

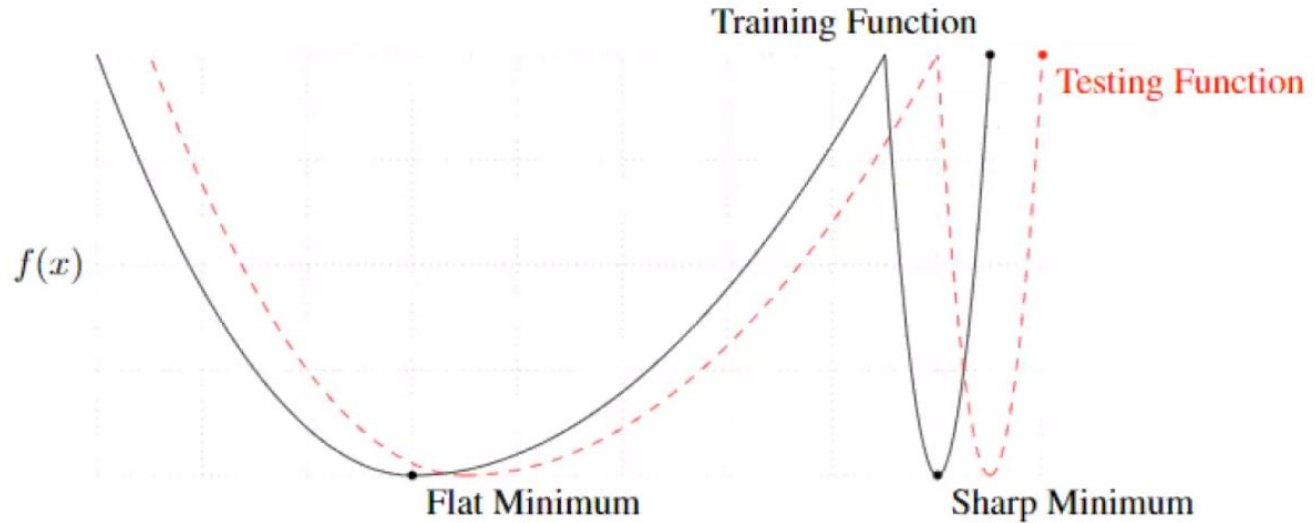
Широкие и узкие ОПТИМУМЫ

Методы получения более широких оптимумов (SAM и SWA)

План

1. Разобраться в широких и узких оптимумах, определить ширину оптимума, понять почему широкие оптимумы лучше
2. Разобраться в методах SAM и SWA

Широкие и узкие оптимумы



по оси Y отмечаются значения функции потерь, а по оси X значения параметров

Проблема: как определить ширину оптимума в многомерном случае

Определение ширины оптимума

1. Объем области, в которой значение функции потерь не сильно отличаются друг от друга (Hochreiter & Schmidhuber, 1997)
2. Максимальное значение функции потери в окрестности найденного оптимума (Кескар и др., 2017)
3. Собственное значение Гессиана в точках минимума (Яо и др., 2018)

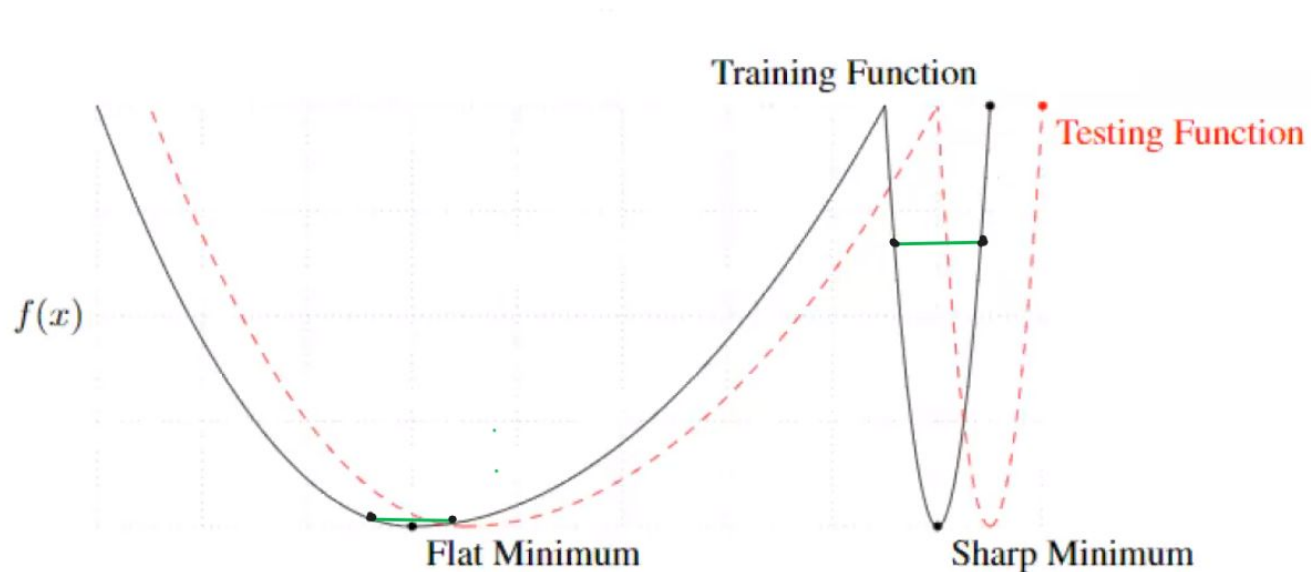
Очевидное наблюдение

Чем оптимум шире — тем выше обобщающая способность и адаптивность модели

Глобальная идея – искать широкие минимумы

SAM (Sharpness-Aware Minimization)

1. Метод для эффективного повышения обобщенности модели
2. Основная идея – минимизировать функцию потерь в некоторой окрестности



Вводим новую функцию потерь с гиперпараметром $p \in [1, \infty]$

$$L_S^{SAM}(\mathbf{w}) \triangleq \max_{\|\epsilon\|_p \leq \rho} L_S(\mathbf{w} + \epsilon),$$

Задача, лежащая в основе метода SAM:

$$\min_{\mathbf{w}} L_S^{SAM}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

Эффективный способ посчитать градиент

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

где

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

$$\epsilon^*(\mathbf{w}) \triangleq \arg \max_{\|\epsilon\|_p \leq \rho} L_{\mathcal{S}}(\mathbf{w} + \epsilon) \approx \arg \max_{\|\epsilon\|_p \leq \rho} L_{\mathcal{S}}(\mathbf{w}) + \epsilon^T \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w}) = \arg \max_{\|\epsilon\|_p \leq \rho} \epsilon^T \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w}).$$

Реализация метода SAM

Input: Training set $\mathcal{S} \triangleq \cup_{i=1}^n \{(x_i, y_i)\}$, 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 $w_0, t = 0$;

while not converged **do**

 Sample batch $\mathcal{B} = \{(x_1, y_1), \dots (x_b, y_b)\}$;

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

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

 Compute gradient approximation for the SAM objective (equation 3): $g = \nabla_w L_{\mathcal{B}}(w)|_{w+\hat{e}(w)}$;

 Update weights: $w_{t+1} = w_t - \eta g$;

$t = t + 1$;

end

return w_t

Algorithm 1: SAM algorithm

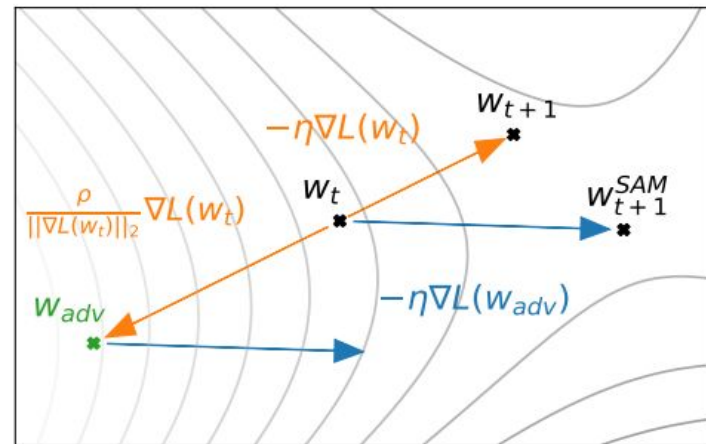
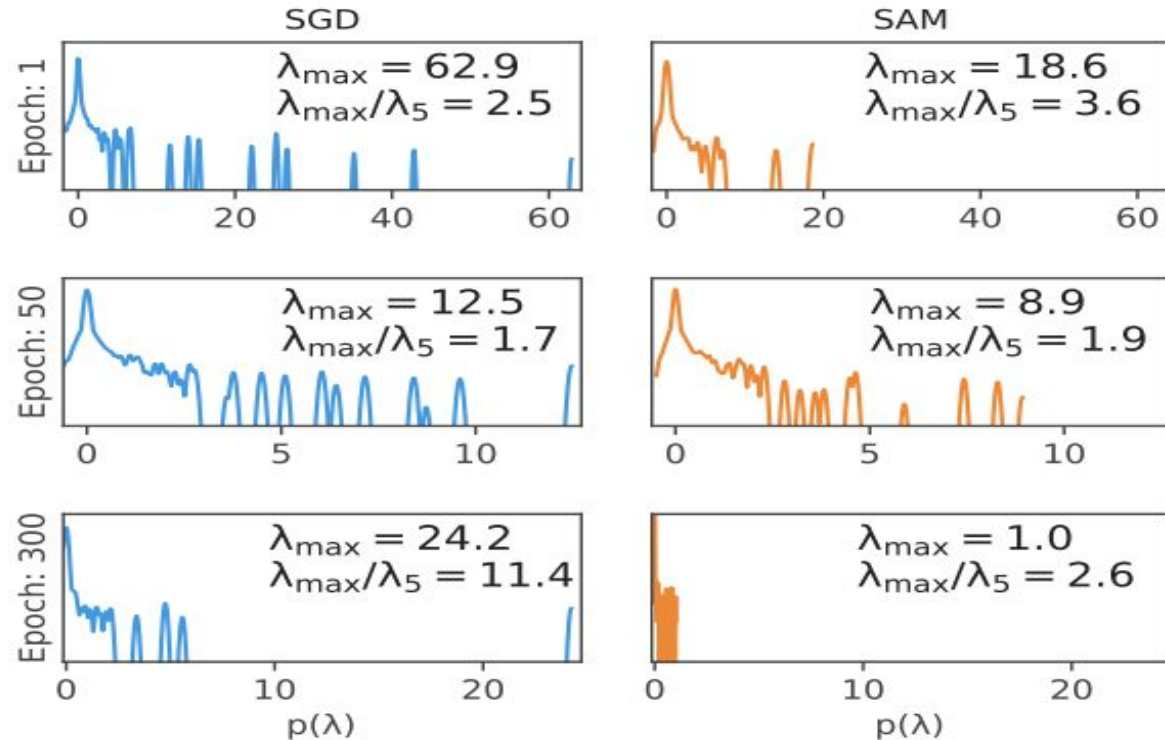


Figure 2: Schematic of the SAM parameter update.

Сравнение собственных значений Гессиан



		CIFAR-10		CIFAR-100	
Model	Augmentation	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

ImageNet

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

Table 2: Test error rates for ResNets trained on ImageNet, with and without SAM.

SWA (Stochastic Weight Averaging)

1. Вдохновимся идеей FGE и будем усреднять веса вдоль кривых между локальными минимумами

Преимущества:

1. Всего одна модель
2. Решает проблему большого инференса

Loss landscape

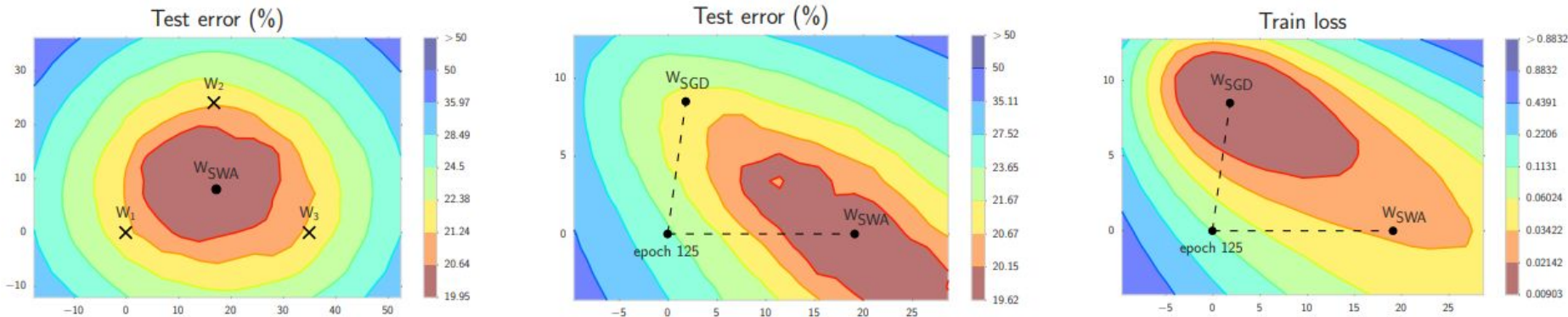
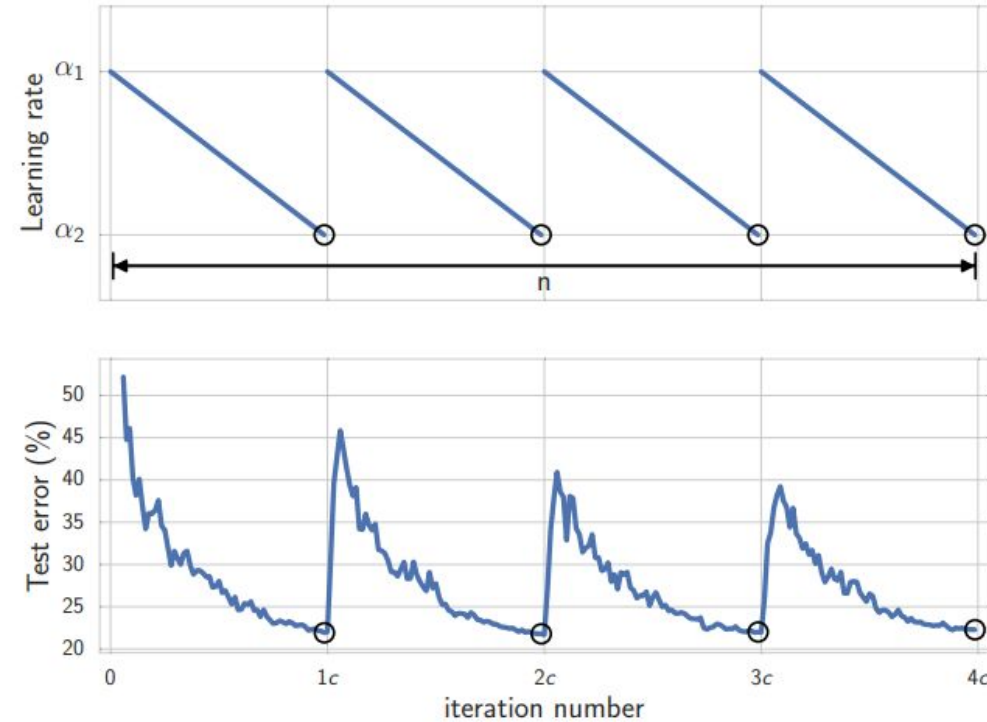


Figure 1: Illustrations of SWA and SGD with a Preactivation ResNet-164 on CIFAR-100¹. **Left:** test error surface for three FGE samples and the corresponding SWA solution (averaging in weight space). **Middle and Right:** test error and train loss surfaces showing the weights proposed by SGD (at convergence) and SWA, starting from the same initialization of SGD after 125 training epochs.

Learning rate



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

Зависимость ошибки от итерации модели ResNet164 на датасете CIFAR-100

Реализация метода SWA

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}

Связь с широким оптимумом

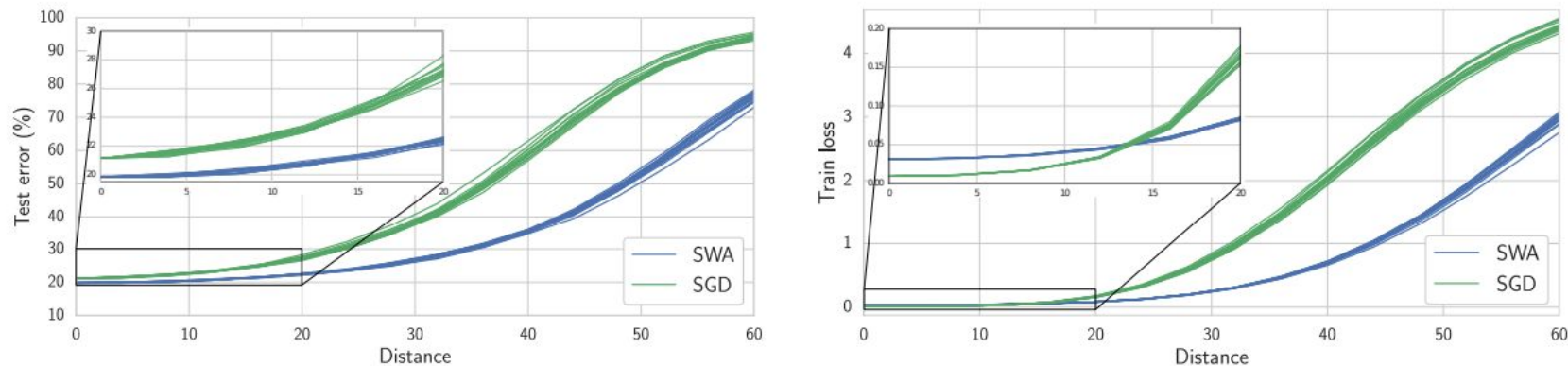


Figure 4: **(Left)** Test error and **(Right)** L_2 -regularized cross-entropy train loss as a function of a point on a random ray starting at SWA (blue) and SGD (green) solutions for Preactivation ResNet-164 on CIFAR-100. Each line corresponds to a different random ray.

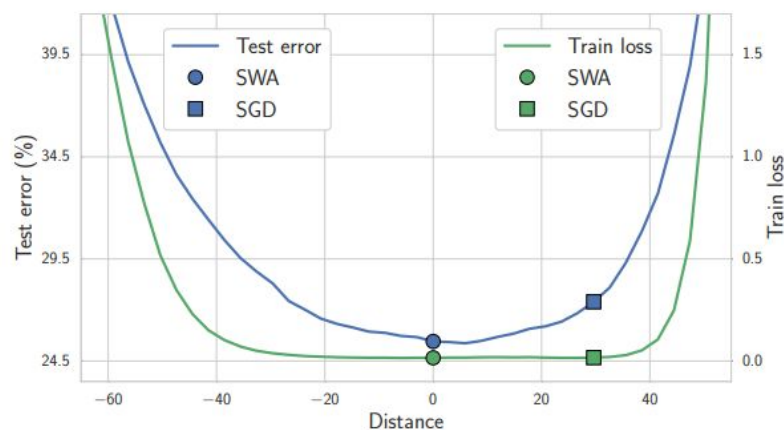
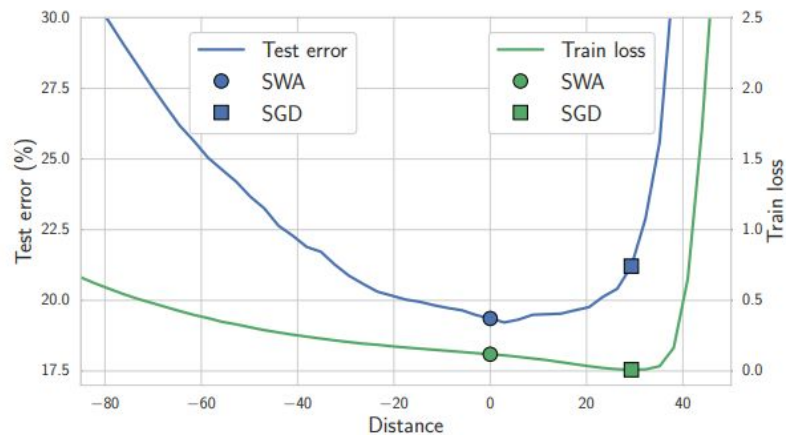


Figure 5: L_2 -regularized cross-entropy train loss and test error as a function of a point on the line connecting SWA and SGD solutions on CIFAR-100. **Left:** Preactivation ResNet-164. **Right:** VGG-16.

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

Table 1: Accuracies (%) of SWA, SGD and FGE methods on CIFAR-100 and CIFAR-10 datasets for different training budgets. Accuracies for the FGE ensemble are from Garipov et al. [2018].

Связь SWA и FGE

$$w_{SWA} = \frac{1}{n} \sum w_i, \quad w_{FGE} = \frac{1}{n} \sum f(w_i)$$

Пусть $\Delta_i = w_i - w_{SWA}$, тогда $\sum \Delta_i = 0$

Линейизируем f в точке w_{SWA}

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

тогда

$$\bar{f} - f(w_{SWA}) = \frac{1}{n} \sum_{i=1}^n (\langle \nabla f(w_{SWA}), \Delta_i \rangle + O(\|\Delta_i\|^2)) = \langle \nabla f(w_{SWA}), \frac{1}{n} \sum_{i=1}^n \Delta_i \rangle + O(\Delta^2) = O(\Delta^2),$$

Преимущества триангулярного lr

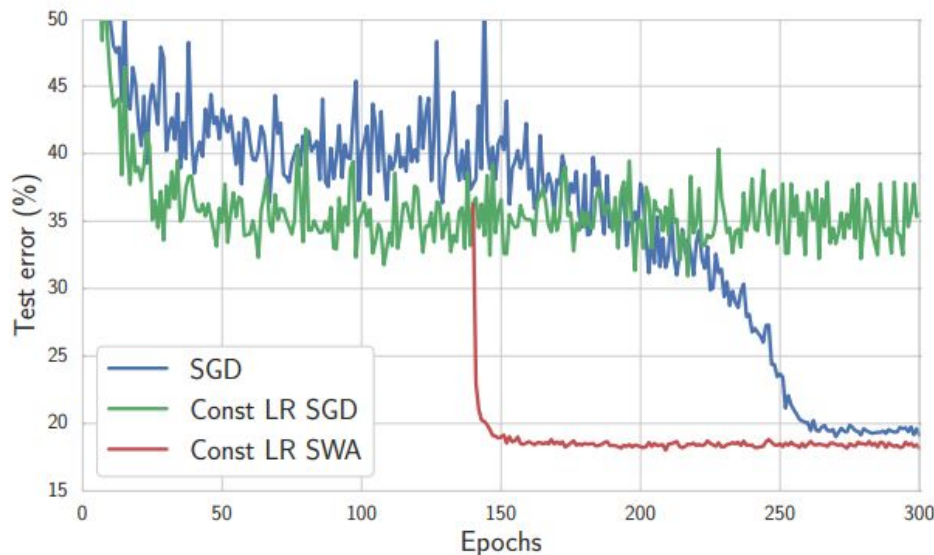


Figure 7: Test error as a function of training epoch for constant (green) and decaying (blue) learning rate schedules for a Wide ResNet-28-10 on CIFAR-100. In red we average the points along the trajectory of SGD with constant learning rate starting at epoch 140.

Используемые ресурсы

1. [Статья про SAM](#)
2. [Статья про SWA](#)