



ADAM, do not?



$$\Sigma_t = \frac{1}{S} \Big[\frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{\theta}_t) \nabla f_i(\boldsymbol{\theta}_t)^T - \nabla \boldsymbol{F}(\boldsymbol{\theta}_t) \nabla f_i(\boldsymbol{\theta}_t)^T - \nabla \boldsymbol{F}(\boldsymbol{\theta}_t) \nabla f_i(\boldsymbol{\theta}_t)^T \Big]$$

$$\begin{cases} \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \boldsymbol{m}_t / (1 - \beta_1^t) / \left(\sqrt{\boldsymbol{v}_t / (1 - \beta_2^t)} + \epsilon \right), \\ \boldsymbol{m}_t = \beta_1 \boldsymbol{m}_{t-1} + (1 - \beta_1) \nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t), \\ \boldsymbol{v}_t = \beta_2 \boldsymbol{v}_{t-1} + (1 - \beta_2) [\nabla f_{\mathcal{S}_t}(\boldsymbol{\theta}_t)]^2, \end{cases} \xrightarrow{\text{ADAM}} \begin{cases} \mathsf{d}\boldsymbol{\theta}_t = \mathbf{d}_t \\ \mathsf{d}\boldsymbol{m}_t = \mathbf{d}_t \\ \mathsf{d}\boldsymbol{v}_t = \mathbf{d}_t \end{cases}$$

constants to correct the bias in m_t and v_t .

Towards Theoretically Understanding Why SGD Generalizes Better Than ADAM in Deep Learning

Pan Zhou*, Jiashi Feng[†], Chao Ma[‡], Caiming Xiong*, Steven HOI*, Weinan E[‡] *Salesforce Research, [†] National University of Singapore, [‡] Princeton University

> ຮ 1.4

of
$$\mathcal{W}$$
 as
 $\{ \mathbf{y} \in \mathbb{R}^d \mid \mathbf{Q}_{\theta^*}^{-1} \Sigma_{\theta^*} \mathbf{y} \in \Omega^{-\varepsilon^{\gamma}} \}.$
 $m(\mathcal{W}^c) \longrightarrow \text{small } m(\mathcal{W}) (m(\mathcal{W} \cup \mathcal{W}^c) = \text{constant})$







with a basin height $h(\theta^*)$ and Hessian matrix $H(\theta^*)$ at θ^* .

Compariso
basin appro
and it satieti
where \mathcal{W}_{SG}
SGD and $ar{\Sigma}_{ heta}$

Conclusion: SGD could better escape from sharp minima and converge to flatter minima, since

- could not.



Result 2: Better Sharp Minima Escaping Ability of SGD over ADAM

 \blacktriangleright with same learning rate ε in ADAM and SGD, the smaller tail index α , the smaller the escaping time Γ .

 \blacktriangleright for some iterations, SGD has smaller α , as exponential gradient average in ADAM smooths noise.

(a) MNIST (over-parameterized fully connected networks)

(b) CIFAR10 (over-parameterized fully connected networks)

Factor 2 $m(\mathcal{W})$ that positively depends on the volume of escaping set \mathcal{W} .

Approximating Ω as a quadratic basin with center θ^* , i.e.

$$\Omega = \left\{ \boldsymbol{y} \mid \boldsymbol{F}(\boldsymbol{\theta}^*) + \frac{1}{2} \boldsymbol{y}^T \boldsymbol{H}(\boldsymbol{\theta}^*) \boldsymbol{y} \leq h(\boldsymbol{\theta}^*) \right\}$$

on of Escaping Sets of SGD and ADAM. Under the quadratic oximation, the escaping set \mathcal{W} of ADAM is

$$\mathcal{N}_{\mathsf{ADAM}} pprox \left\{ oldsymbol{y} \in \mathbb{R}^d \mid oldsymbol{y}^{\mathsf{T}}oldsymbol{H}(oldsymbol{ heta}^*)oldsymbol{y} \geq S^2 h_f^*
ight\}.$$

les

 $m({m {\cal W}}_{ extsf{ADAM}}) < m({m {\cal W}}_{ extsf{SGD}})$ $_{ ext{GD}} = \left\{ m{y} \in \mathbb{R}^d \; ig| m{y}^T ar{\Sigma}_{ heta^*} m{H}(m{ heta}^*) ar{\Sigma}_{ heta^*} m{y} \geq S^2 h_f^*
ight\}$ is escaping set of $\Sigma_{\theta^*} = \frac{1}{n} \sum_{i=1}^n \nabla f_i(\theta^*) \nabla f_i(\theta^*)^T.$

from Factors 1 and 2, SGD has smaller escaping time and is much more unstable. For the same basin, ADAM could stuck in one basin, but SGD

both SGD and ADAM prefers to converge to flat minima