tau\|} \right) d\tau$$

ROA for ODE SAM:

$$\left\{ \mathbf{w}_\tau \left| \rho \geq -\frac{\|\mathbf{H} \mathbf{w}_\tau\|}{\lambda_{\min}} \right. \right\}$$

λ_{\min} is the minimum eigenvalue of \mathbf{H} .

$$d\mathbf{w}_\tau = -\mathbf{H} \left(\mathbf{w}_\tau - \eta'_\tau \mathbf{H} \left(\mathbf{w}_\tau + \frac{\rho \mathbf{H} \mathbf{w}_\tau}{\|\mathbf{H} \mathbf{w}_\tau\|} \right) + \frac{\rho \mathbf{H} \mathbf{w}_\tau}{\|\mathbf{H} \mathbf{w}_\tau\|} \right) d\tau$$

ROA for ODE Lookahead-SAM:

$$\left\{ \mathbf{w}_\tau \left| (1 + \eta'_\tau \lambda_{\min}) \rho \geq \frac{1}{\lambda_{\min}} \|\mathbf{H} \mathbf{w}_\tau\| (\eta'_\tau \lambda_{\min} - 1) \right. \right\}$$

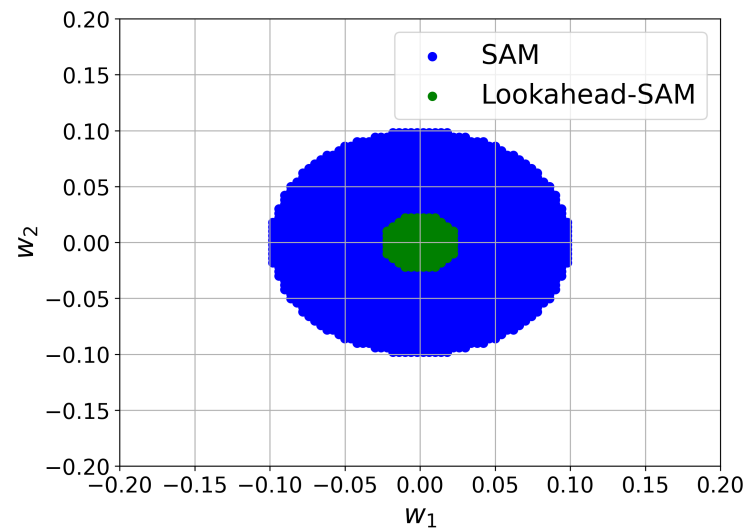
Analysis: Region of attraction (ROA)

ROA for ODE SAM:

$$\left\{ \mathbf{w}_\tau \left| \rho \geq -\frac{\|\mathbf{H}\mathbf{w}_\tau\|}{\lambda_{\min}} \right. \right\}$$

ROA for ODE Lookahead-SAM:

$$\left\{ \mathbf{w}_\tau \left| (1 + \eta'_\tau \lambda_{\min}) \rho \geq \frac{1}{\lambda_{\min}} \|\mathbf{H}\mathbf{w}_\tau\| (\eta'_\tau \lambda_{\min} - 1) \right. \right\}$$



ROAs for SAM and Lookahead-SAM at saddle point

Analysis: Convergence Analysis

Theorem (Informal): under mild assumptions, both Lookahead-SAM, Opt-SAM, and AO-SAM satisfies

$$\frac{1}{T} \sum_{t=0}^T \mathbb{E} \|\nabla_{\mathbf{w}_t} L(\mathbf{w}_t)\|^2 = O\left(\frac{1}{\sqrt{Tb}}\right)$$

.

Lookahead-SAM, Opt-SAM and AO-SAM have the same $O(\frac{1}{\sqrt{Tb}})$ convergence rate as SAM and its variant AESAM.

Experiments: Comparison with EG,OG and the proposed methods

Table 1: Testing accuracy and fraction of SAM updates (%SAM) on *CIFAR-10* using *ResNet-18*. The best accuracy is in bold.

	accuracy	%SAM
SAM	96.52 \pm 0.12	100.0 \pm 0.0
EG	96.45 \pm 0.05	200.0 \pm 0.0
OG	96.52 \pm 0.03	100.0 \pm 0.0
Lookahead-SAM	96.81 \pm 0.01	150.0 \pm 0.0
Opt-SAM	96.79 \pm 0.02	100.0 \pm 0.0
AO-SAM	96.82 \pm 0.04	61.1 \pm 0.0

Table 2: Testing accuracy and fraction of SAM updates (%SAM) on *CIFAR-100* using *ResNet-18*. The best accuracy is in bold.

	accuracy	%SAM
SAM	80.17 \pm 0.05	100.0 \pm 0.0
EG	79.91 \pm 0.16	200.0 \pm 0.0
OG	79.92 \pm 0.08	100.0 \pm 0.0
Lookahead-SAM	80.79 \pm 0.13	150.0 \pm 0.0
Opt-SAM	80.76 \pm 0.15	100.0 \pm 0.0
AO-SAM	80.70 \pm 0.14	61.2 \pm 0.0

Lookahead-SAM has the highest accuracy on CIFAR-100 and the second highest on CIFAR-10, it also has a higher %SAM.

Opt-SAM is as fast as SAM w.r.t. %SAM but is more accurate.

AO-SAM is as accurate as Opt-SAM, but is even faster.

Experiments: Comparison with SAM variants in CIFAR

Table 3: Testing accuracies (mean and standard deviation) and fractions of SAM updates on *CIFAR-10* and *CIFAR-100*. Methods with similar %SAM’s are grouped together for easier comparison. Results of ERM, SAM, and ESAM are from (Jiang et al., 2023), while the other baseline results are obtained with the corresponding authors’ codes. The best accuracy is in bold. * means the improvements over the second-best baseline are statistically significant (achieving a p-value of less than 0.05 in t-test).

		<i>CIFAR-10</i>		<i>CIFAR-100</i>	
		Accuracy	% SAM	Accuracy	% SAM
<i>ResNet-18</i>	ERM	95.41 \pm 0.03	0.0 \pm 0.0	78.17 \pm 0.05	0.0 \pm 0.0
	SAM (Foret et al., 2021)	96.52 \pm 0.12	100.0 \pm 0.0	80.17 \pm 0.15	100.0 \pm 0.0
	ESAM (Du et al., 2022a)	96.56 \pm 0.08	100.0 \pm 0.0	80.41 \pm 0.10	100.0 \pm 0.0
	ASAM (Kwon et al., 2021)	96.55 \pm 0.14	100.0 \pm 0.0	80.52 \pm 0.13	100.0 \pm 0.0
	GSAM (Zhuang et al., 2022)	96.70 \pm 0.01	100.0 \pm 0.0	80.48 \pm 0.11	100.0 \pm 0.0
	Opt-SAM	96.79 \pm 0.02	100.0 \pm 0.0	80.76* \pm 0.15	100.0 \pm 0.0
	SS-SAM (Zhao, 2022)	96.64 \pm 0.02	60.0 \pm 0.0	80.49 \pm 0.10	60.0 \pm 0.0
	AE-SAM (Jiang et al., 2023)	96.66 \pm 0.02	61.3 \pm 0.1	79.96 \pm 0.08	61.3 \pm 0.0
	AO-SAM	96.82* \pm 0.04	61.1 \pm 0.0	80.70 \pm 0.14	61.2 \pm 0.0
<i>WideResNet-28-10</i>	ERM	96.34 \pm 0.12	0.0 \pm 0.0	81.56 \pm 0.14	0.0 \pm 0.0
	SAM (Foret et al., 2021)	97.27 \pm 0.11	100.0 \pm 0.0	83.42 \pm 0.05	100.0 \pm 0.0
	ESAM (Du et al., 2022a)	97.29 \pm 0.11	100.0 \pm 0.0	84.51 \pm 0.02	100.0 \pm 0.0
	ASAM (Kwon et al., 2021)	97.38 \pm 0.09	100.0 \pm 0.0	84.48 \pm 0.10	100.0 \pm 0.0
	GSAM (Zhuang et al., 2022)	97.44 \pm 0.07	100.0 \pm 0.0	84.50 \pm 0.12	100.0 \pm 0.0
	Opt-SAM	97.56* \pm 0.03	100.0 \pm 0.0	84.74 \pm 0.02	100.0 \pm 0.0
	SS-SAM (Zhao, 2022)	97.32 \pm 0.03	60.0 \pm 0.0	84.39 \pm 0.04	60.0 \pm 0.0
	AE-SAM (Jiang et al., 2023)	97.37 \pm 0.08	61.3 \pm 0.0	84.23 \pm 0.08	61.3 \pm 0.0
	AO-SAM	97.49 \pm 0.02	61.2 \pm 0.0	84.80* \pm 0.11	61.2 \pm 0.0

Opt-SAM
and AO-SAM are consistently
more accurate than SAM and
its variants on all datasets
and backbones

Experiments: Comparison with SAM variants in CIFAR

		<i>CIFAR-10</i>		<i>CIFAR-100</i>	
		Accuracy	% SAM	Accuracy	% SAM
<i>PyramidNet-110</i>	ERM	96.62 \pm 0.10	0.0 \pm 0.0	81.89 \pm 0.15	0.0 \pm 0.0
	SAM (Foret et al., 2021)	97.30 \pm 0.10	100.0 \pm 0.0	84.46 \pm 0.05	100.0 \pm 0.0
	ESAM (Du et al., 2022a)	97.81 \pm 0.01	100.0 \pm 0.0	85.56 \pm 0.05	100.0 \pm 0.0
	ASAM (Kwon et al., 2021)	97.71 \pm 0.09	100.0 \pm 0.0	85.55 \pm 0.11	100.0 \pm 0.0
	GSAM (Zhuang et al., 2022)	97.74 \pm 0.02	100.0 \pm 0.0	85.25 \pm 0.11	100.0 \pm 0.0
	Opt-SAM	97.79 \pm 0.04	100.0 \pm 0.0	85.74* \pm 0.14	100.0 \pm 0.0
	SS-SAM (Zhao, 2022)	97.62 \pm 0.03	60.0 \pm 0.0	85.41 \pm 0.11	60.0 \pm 0.0
	AE-SAM (Jiang et al., 2023)	97.52 \pm 0.07	61.4 \pm 0.1	85.43 \pm 0.08	61.4 \pm 0.1
	AO-SAM	97.87* \pm 0.02	61.2 \pm 0.0	85.60 \pm 0.07	61.2 \pm 0.12
<i>ViT-S16</i>	ERM	86.69 \pm 0.11	0.0 \pm 0.0	62.42 \pm 0.22	0.0 \pm 0.0
	SAM (Foret et al., 2021)	87.37 \pm 0.09	100.0 \pm 0.0	63.23 \pm 0.25	100.0 \pm 0.0
	ESAM (Du et al., 2022a)	84.27 \pm 0.11	100.0 \pm 0.0	62.11 \pm 0.15	100.0 \pm 0.0
	ASAM (Kwon et al., 2021)	82.25 \pm 0.41	100.0 \pm 0.0	63.26 \pm 0.18	100.0 \pm 0.0
	GSAM (Zhuang et al., 2022)	83.62 \pm 0.11	100.0 \pm 0.0	59.82 \pm 0.12	100.0 \pm 0.0
	Opt-SAM	87.91 \pm 0.26	100.0 \pm 0.0	63.78 \pm 0.22	100.0 \pm 0.0
	SS-SAM (Zhao, 2022)	83.36 \pm 0.04	60.0 \pm 0.0	54.04 \pm 5.09	60.0 \pm 0.0
	AE-SAM (Jiang et al., 2023)	77.37 \pm 0.07	61.4 \pm 0.0	57.13 \pm 2.87	61.3 \pm 0.0
	AO-SAM	88.27* \pm 0.12	61.3 \pm 0.0	64.45* \pm 0.23	61.2 \pm 0.0

Opt-SAM and AO-SAM are consistently more accurate than SAM and its variants on all datasets and backbones

Experiments: Comparison with SAM variants in CIFAR with noise

Table 8: Testing accuracies and fractions of SAM updates on *CIFAR-10* with different levels of label noise. Results of ERM, SAM, and ESAM with *ResNet-18* and *ResNet-32* are from (Jiang et al., 2023) (standard derivations for some baselines are not provided in (Jiang et al., 2023)), while the other baseline results are obtained with the authors’ codes. The best accuracy is in bold.

		noise = 20%		noise = 40%		noise = 60%		noise = 80%	
		accuracy	%SAM	accuracy	%SAM	accuracy	%SAM	accuracy	%SAM
<i>ResNet-18</i>	ERM	87.92	0.0	70.82	0.0	49.61	0.0	28.23	0.0
	SAM (Foret et al., 2021)	94.80	100.0	91.50	100.0	88.15	100.0	77.40	100.0
	ESAM (Du et al., 2022a)	94.19	100.0	91.46	100.0	81.30	100.0	15.00	100.0
	ASAM (Kwon et al., 2021)	91.17 \pm 0.19	100.0	87.38 \pm 0.61	100.0	83.22 \pm 0.41	100.0	71.03 \pm 0.88	100.0
	GSAM (Zhuang et al., 2022)	94.54 \pm 0.18	100.0	91.72 \pm 0.05	100.0	87.70 \pm 0.02	100.0	24.70 \pm 10.69	100.0
	Opt-SAM	95.12 \pm 0.12	100.0	92.16 \pm 0.35	100.0	88.45 \pm 0.53	100.0	77.47 \pm 0.65	100.0
	SS-SAM (Zhao, 2022)	94.61 \pm 0.16	60.0	91.81 \pm 0.13	60.0	78.67 \pm 0.42	60.0	62.94 \pm 1.01	60.0
	AE-SAM (Jiang et al., 2023)	92.13 \pm 0.14	61.4	86.02 \pm 0.62	61.4	75.95 \pm 1.30	61.4	67.28 \pm 1.66	61.4
	AO-SAM	95.02 \pm 0.04	61.2	92.62 \pm 0.18	61.3	89.36 \pm 0.12	61.2	78.12 \pm 0.38	61.2
<i>ResNet-32</i>	ERM	87.43	0.0	70.82	0.0	46.26	0.0	29.00	0.0
	SAM (Foret et al., 2021)	95.08	100.0	91.01	100.0	88.90	100.0	77.32	100.0
	ESAM (Du et al., 2022a)	93.42	100.0	91.63	100.0	82.73	100.0	10.09	100.0
	ASAM (Kwon et al., 2021)	92.04 \pm 0.09	100.0	88.83 \pm 0.11	100.0	83.90 \pm 0.56	100.0	75.64 \pm 0.75	100.0
	GSAM (Zhuang et al., 2022)	94.12 \pm 0.09	100.0	91.74 \pm 0.05	100.0	89.23 \pm 0.06	100.0	31.16 \pm 2.77	100.0
	Opt-SAM	95.25 \pm 0.04	100.0	92.11 \pm 0.07	100.0	88.36 \pm 0.22	100.0	77.61 \pm 0.39	100.0
	SS-SAM (Zhao, 2022)	95.03 \pm 0.23	60.0	90.59 \pm 0.30	60.0	87.22 \pm 0.46	60.0	48.89 \pm 1.02	60.0
	AE-SAM (Jiang et al., 2023)	92.04 \pm 0.27	61.3	86.83 \pm 0.49	61.3	73.90 \pm 0.44	61.2	67.64 \pm 1.34	61.3
	AO-SAM	95.32 \pm 0.12	61.2	91.73 \pm 0.65	61.2	89.40 \pm 0.44	61.2	77.78 \pm 0.84	61.2

AO-SAM and Opt-SAM outperform all baselines at all label noise ratios.

Experiments: Comparison with SAM variants in ImageNet

Table 7: Testing accuracies and fractions of SAM updates (%SAM) on *ImageNet*. Results of ERM, SAM and ESAM on *ResNet-50* are from (Jiang et al., 2023), ASAM is from (Kwon et al., 2021), GSAM is from (Zhuang et al., 2022), while the other baseline results are obtained by the corresponding authors’ codes. The best accuracy is in bold. † means that the original papers do not provide standard deviation. We do not report ASAM on *ResNet-101* and *Vit-S/32*, and GSAM on *Vit-S/32* because they are not provided in the original papers.

		Accuracy	%SAM
<i>ResNet-50</i>	ERM	77.11 ± 0.14	0.0
	SAM	77.47 ± 0.12	100.0
	ESAM	77.25 ± 0.75	100.0
	ASAM	76.63 ± 0.18	100.0
	GSAM	77.2^\dagger	100.0
	AO-SAM	77.68 ± 0.04	61.1

AO-SAM again outperforms all the baselines.

		Accuracy	%SAM
<i>ResNet-101</i>	ERM	77.80^\dagger	0.0
	SAM	78.90^\dagger	100.0
	ESAM	79.09^\dagger	100.0
	GSAM	78.9^\dagger	100.0
	AO-SAM	79.38 ± 0.10	61.2
<i>Vit-S/32</i>	ERM	67.0^\dagger	0.0
	SAM	69.1^\dagger	100.0
	ESAM	66.1^\dagger	100.0
	AO-SAM	69.38 ± 0.24	61.6

Experiments: Flat Minima

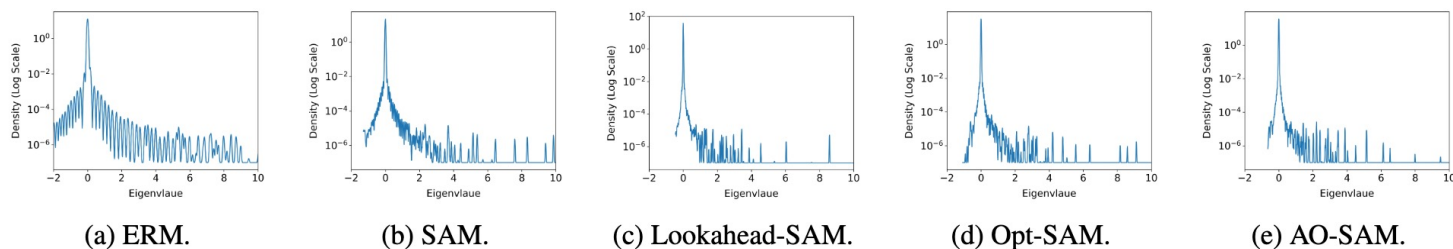


Figure 4: Hessian spectra obtained by ERM, SAM, Lookahead-SAM, Opt-SAM, and AO-SAM on *CIFAR-10* with *ResNet18*.

Table 5: Eigenvalues of the Hessian on *CIFAR-10* with *ResNet18* backbone. The smallest is in bold.

	λ_1	λ_1/λ_5
ERM	88.8	3.3
SAM	29.6	3.3
Lookahead-SAM	10.2	1.8
Opt-SAM	13.1	2.0
AO-SAM	11.1	1.8

The eigenvalues of Lookahead-SAM, Opt-SAM, and AO-SAM are smaller than ERM and SAM.

This indicates the loss landscapes at the converged solutions of these SAM variants are flatter.

Conclusion

1. Incorporate the idea of extrapolation into SAM to gain more information about the landscape, and thus help convergence.
2. Develop a method that combines SAM's approximate maximizer to its inner optimization subproblem with lookahead.
3. Provide theoretical guarantees that they converge to stationary points at the same rate as SAM, and are not easily trapped at saddle points..