

SASSHA: Sharpness-aware Adaptive Second-order Optimization with Stable Hessian Approximation

Dahun Shin*¹ Dongyeop Lee*¹ Jinseok Chung¹ Namhoon Lee¹

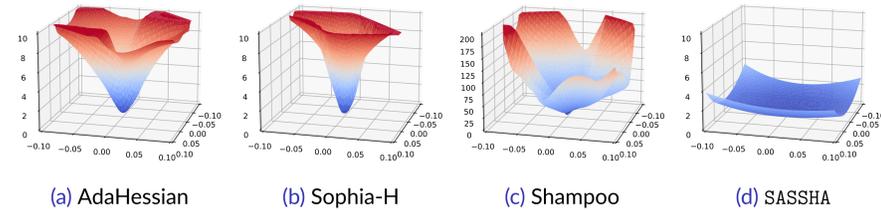
¹Pohang University of Science and Technology



Main Contribution

- observe that approximate second-order methods tend to **converge to sharp minima**, which may explain their **degraded generalization performance**.
- propose SASSHA, an adaptive second-order method designed to **stably converge to flat solutions**, thereby **improving generalization**.
- SASSHA demonstrates strong performance across diverse deep learning tasks, including vision, language, and robustness to label noise.
- SASSHA achieves this **efficiently via lazy Hessian updates without performance degradation**, due to its trajectory staying within regions of slowly changing curvature.

Second-order methods converge to sharp minima



	Sharpness				Generalization	
	$\lambda_{max}(H)$	$tr(H) \times 10^3$	δL_{grad}	$\delta L_{avg} \times 10^{-3}$	L_{val}	Acc _{val} (%)
SGD	265 ± 25.00	7.290 ± 0.300	0.703 ± 0.132	1.310 ± 1.030	1.260 ± 0.001	69.32 ± 0.19
AdaHessian	11992 ± 5779	46.94 ± 17.60	4.119 ± 1.136	12.50 ± 6.080	1.377 ± 0.070	68.06 ± 0.22
Sophia-H	22797 ± 10857	68.15 ± 20.19	8.130 ± 3.082	19.19 ± 6.380	1.463 ± 0.022	67.76 ± 0.37
Shampoo	436374 ± 9017	6823 ± 664.7	73.27 ± 12.49	307489 ± 56979794	1.386 ± 0.010	64.08 ± 0.46
SASSHA	107 ± 40.00	1.870 ± 0.650	0.238 ± 0.088	0.650 ± 0.860	0.961 ± 0.005	72.14 ± 0.16

- Approximate second-order optimizers tend to yield minima of **high sharpness** and **worse generalization** compared to SGD. However, SASSHA effectively recovers this^a.
- provide theoretical explanation for this phenomenon through linear stability analysis:

Corollary 4.6: The linearly stable fixed point x^* of SASSHA satisfies the following necessary conditions:

$$0 \leq a(1 + \rho a) \leq \frac{2\epsilon}{\eta}, \quad 0 \leq s_2^2 \leq \frac{\epsilon^2}{\eta^2 - 2\eta\rho\epsilon}, \quad 0 \leq s_3^2 \leq \frac{\epsilon^2}{2\eta^2\rho}, \quad 0 \leq s_4^2 \leq \frac{\epsilon^2}{\eta^2\rho^2},$$

where $a = \lambda_{max}(\mathbb{E}[H_\xi])$ and $s_k = \lambda_{max}((\mathbb{E}[H_\xi^k] - \mathbb{E}[H_\xi]^k)^{1/k})$ are the sharpness and the non-uniformity of the stochastic Hessian measured with the k -th moment, respectively.

In contrast, standard approximate second-order methods with $P \approx H^{-1}$ impose **no constraint on sharpness**.

Method: SASSHA

Update rule: We propose SASSHA, which performs the following update:

$$x_{t+1} = x_t - \widehat{H}_t^{-1} \nabla f_t(x_t)$$

where \widehat{H} denotes the diagonal of the Hessian, approximated via Hutchinson method.

^aWe find that the same trend holds for other workloads.

- Sharpness minimization:** $\tilde{x}_t = x_t + \rho \frac{\nabla f(x_t)}{\|\nabla f(x_t)\|_2}$.
 - encourage convergence to flat minima.
 - pose a risk of divergence by penalizing Hessian entries on average.
- Stable Hessian approximation:** $\widehat{H} \rightarrow |\widehat{H}|^{1/2}$
 - alleviate divergence and preserve the relative scale of the Hessian.
 - can be interpreted as a geometrical interpolation $H^\alpha (0 \leq \alpha \leq 1)$ which can balance between bias and variance of the population risk.
 - avoid saddle points or local maxima.
- Lazy Hessian update:** reusing previously computed Hessian for several iterations without performance degradation.

Algorithm 1 SASSHA algorithm

```

1: Input: Initial parameter  $x_0$ , learning rate  $\{\eta_t\}$ , moving average parameters  $\beta_1, \beta_2$ , Hessian update interval  $k$ , weight decay parameter  $\lambda$ 
2: Set  $m_{-1} = 0, D_{-1} = 0$ 
3: for  $t = 1$  to  $T$  do
4:    $g_t = \nabla f_B(x_t)$ 
5:    $\epsilon_t^* = \rho g_t / \|g_t\|_2$ 
6:    $\tilde{g}_t = \nabla f_B(x_t + \epsilon_t^*)$ 
7:    $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \tilde{g}_t$ 
8:    $\bar{m}_t = m_t / (1 - \beta_1^t)$ 
9:   if  $t \bmod k = 1$  then
10:     $\widehat{H}_t = \widehat{H}(x_t + \epsilon_t^*)$ 
11:     $D_t = \beta_2 D_{t-1} + (1 - \beta_2) |\widehat{H}_t|$ 
12:     $\bar{D}_t = \sqrt{D_t / (1 - \beta_2^t)}$ 
13:   else
14:     $\bar{D}_t = \bar{D}_{t-1}$ 
15:   end if
16:    $x_{t+1} = x_t - \eta_t \bar{D}_t^{-1} \bar{m}_t - \eta_t \lambda x_t$ 
17: end for
    
```

Theorem 4.4 (Convergence guarantee): Under standard smoothness and convexity assumptions, let $\{x_t\}$ be generated by SASSHA from given any initial point $x_0 \in \mathbb{R}^d$, with step sizes η_t and perturbation radii ρ_t satisfying $\sum_{t=1}^{\infty} \eta_t = \infty$, $\sum_{t=1}^{\infty} \eta_t^2 < \infty$, $\sum_{t=1}^{\infty} \rho_t^2 \eta_t < \infty$. Then,

$$\liminf_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$$

Experiments

Table 1. Image classification results.

Category	Method	CIFAR-10		CIFAR-100		ImageNet	
		ResNet-20	ResNet-32	ResNet-32	WRN-28-10	ResNet-50	ViT-s-32
First-order	SGD	92.03 ± 0.32	92.69 ± 0.06	69.32 ± 0.19	80.06 ± 0.15	75.58 ± 0.05	62.90 ± 0.36
	AdamW	92.04 ± 0.11	92.42 ± 0.13	68.78 ± 0.22	79.09 ± 0.35	75.38 ± 0.08	66.46 ± 0.15
	SAM _{SGD}	92.85 ± 0.07	93.89 ± 0.13	71.99 ± 0.20	83.14 ± 0.13	76.36 ± 0.16	64.54 ± 0.63
	SAM _{AdamW}	92.77 ± 0.29	93.45 ± 0.24	71.15 ± 0.37	82.88 ± 0.31	76.35 ± 0.16	68.31 ± 0.17
Second-order	AdaHessian	92.00 ± 0.17	92.48 ± 0.15	68.06 ± 0.22	76.92 ± 0.26	73.64 ± 0.16	66.42 ± 0.23
	Sophia-H	91.81 ± 0.27	91.99 ± 0.08	67.76 ± 0.37	79.35 ± 0.24	72.06 ± 0.49	62.44 ± 0.36
	Shampoo	88.55 ± 0.83	90.23 ± 0.24	64.08 ± 0.46	74.06 ± 1.28	*	*
	SASSHA	92.98 ± 0.05	94.09 ± 0.24	72.14 ± 0.16	83.54 ± 0.08	76.43 ± 0.18	69.20 ± 0.30

- SASSHA consistently **outperforms** the other methods for all workloads.

Table 2. Language pretraining and finetuning results.

Pretrain / GPT1-mini	Perplexity	Finetune / SqueezeBERT							
		SST-2 Acc	MRPC Acc / F1	STS-B S/P corr.	QQP F1 / Acc	MNLI mat/m.mat	QNLI Acc	RTE Acc	
AdamW	175.06	90.29 ± 0.32	84.56 ± 0.25 / 88.99 ± 0.11	88.34 ± 0.15 / 88.48 ± 0.20	89.92 ± 0.05 / 86.58 ± 0.11	81.22 ± 0.07 / 82.26 ± 0.05	89.93 ± 0.14	68.95 ± 0.72	
SAM _{AdamW}	158.06	90.52 ± 0.27	83.25 ± 0.29 / 87.90 ± 0.21	88.38 ± 0.01 / 88.79 ± 0.09	90.26 ± 0.28 / 86.99 ± 0.31	81.56 ± 0.18 / 82.46 ± 0.19	90.38 ± 0.05	68.83 ± 1.46	
AdaHessian	407.69	89.64 ± 0.13	79.74 ± 0.10 / 85.26 ± 0.30	86.08 ± 0.04 / 86.46 ± 0.06	90.37 ± 0.05 / 87.07 ± 0.05	81.33 ± 0.17 / 82.08 ± 0.02	89.94 ± 0.12	71.00 ± 1.04	
Sophia-H	157.60	90.44 ± 0.46	85.78 ± 0.07 / 89.90 ± 0.82	88.17 ± 0.07 / 88.53 ± 0.13	90.70 ± 0.04 / 87.60 ± 0.06	81.77 ± 0.18 / 82.36 ± 0.22	90.12 ± 0.14	70.76 ± 1.44	
SASSHA	122.40	90.44 ± 0.08	86.28 ± 0.28 / 90.13 ± 0.10	88.72 ± 0.75 / 89.10 ± 0.70	90.91 ± 0.06 / 87.85 ± 0.09	81.61 ± 0.25 / 81.71 ± 0.11	89.85 ± 0.20	72.08 ± 0.56	

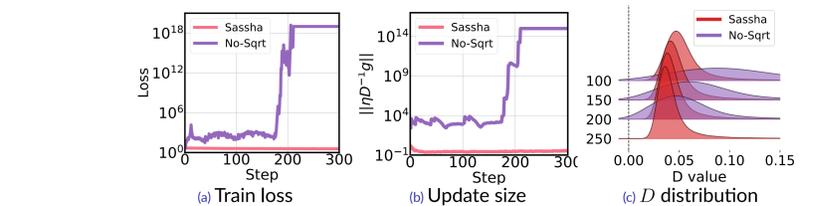
- (pretraining; left) SASSHA achieves the **lowest perplexity**.
- (finetuning; right) **better** than others, but competitive to Sophia-H.

Table 3. Robustness to label noise.

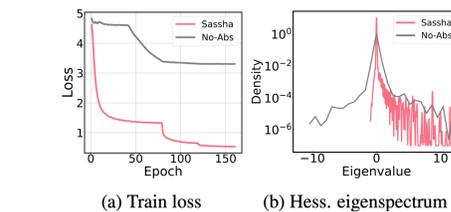
Noise level	CIFAR-10				CIFAR-100			
	0%	20%	40%	60%	0%	20%	40%	60%
SGD	92.69 ± 0.06	89.91 ± 0.87	87.26 ± 0.40	82.72 ± 1.59	69.32 ± 0.19	62.18 ± 0.06	55.78 ± 0.55	45.53 ± 0.78
SAM _{SGD}	93.89 ± 0.13	92.27 ± 0.14	90.11 ± 0.25	85.79 ± 0.30	71.99 ± 0.20	65.53 ± 0.11	61.20 ± 0.17	51.93 ± 0.47
AdaHessian	92.48 ± 0.15	90.11 ± 0.01	86.88 ± 0.04	83.25 ± 0.01	68.06 ± 0.22	63.06 ± 0.25	58.37 ± 0.13	46.02 ± 1.96
Sophia-H	91.99 ± 0.08	89.93 ± 0.01	87.30 ± 0.51	82.78 ± 1.43	67.76 ± 0.37	62.34 ± 0.47	56.54 ± 0.28	45.37 ± 0.27
Shampoo	90.23 ± 0.83	88.14 ± 0.29	85.15 ± 0.61	81.16 ± 0.30	64.08 ± 0.46	58.85 ± 0.66	53.82 ± 0.71	42.91 ± 0.99
SASSHA	94.09 ± 0.24	92.49 ± 0.11	90.29 ± 0.11	86.50 ± 0.08	72.14 ± 0.16	66.78 ± 0.47	61.97 ± 0.27	53.98 ± 0.57

- SASSHA shows **much robust performance** under label noise.

Ablation: stable Hessian approximation

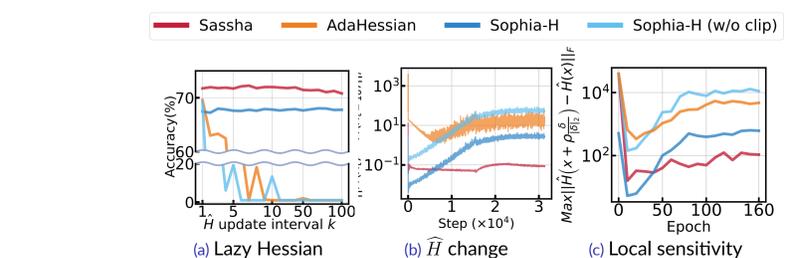


- (a) **without sqrt** training diverges when (b) the update size starts to spike.
- (c) individual entries in the diagonal Hessian; they gradually **shifts toward zero values**, as a result, D^{-1} becomes too large, creating overly large steps.



- (left) **without absolute** SASSHA does not converge well.
- (right) **negative values** in the Hessian distribution driving toward saddle points.

Robust to Lazy Hessian updates



- (a) SASSHA remains relatively effective with prolonged Hessian reusing
- (b) **The difference** between the current and previous Hessian is **small**.
- (c) Optimization is biased **toward regions with less curvature change**.

Cost analysis

Table 4. Average wall-clock time per epoch (s) and the theoretical cost of different methods, with the **lazy Hessian interval is set to 10**.

Method	Cost				CIFAR10	CIFAR100	ImageNet
	Descent	Sharpness	Hessian	Total	ResNet32	WRN28-10	ViT-small
AdamW	1 GC	0 GC	0 HVP	1 GC	5.03	59.29	976.56
SAM	1 GC	1 GC	0 HVP	2 GC	9.16	118.46	1302.08
AdaHessian	1 GC	0 GC	1 HVP	4 GC	33.75	296.63	2489.07
SASSHA	1 GC	1 GC	0.1 HVP	2.3 GC	12.00	142.06	1377.20

- SASSHA is **faster** than other approximate second-order methods.