

Широкие и узкие ОПТИМУМЫ

Методы обучения широких оптимумов (SWA и SAM)

План

- Разобраться что такое широкие и узкие оптимумы и почему мы хотим получать именно широкие
- Рассмотреть модели SWA и SAM которые позволяют нам обучаться под широкие оптимумы

Что такое широкие и узкие оптимумы?



Narrow optimum

Eigenvalues of Hessian very large

Loss increases rapidly as parameters change



Broad optimum

Eigenvalues of Hessian close to zero

Loss increases slowly as parameters change

Чем хороши широкие оптимумы?

- Широкие оптимумы повышают обобщающую способность нашей модели
- Узкие оптимумы наоборот крайне неустойчивы и поэтому обладают более низкой обобщающей способностью

SWA (Stochastic Weight Averaging): Основные идеи

- Вспоминаем FGE (Fast Geometric Ensembling)
- Результаты были хорошие, но тяжёлый инференс
- А что будет, если усреднять не предсказания моделей, а их веса?
- А будет чудо: получим хорошие предсказания и всего одну предсказывающую модель на выходе

SWA

Learning Rate

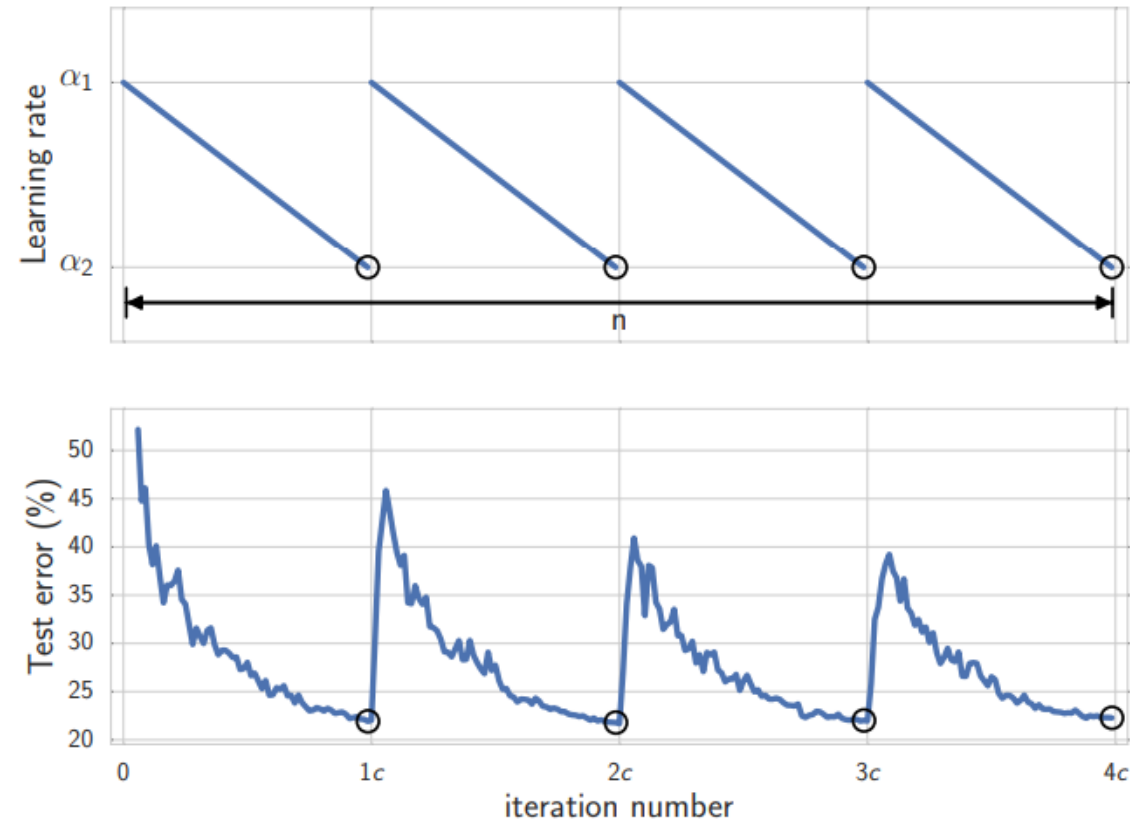
Возьмём LearningRate как у FGE:

$$\alpha(i) = (1 - t(i))\alpha_1 + t(i)\alpha_2,$$
$$t(i) = \frac{1}{c} (\text{mod}(i - 1, c) + 1) .$$

То есть наш LR будет циклично снижаться от α_1 до α_2

SWA

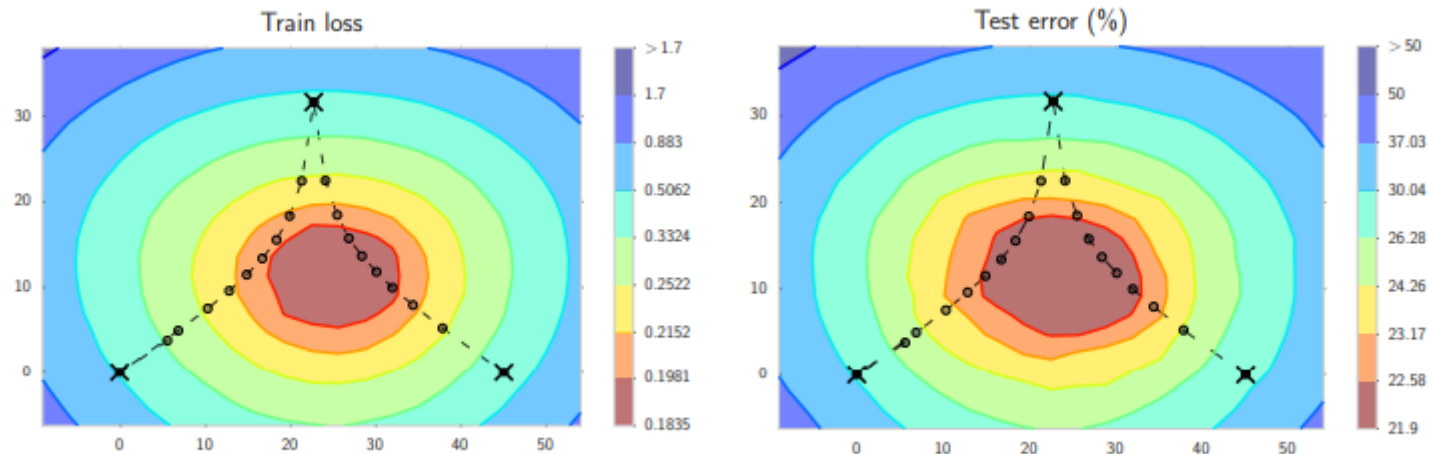
Learning Rate



SWA

Learning Rate (мотивация)

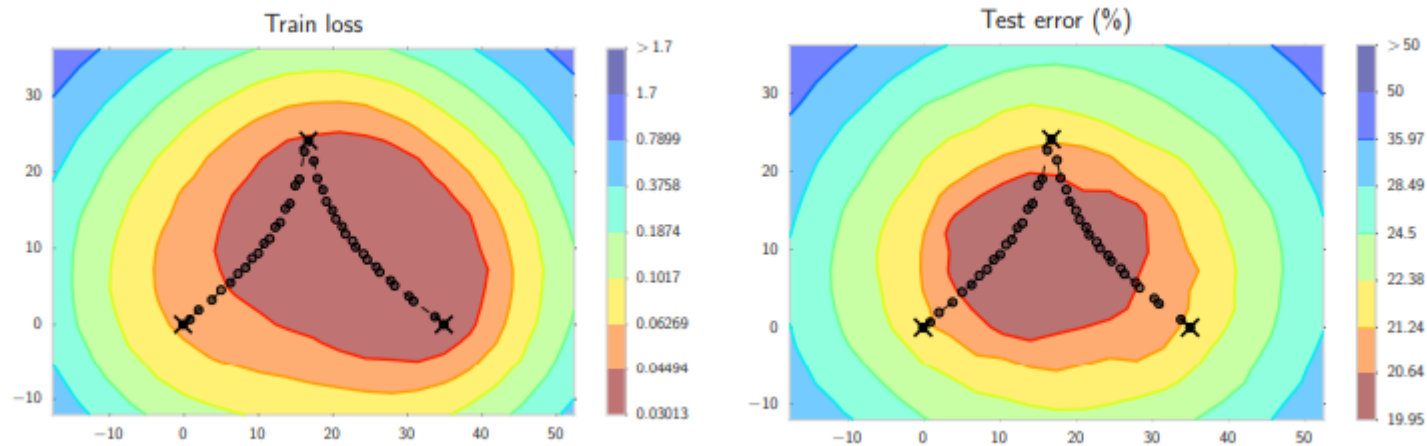
- Если взять константный LR, то мы будем быстро приближаться к минимуму, однако потом будем “перескакивать”



SWA

Learning Rate (мотивация)

- С другой стороны, если взять наш циклический LR, то мы будем дольше идти к минимуму, но зато за несколько циклов найдём сразу несколько локальных минимумов и вряд ли перескочим искомый



SWA

- Давайте попытаемся совместить плюсы каждого подхода
- Для начала будем использовать constant LR
- Затем с какого-то момента возьмём наш cyclic LR и дообучим модель с ним
- Таким образом получим быстро и хорошо обучаемую модель

SWA

Алгоритм

Algorithm 1 Stochastic Weight Averaging

Require:

weights \hat{w} , LR bounds α_1, α_2 ,
cycle length c (for constant learning rate $c = 1$), number of iterations n

Ensure: w_{SWA}

$w \leftarrow \hat{w}$ {Initialize weights with \hat{w} }

$w_{\text{SWA}} \leftarrow w$

for $i \leftarrow 1, 2, \dots, n$ **do**

$\alpha \leftarrow \alpha(i)$ {Calculate LR for the iteration}

$w \leftarrow w - \alpha \nabla \mathcal{L}_i(w)$ {Stochastic gradient update}

if $\text{mod}(i, c) = 0$ **then**

$n_{\text{models}} \leftarrow i/c$ {Number of models}

$w_{\text{SWA}} \leftarrow \frac{w_{\text{SWA}} \cdot n_{\text{models}} + w}{n_{\text{models}} + 1}$ {Update average}

end if

end for

{Compute BatchNorm statistics for w_{SWA} weights}

SWA и FGE

Покажем, что SWA и FGE выдают близкие результаты

$$w_{SWA} = \frac{1}{n} \sum w_i$$

Пускай $\Delta_i = w_i - w_{SWA}$ тогда получаем, что $\sum \Delta_i = 0$

Линеаризуем f в точке w_{SWA} $f(w_j) = f(w_{SWA}) + \langle \nabla f(w_{SWA}), \Delta_j \rangle + O(\|\Delta_j\|^2)$

$$f_{FGE} = \frac{1}{n} \sum f(w_i)$$

$$f_{FGE} - f(w_{SWA}) = \frac{1}{n} (\langle \nabla f(w_{SWA}), \sum \Delta_j \rangle) + O(\|\Delta_j\|^2) = O(\|\Delta_j\|^2)$$

SWA

Сравнение результатов

$$w_{\text{SWA}}(t, d) = w_{\text{SWA}} + t \cdot d,$$

$$w_{\text{SGD}}(t, d) = w_{\text{SGD}} + t \cdot d,$$

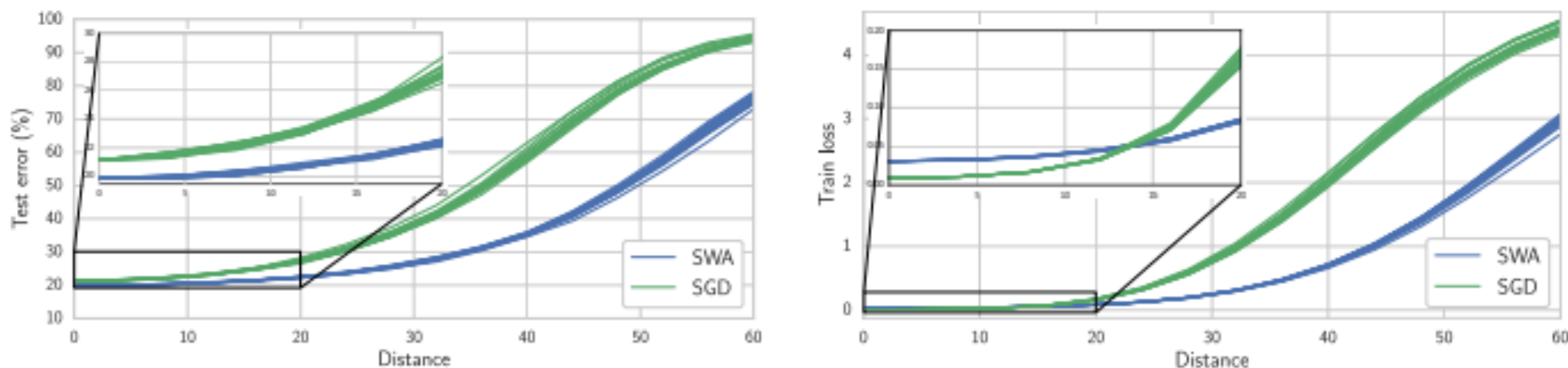


Figure 4: **(Left)** Test error and **(Right)** L_2 -regularized cross-entropy train loss as a function of a point on a random ray starting at SWA (blue) and SGD (green) solutions for Preactivation ResNet-164 on CIFAR-100. Each line corresponds to a different random ray.

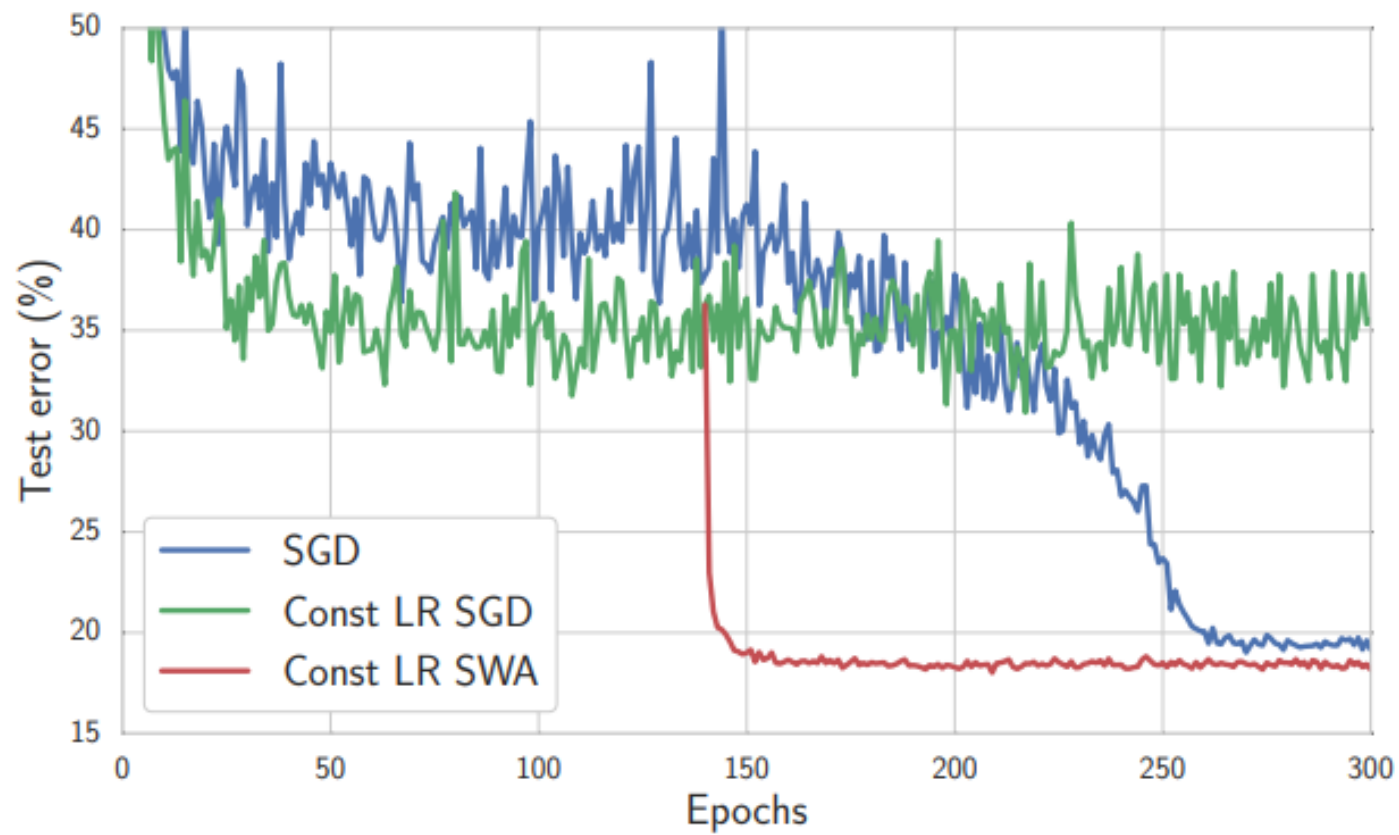
SWA

Сравнение результатов

DNN (Budget)	SGD	FGE (1 Budget)	SWA		
			1 Budget	1.25 Budgets	1.5 Budgets
CIFAR-100					
VGG-16 (200)	72.55 ± 0.10	74.26	73.91 ± 0.12	74.17 ± 0.15	74.27 ± 0.25
ResNet-164 (150)	78.49 ± 0.36	79.84	79.77 ± 0.17	80.18 ± 0.23	80.35 ± 0.16
WRN-28-10 (200)	80.82 ± 0.23	82.27	81.46 ± 0.23	81.91 ± 0.27	82.15 ± 0.27
PyramidNet-272 (300)	83.41 ± 0.21	–	–	83.93 ± 0.18	84.16 ± 0.15
CIFAR-10					
VGG-16 (200)	93.25 ± 0.16	93.52	93.59 ± 0.16	93.70 ± 0.22	93.64 ± 0.18
ResNet-164 (150)	95.28 ± 0.10	95.45	95.56 ± 0.11	95.77 ± 0.04	95.83 ± 0.03
WRN-28-10 (200)	96.18 ± 0.11	96.36	96.45 ± 0.11	96.64 ± 0.08	96.79 ± 0.05
ShakeShake-2x64d (1800)	96.93 ± 0.10	–	–	97.16 ± 0.10	97.12 ± 0.06

SWA

Сравнение результатов



SAM(Sharpness-Aware Minimization)

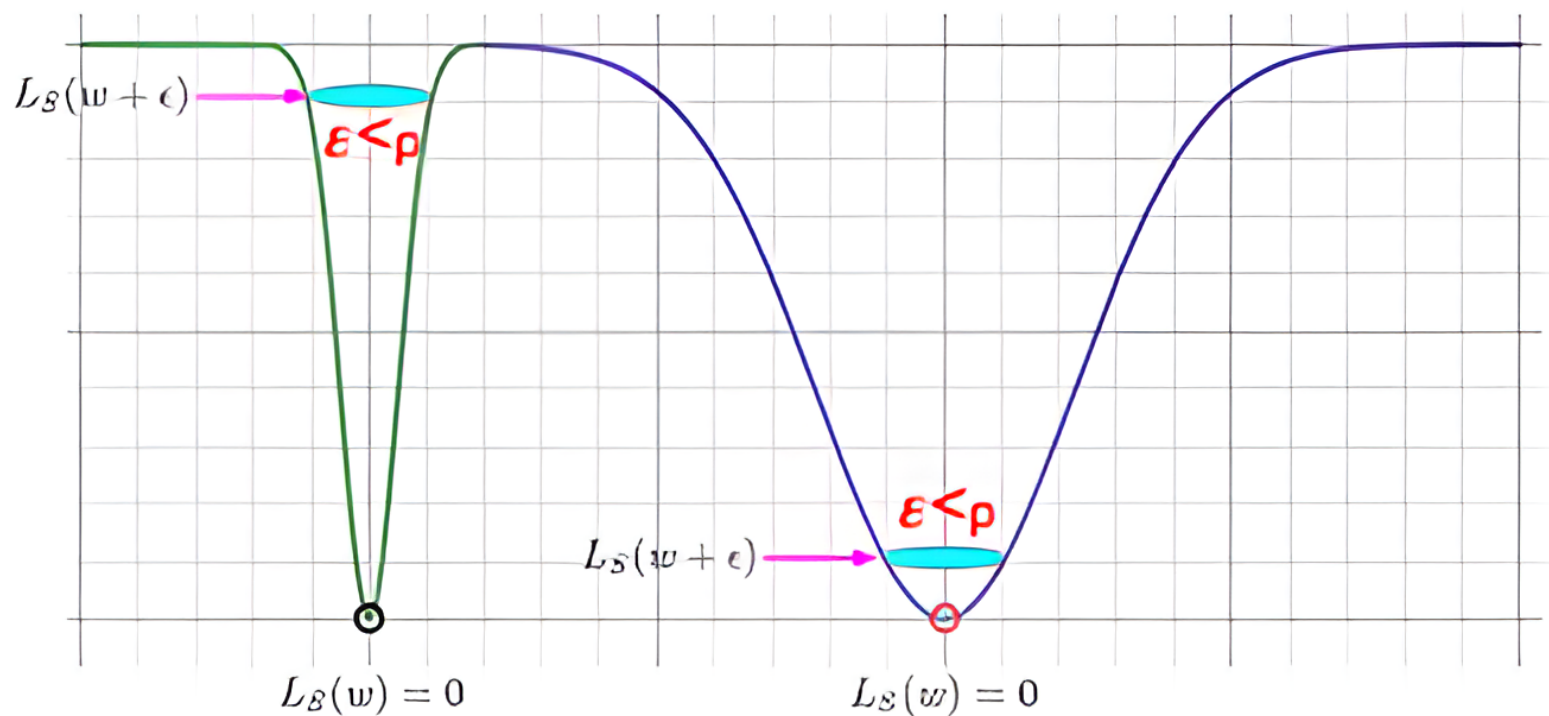
Основные идеи

- Сейчас мы минимизируем функцию потерь в точке, давайте попробуем минимизировать функцию потерь в некоторой окрестности
- Будем искать такую область максимум в которой минимален, среди всех областей
- Таким образом получим целую область на которой значение функции потерь маленькое, а, значит, найдём широкий оптимум

SAM

Функция потерь

$$\min_{\mathbf{w}} L_S^{SAM}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 \quad \text{where} \quad L_S^{SAM}(\mathbf{w}) \triangleq \max_{\|\epsilon\|_p \leq \rho} L_S(\mathbf{w} + \epsilon),$$



SAM

Формулы, формулы, формулы...

$$\hat{\epsilon}(\boldsymbol{w}) = \rho \operatorname{sign}(\nabla_{\boldsymbol{w}} L_{\mathcal{S}}(\boldsymbol{w})) |\nabla_{\boldsymbol{w}} L_{\mathcal{S}}(\boldsymbol{w})|^{q-1} / \left(\|\nabla_{\boldsymbol{w}} L_{\mathcal{S}}(\boldsymbol{w})\|_q^q \right)^{1/p} \quad (2)$$

$$\nabla_{\boldsymbol{w}} L_{\mathcal{S}}^{SAM}(\boldsymbol{w}) \approx \nabla_{\boldsymbol{w}} L_{\mathcal{S}}(\boldsymbol{w})|_{\boldsymbol{w} + \hat{\epsilon}(\boldsymbol{w})}. \quad (3)$$

SAM

Алгоритм

Input: Training set $\mathcal{S} \triangleq \cup_{i=1}^n \{(\mathbf{x}_i, \mathbf{y}_i)\}$, Loss function $l : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, Batch size b , Step size $\eta > 0$, Neighborhood size $\rho > 0$.

Output: Model trained with SAM

Initialize weights $\mathbf{w}_0, t = 0$;

while *not converged* **do**

 Sample batch $\mathcal{B} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_b, \mathbf{y}_b)\}$;

 Compute gradient $\nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})$ of the batch's training loss;

 Compute $\hat{\epsilon}(\mathbf{w})$ per equation 2;

 Compute gradient approximation for the SAM objective

 (equation 3): $\mathbf{g} = \nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})}$;

 Update weights: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}$;

$t = t + 1$;

end

return \mathbf{w}_t

Algorithm 1: SAM algorithm

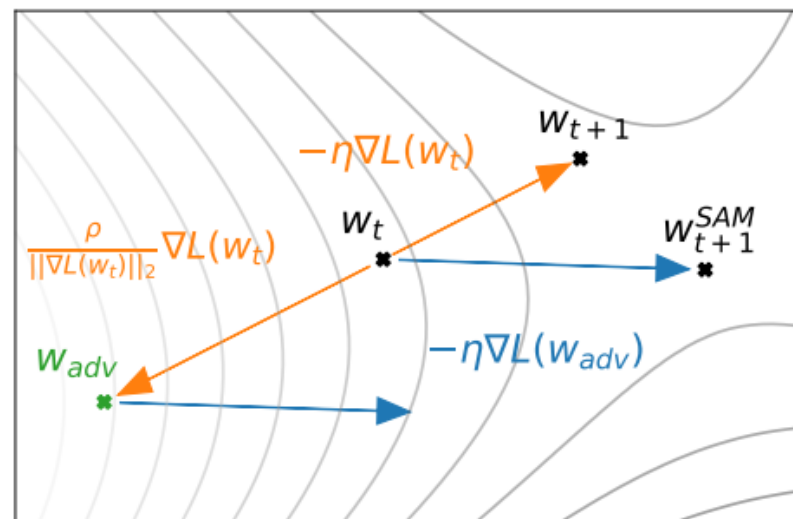


Figure 2: Schematic of the SAM parameter update.

SAM

Сравнение результатов

Model	Augmentation	CIFAR-10		CIFAR-100	
		SAM	SGD	SAM	SGD
WRN-28-10 (200 epochs)	Basic	2.7 ± 0.1	3.5 ± 0.1	16.5 ± 0.2	18.8 ± 0.2
WRN-28-10 (200 epochs)	Cutout	2.3 ± 0.1	2.6 ± 0.1	14.9 ± 0.2	16.9 ± 0.1
WRN-28-10 (200 epochs)	AA	2.1 $\pm <0.1$	2.3 ± 0.1	13.6 ± 0.2	15.8 ± 0.2
WRN-28-10 (1800 epochs)	Basic	2.4 ± 0.1	3.5 ± 0.1	16.3 ± 0.2	19.1 ± 0.1
WRN-28-10 (1800 epochs)	Cutout	2.1 ± 0.1	2.7 ± 0.1	14.0 ± 0.1	17.4 ± 0.1
WRN-28-10 (1800 epochs)	AA	1.6 ± 0.1	2.2 $\pm <0.1$	12.8 ± 0.2	16.1 ± 0.2
Shake-Shake (26 2x96d)	Basic	2.3 $\pm <0.1$	2.7 ± 0.1	15.1 ± 0.1	17.0 ± 0.1
Shake-Shake (26 2x96d)	Cutout	2.0 $\pm <0.1$	2.3 ± 0.1	14.2 ± 0.2	15.7 ± 0.2
Shake-Shake (26 2x96d)	AA	1.6 $\pm <0.1$	1.9 ± 0.1	12.8 ± 0.1	14.1 ± 0.2
PyramidNet	Basic	2.7 ± 0.1	4.0 ± 0.1	14.6 ± 0.4	19.7 ± 0.3
PyramidNet	Cutout	1.9 ± 0.1	2.5 ± 0.1	12.6 ± 0.2	16.4 ± 0.1
PyramidNet	AA	1.6 ± 0.1	1.9 ± 0.1	11.6 ± 0.1	14.6 ± 0.1
PyramidNet+ShakeDrop	Basic	2.1 ± 0.1	2.5 ± 0.1	13.3 ± 0.2	14.5 ± 0.1
PyramidNet+ShakeDrop	Cutout	1.6 $\pm <0.1$	1.9 ± 0.1	11.3 ± 0.1	11.8 ± 0.2
PyramidNet+ShakeDrop	AA	1.4 $\pm <0.1$	1.6 $\pm <0.1$	10.3 ± 0.1	10.6 ± 0.1

SAM

Сравнение результатов

Model	Epoch	SAM		Standard Training (No SAM)	
		Top-1	Top-5	Top-1	Top-5
ResNet-50	100	22.5 ± 0.1	6.28 ± 0.08	22.9 ± 0.1	6.62 ± 0.11
	200	21.4 ± 0.1	5.82 ± 0.03	22.3 ± 0.1	6.37 ± 0.04
	400	20.9 ± 0.1	5.51 ± 0.03	22.3 ± 0.1	6.40 ± 0.06
ResNet-101	100	20.2 ± 0.1	5.12 ± 0.03	21.2 ± 0.1	5.66 ± 0.05
	200	19.4 ± 0.1	4.76 ± 0.03	20.9 ± 0.1	5.66 ± 0.04
	400	19.0 $\pm <0.01$	4.65 ± 0.05	22.3 ± 0.1	6.41 ± 0.06
ResNet-152	100	19.2 $\pm <0.01$	4.69 ± 0.04	20.4 $\pm <0.0$	5.39 ± 0.06
	200	18.5 ± 0.1	4.37 ± 0.03	20.3 ± 0.2	5.39 ± 0.07
	400	18.4 $\pm <0.01$	4.35 ± 0.04	20.9 $\pm <0.0$	5.84 ± 0.07

SAM

Сравнение результатов

