

Shapeless-Aware Minimization

for Efficiently Improving Generalization

Сафонов Иван (докладчик)
Рахматуллин Рамазан (рецензент)
Терехова Юлия (практик-исследователь)
Сухоросов Алексей (хакер)
БПМИ 182

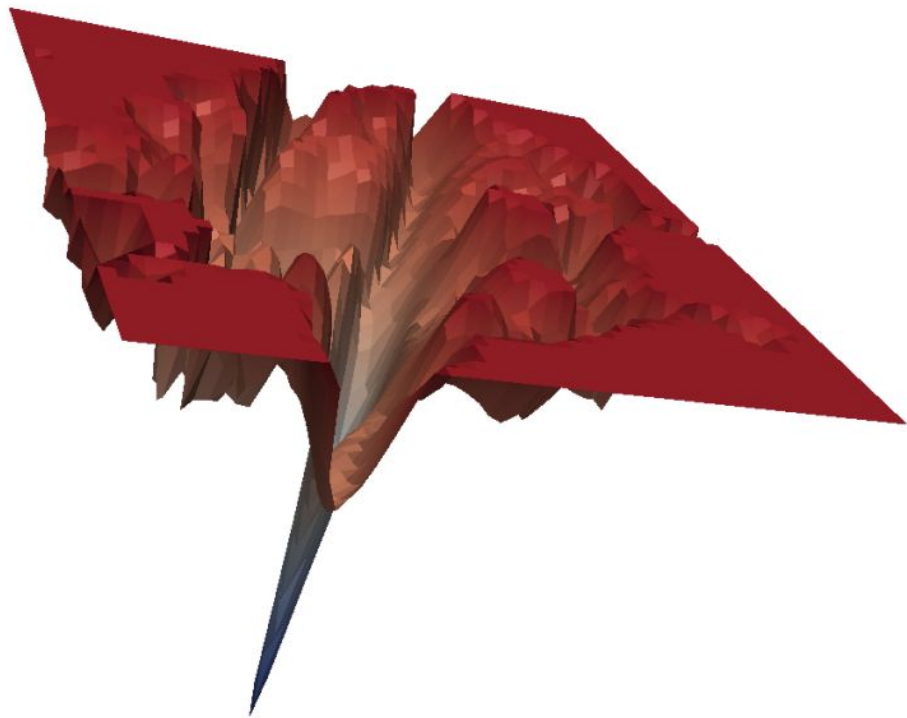
06.10.2021

План

- Мотивация метода
- Описание метода и алгоритма
- Результаты
- Исследования вокруг алгоритма

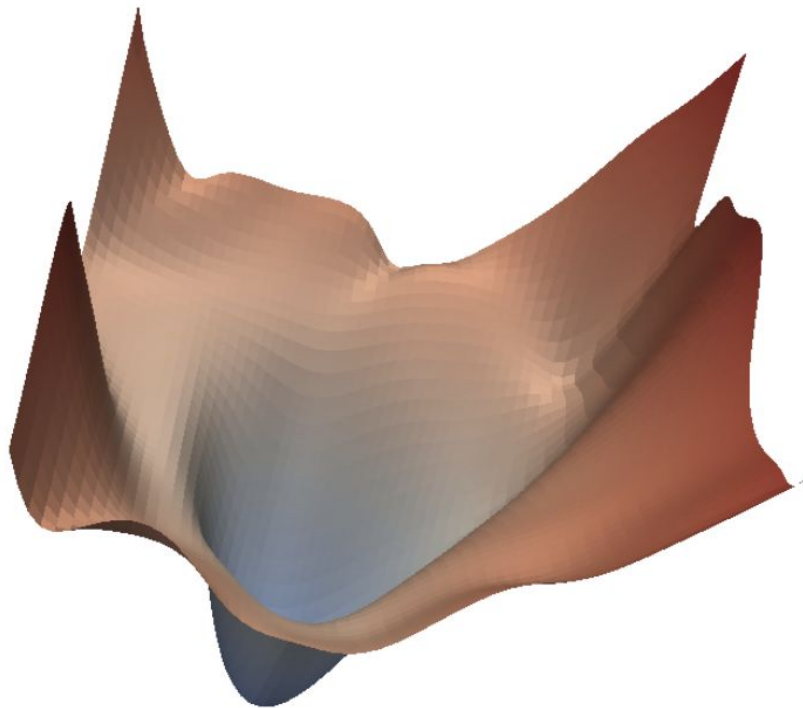
Мотивация

- loss функция крайне сложна и имеет множество точек оптимума
- В разных точках модель получает разные generalization свойства
- Выбор оптимизатора/регуляризации является очень важным шагом



Идея

- Давайте кроме самого loss минимизировать его заостренность в окрестности точки
- Тогда (как утверждается), наш алгоритм не сможет упасть в какой-то локальный минимум, приводящий к переобучению



Определения

$$\mathcal{S} \triangleq \cup_{i=1}^n \{(\mathbf{x}_i, \mathbf{y}_i)\}$$

- выборка

$$l : \mathcal{W} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$$

- loss function

$$L_{\mathcal{S}}(\mathbf{w}) \triangleq \frac{1}{n} \sum_{i=1}^n l(\mathbf{w}, \mathbf{x}_i, \mathbf{y}_i)$$

- training set loss

$$L_{\mathcal{D}}(\mathbf{w}) \triangleq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim D} [l(\mathbf{w}, \mathbf{x}, \mathbf{y})]$$

- population loss

Основная теорема

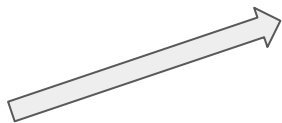
Для любого $\rho > 0$ с большой вероятностью при генерации случайной выборки:

$$L_{\mathcal{D}}(\mathbf{w}) \leq \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} L_{\mathcal{S}}(\mathbf{w} + \boldsymbol{\epsilon}) + h(\|\mathbf{w}\|_2^2 / \rho^2)$$

$h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ - некоторая строго возрастающая функция

SAM

$$\underbrace{\left[\max_{\|\epsilon\|_2 \leq \rho} L_S(\mathbf{w} + \epsilon) - L_S(\mathbf{w}) \right] + L_S(\mathbf{w}) + h(\|\mathbf{w}\|_2^2 / \rho^2)}_{\text{метрика заостренности}}$$



метрика заостренности

SAM optimization problem:

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

Решение задачи оптимизации

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

Решение:

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

для $p = 2$: градиент, нормированный на длину ρ

Решение задачи оптимизации

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

$$\begin{aligned} \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w} + \hat{\epsilon}(\mathbf{w})) &= \frac{d(\mathbf{w} + \hat{\epsilon}(\mathbf{w}))}{d\mathbf{w}} \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})} = \\ &= \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})} + \frac{d\hat{\epsilon}(\mathbf{w})}{d\mathbf{w}} \nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})|_{\mathbf{w} + \hat{\epsilon}(\mathbf{w})} \end{aligned}$$

В результате, для простоты:

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

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

 Compute gradient approximation for the SAM objective

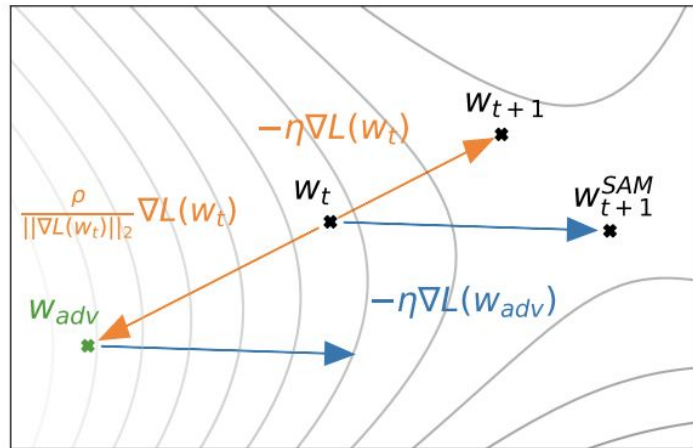
 (equation 3): $\mathbf{g} = \nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})|_{\mathbf{w} + \hat{\mathbf{e}}(\mathbf{w})}$;

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

$t = t + 1$;

end

return \mathbf{w}_t



Результаты, CIFAR-10, 100

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

Детали экспериментов

- Гиперпараметр ρ искали с помощью grid-search
- В качестве baseline запускали метод, давая ему в 2 раза больше эпох для обучения. При этом берется лучший результат среди такого же количества эпох и в 2 раза большего
- 5 независимых запусков для каждого из случаев

Результаты, 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

Результаты, finetuning, pretrained on ImageNet

Dataset	EffNet-b7 + SAM	EffNet-b7	Prev. SOTA (ImageNet only)	EffNet-L2 + SAM	EffNet-L2	Prev. SOTA
FGVC_Aircraft	6.80 \pm 0.06	8.15 \pm 0.08	5.3 (TBMSL-Net)	4.82 \pm 0.08	5.80 \pm 0.1	5.3 (TBMSL-Net)
Flowers	0.63 \pm 0.02	1.16 \pm 0.05	0.7 (BiT-M)	0.35 \pm 0.01	0.40 \pm 0.02	0.37 (EffNet)
Oxford_IIT_Pets	3.97 \pm 0.04	4.24 \pm 0.09	4.1 (Gpipe)	2.90 \pm 0.04	3.08 \pm 0.04	4.1 (Gpipe)
Stanford_Cars	5.18 \pm 0.02	5.94 \pm 0.06	5.0 (TBMSL-Net)	4.04 \pm 0.03	4.93 \pm 0.04	3.8 (DAT)
CIFAR-10	0.88 \pm 0.02	0.95 \pm 0.03	1 (Gpipe)	0.30 \pm 0.01	0.34 \pm 0.02	0.63 (BiT-L)
CIFAR-100	7.44 \pm 0.06	7.68 \pm 0.06	7.83 (BiT-M)	3.92 \pm 0.06	4.07 \pm 0.08	6.49 (BiT-L)
Birdsnap	13.64 \pm 0.15	14.30 \pm 0.18	15.7 (EffNet)	9.93 \pm 0.15	10.31 \pm 0.15	14.5 (DAT)
Food101	7.02 \pm 0.02	7.17 \pm 0.03	7.0 (Gpipe)	3.82 \pm 0.01	3.97 \pm 0.03	4.7 (DAT)
ImageNet	15.14 \pm 0.03	15.3	14.2 (KDforAA)	11.39 \pm 0.02	11.8	11.45 (ViT)

Устойчивость к label noise

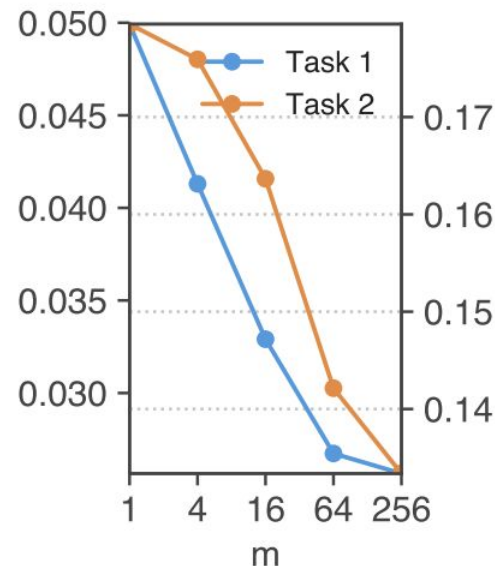
- SAM более устойчив к label noise, потому что он может вызывать неверные изменения поверхности loss-a

Method	Noise rate (%)			
	20	40	60	80
Sanchez et al. (2019)	94.0	92.8	90.3	74.1
Zhang & Sabuncu (2018)	89.7	87.6	82.7	67.9
Lee et al. (2019)	87.1	81.8	75.4	-
Chen et al. (2019)	89.7	-	-	52.3
Huang et al. (2019)	92.6	90.3	43.4	-
MentorNet (2017)	92.0	91.2	74.2	60.0
Mixup (2017)	94.0	91.5	86.8	76.9
MentorMix (2019)	95.6	94.2	91.3	81.0
SGD	84.8	68.8	48.2	26.2
Mixup	93.0	90.0	83.8	70.2
Bootstrap + Mixup	93.3	92.0	87.6	72.0
SAM	95.1	93.4	90.5	77.9
Bootstrap + SAM	95.4	94.2	91.8	79.9

m-sharpness

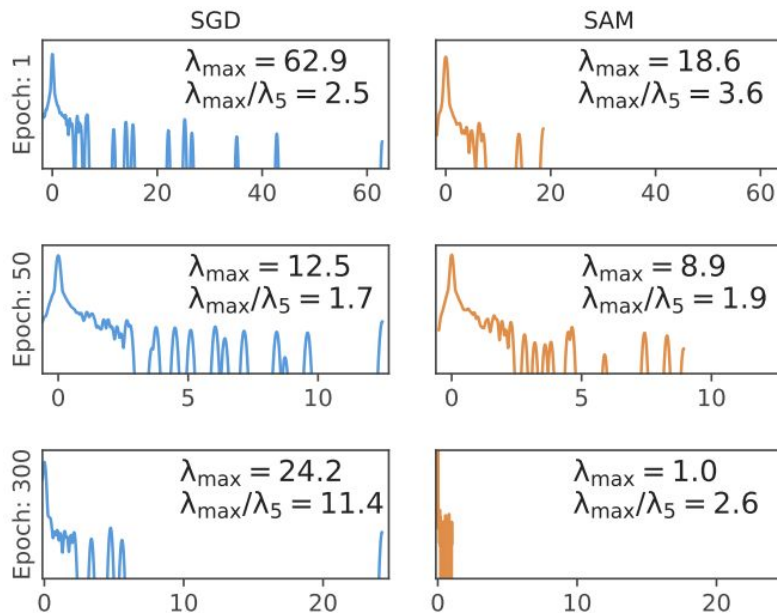
- Выберем несколько случайных подмножеств батча размера m . Для каждого из них вычислим $\hat{\epsilon}(\mathbf{w})$, а затем возьмем среднее SAM update-ов в качестве общего SAM update
- С меньшими значениями m получаются лучше

Эксперимент: маленький ResNet на CIFAR-10



Спектр Гессиана

- Посмотрим на спектр Гессиана в точке оптимума. Ожидаем, собственные числа близки к 1 в случае оптимизации с помощью SAM



Выводы

Плюсы:

- Новый метод регуляризации, применимый к произвольным задачам
- Есть теоретическое обоснование, геометрическая интерпретация
- Обновлено state of the art в задачах компьютерного зрения
- Есть исходный код для применения в популярных фреймворках, замедляет вычисление в 2 раза

Минусы:

- Введение m-sharpness оставляет некоторые вопросы, данная часть статьи выглядит неполной
- Улучшение оценки SAM-loss на практике не работает
- SOTA на finetune бьется бейзлайном

Информация о статье

- ICLR 2021 (можно [посмотреть](#) рассказ авторов)
- работа сделана как часть Google AI Residency program
- более 60 цитирований

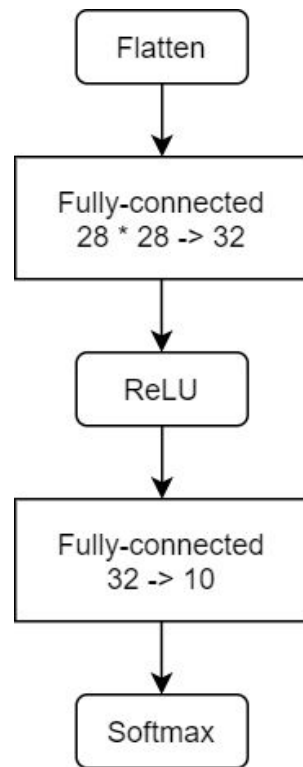
Другие статьи

- Exploring the Vulnerability of Deep Neural Networks: A Study of Parameter Corruption ([ссылка](#)) - представлен индикатор устойчивости модели
- Improved Sample Complexities for Deep Networks and Robust Classification via an All-Layer Margin ([ссылка](#)) - all-layer margin

В практической части попробовали реализовать:

- Обучение полносвязной сети на MNIST с большим значением label noise
- Обучение свёрточной сети на CIFAR-10 с большим значением label noise
- Сравнение направления спуска на избранных функционалах с двумерным доменом

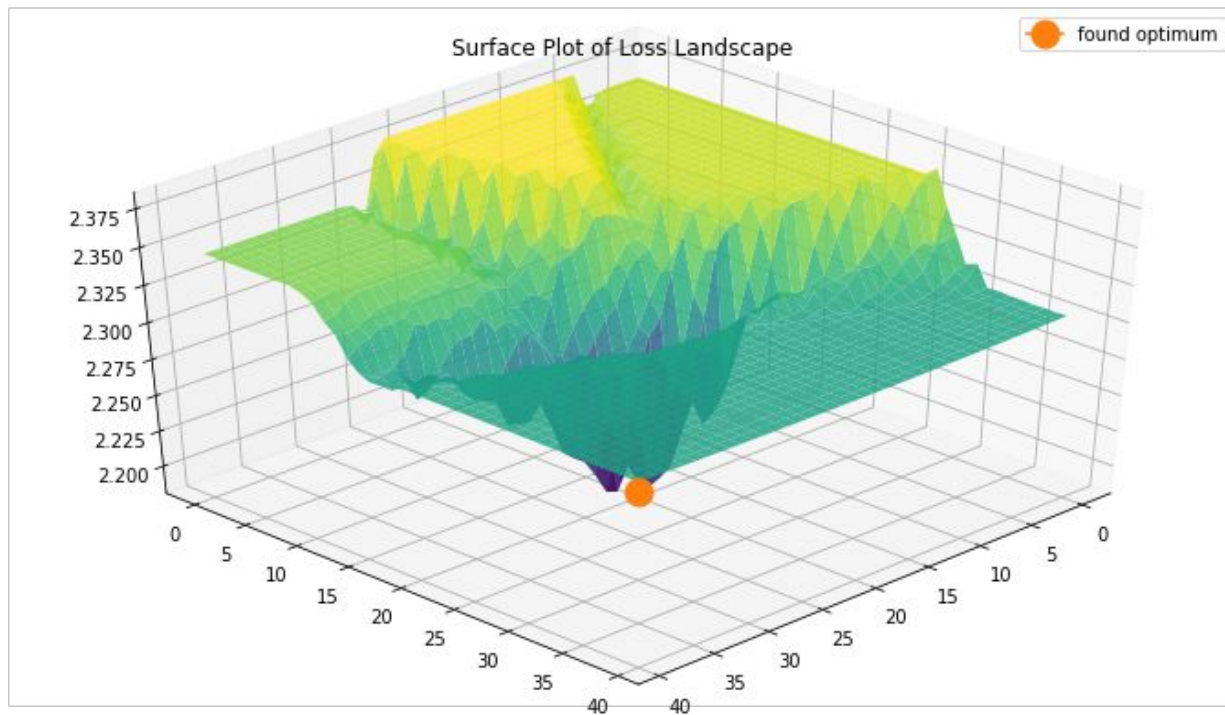
Обучение полносвязной сети на MNIST с большим значением label noise (0.8)



Эффект от SAM негативный

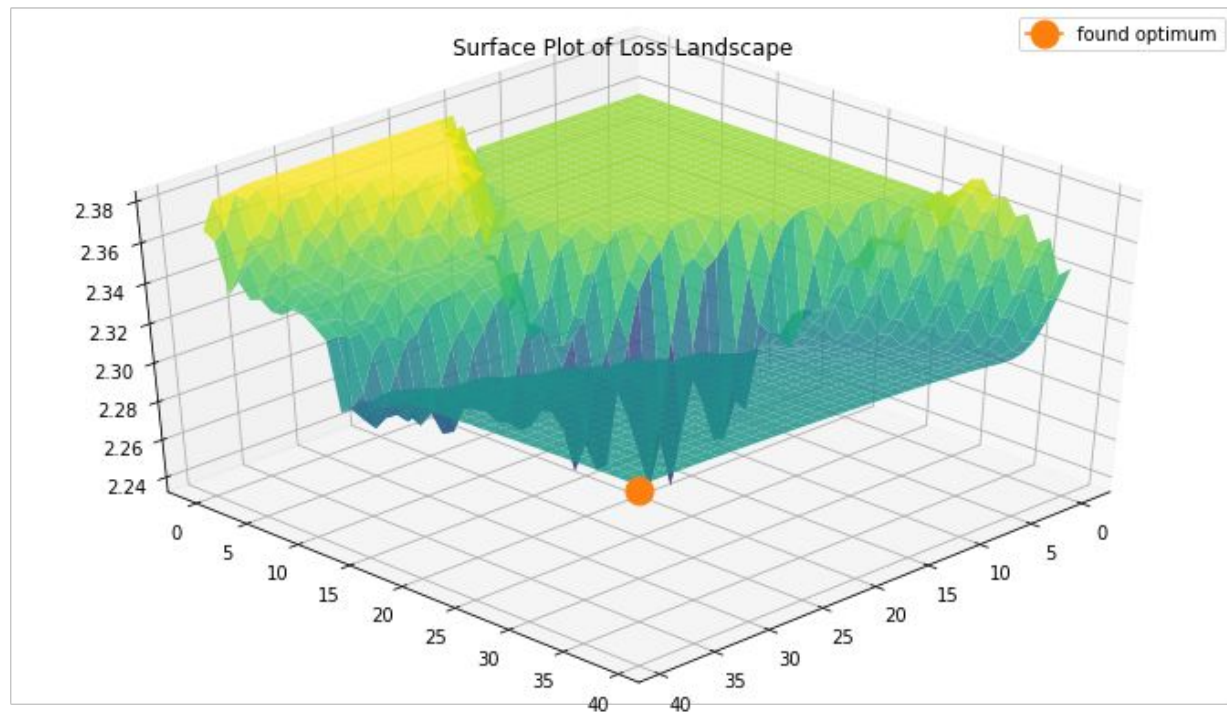
	Adam	Adam with SAM
accuracy	0.88	0.46

Эффект от SAM негативный



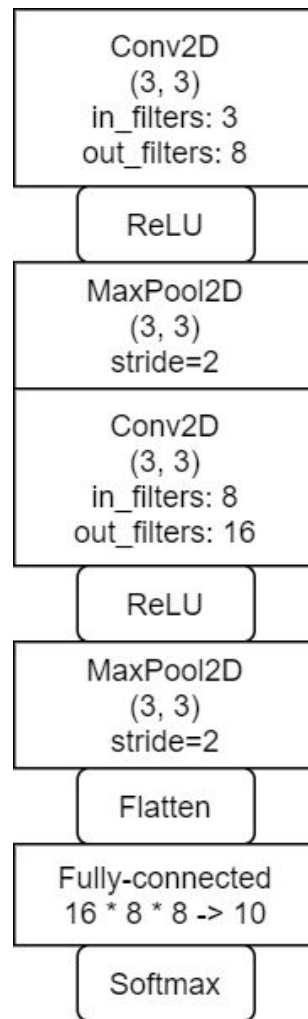
Визуализация
ландшафта
функции потерь
модели без SAM

Эффект от SAM негативный



Визуализация
ландшафта
функции потерь
модели с
использованием
SAM

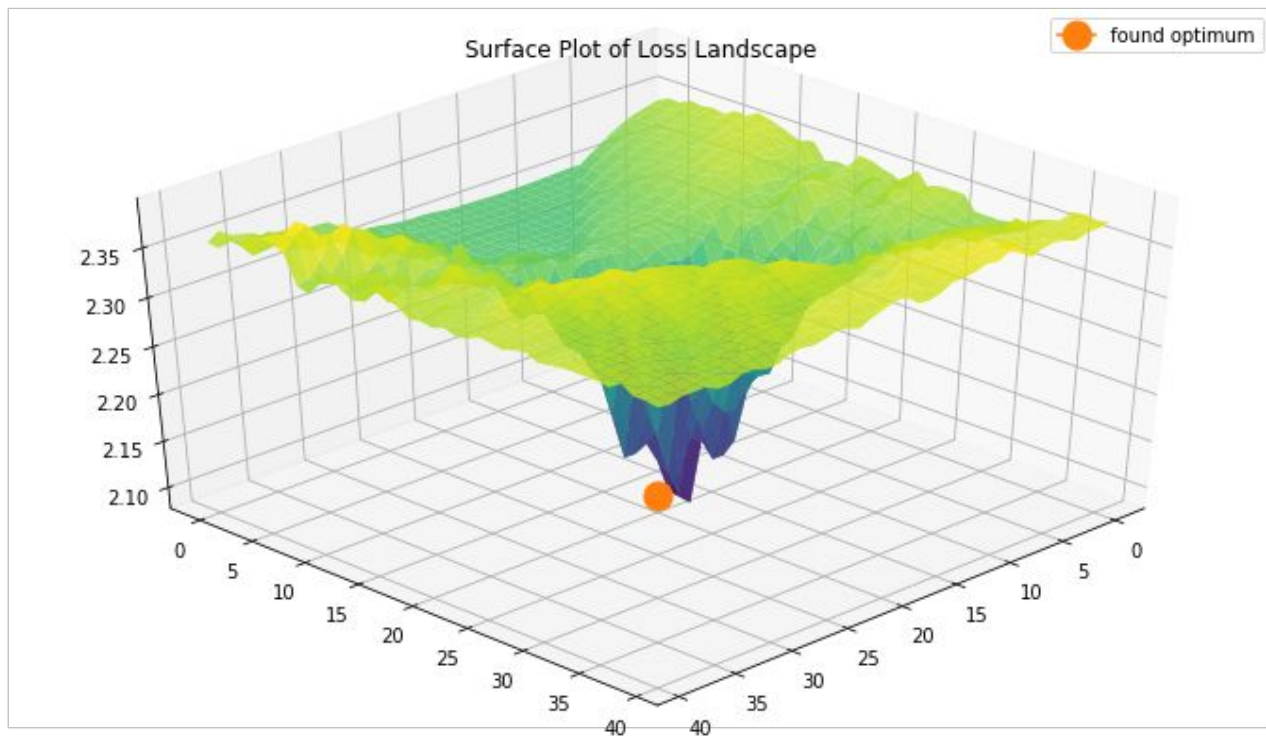
Обучение свёрточной сети на CIFAR-10 с большим значением label noise (0.5)



Эффект от SAM нейтральный

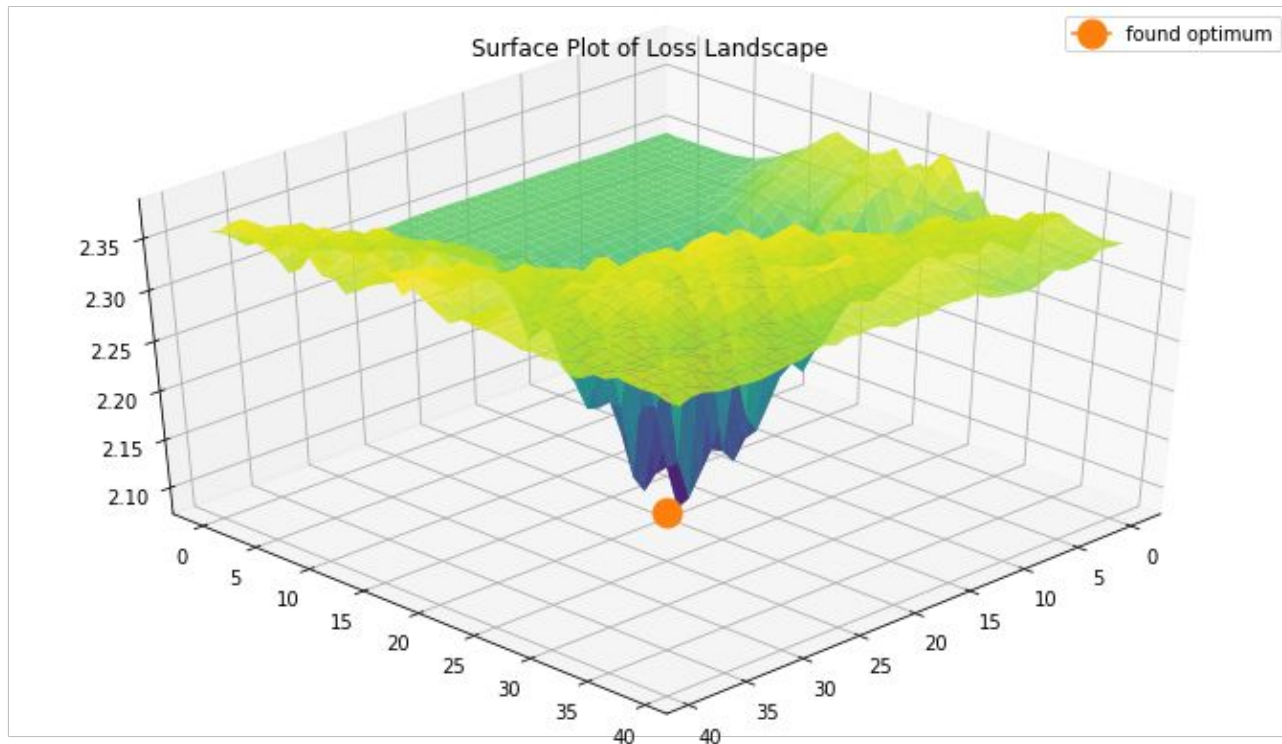
	Adam	Adam with SAM
accuracy	0.47	0.48

Эффект от SAM нейтральный



Визуализация
ландшафта
функции потерь
модели без SAM

Эффект от SAM нейтральный



Визуализация
ландшафта
функции потерь
модели с
использованием
SAM

Сравнение на функциях $\mathbb{R}^2 \rightarrow \mathbb{R}$

- Easom function

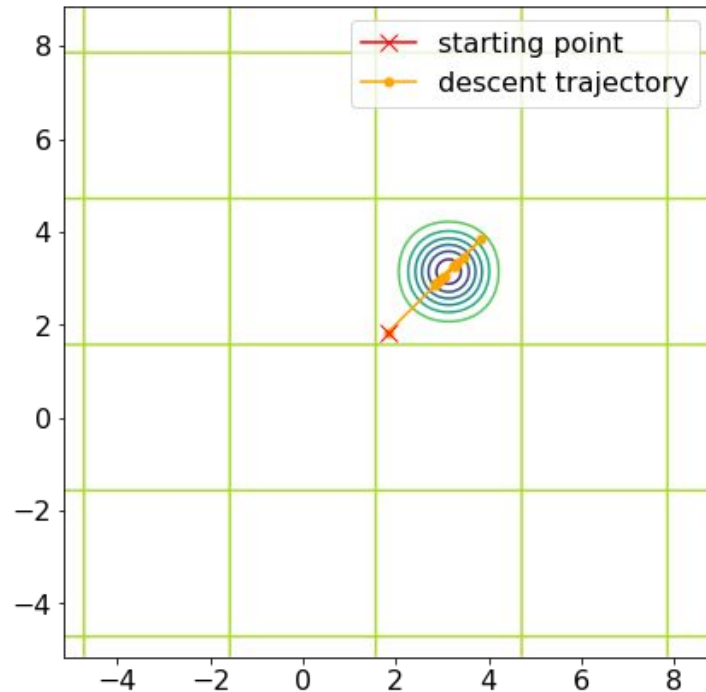
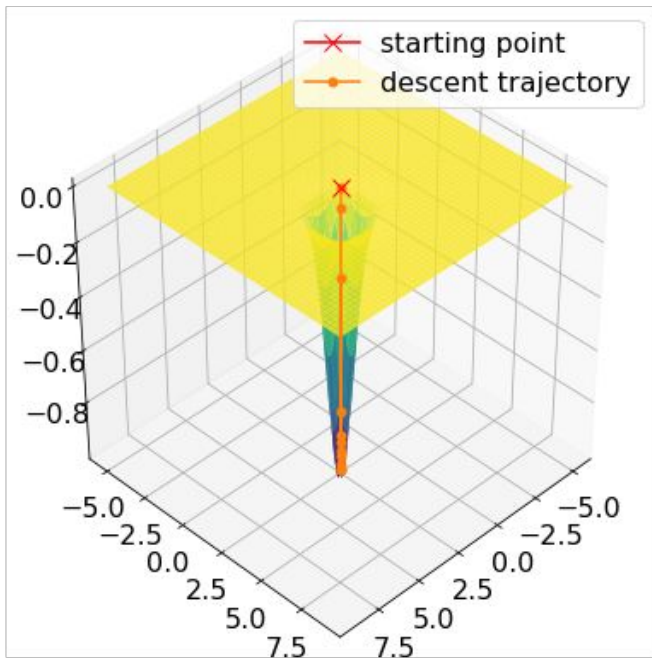
$$f(x, y) = -\cos(x) \cos(y) \exp\left(-\left((x - \pi)^2 + (y - \pi)^2\right)\right)$$

- Eggholder function

$$f(x, y) = -(y + 47) \sin \sqrt{\left|\frac{x}{2} + (y + 47)\right|} - x \sin \sqrt{|x - (y + 47)|}$$

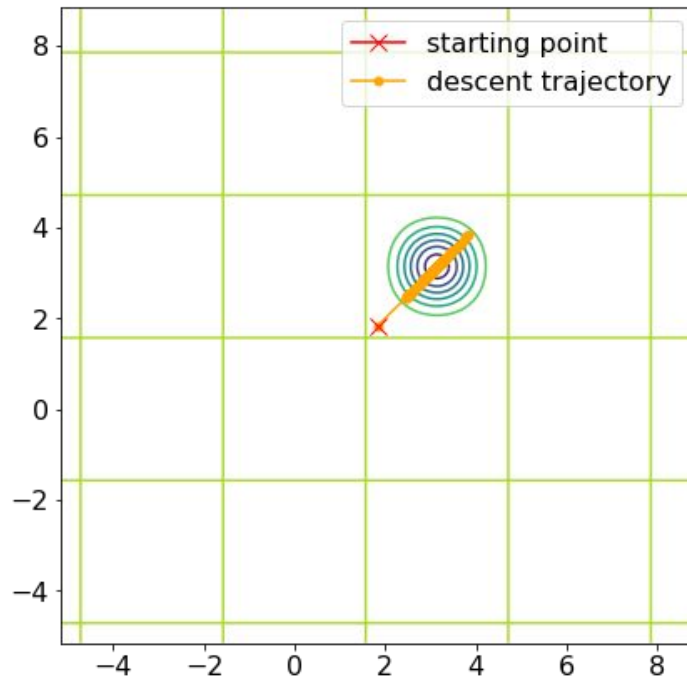
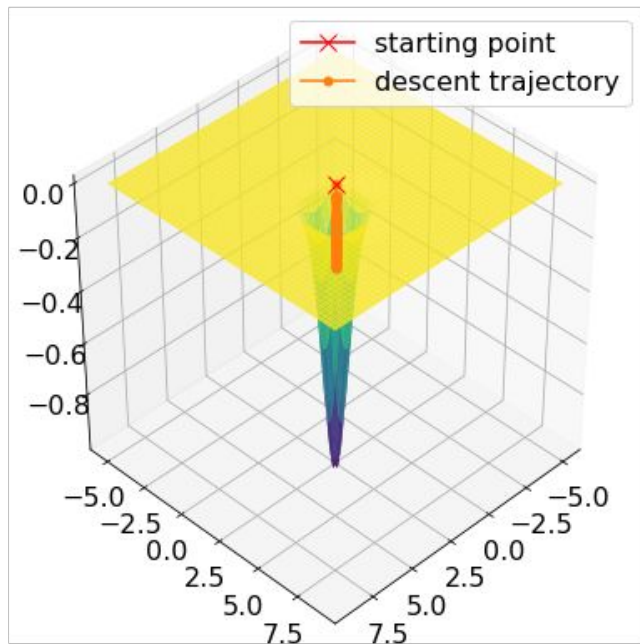
Easom function

Loss Landscape and Adagrad descent trajectory



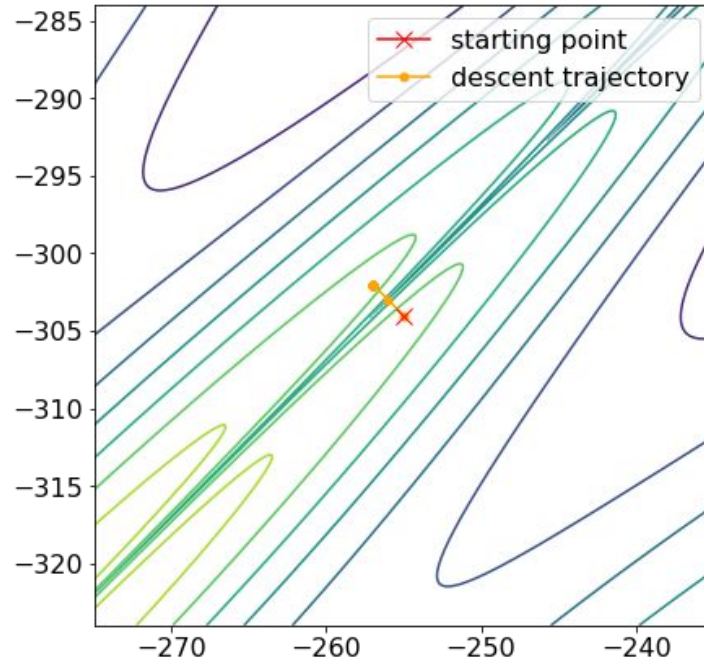
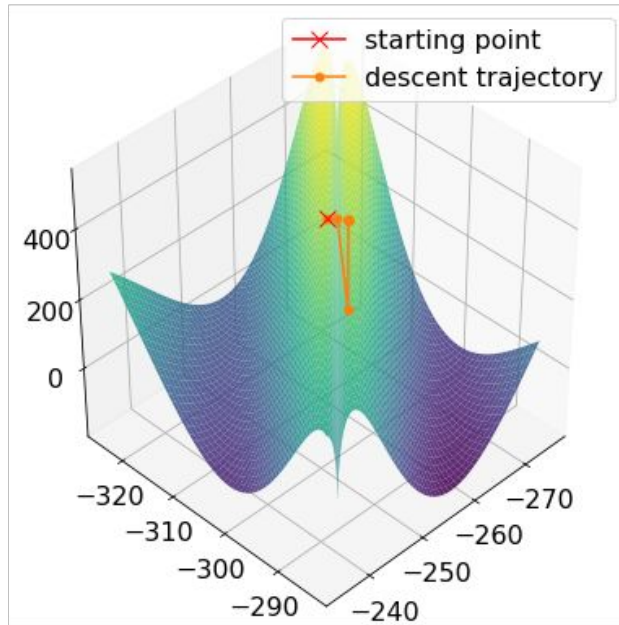
Easom function

Loss Landscape and SAM Adagrad descent trajectory



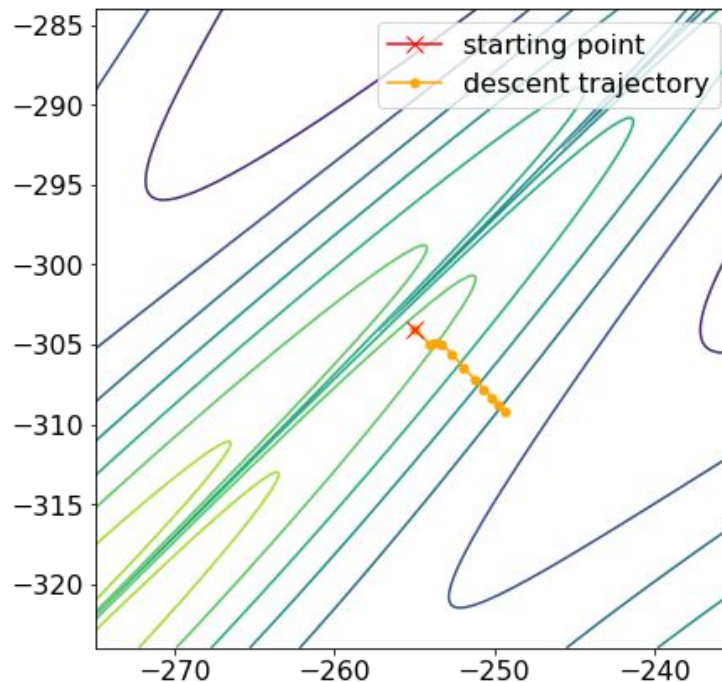
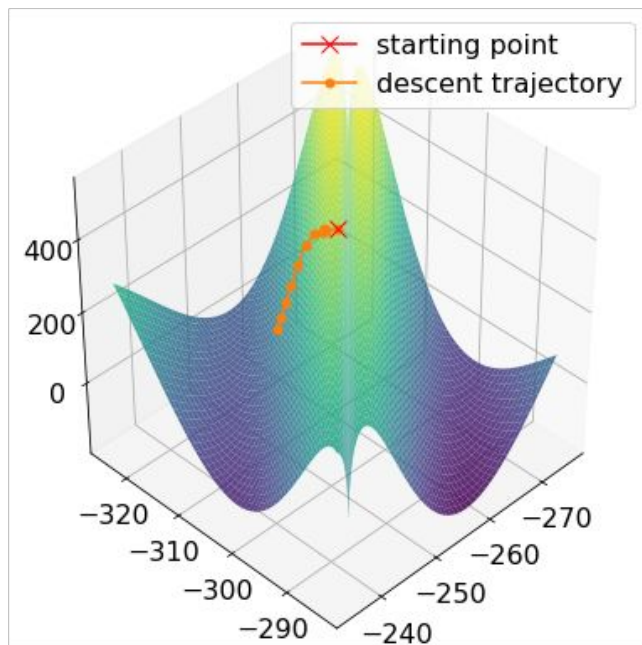
Eggholder function

Loss Landscape and Adagrad descent trajectory



Eggholder function

Loss Landscape and SAM Adagrad descent trajectory



Выводы практической части

- SAM действительно не позволяет упасть в “острые” минимумы.
- Применимость SAM для простых моделей под вопросом (даже с сильным label noise).

ИСТОЧНИКИ

- <https://arxiv.org/abs/2010.01412>
- <https://github.com/google-research/sam>