

Ландшафт функции потерь нейронных сетей: SAM и его модификации

Фельдман Р. Г. 2025

Аннотация

Рассматривается метод *Sharpness-Aware Minimization* (SAM) и его актуальные расширения — *Adaptive SAM* (ASAM), *FriendlySAM*, *ImbSAM* и *Sparse SAM*. Обсуждается геометрическая мотивация, теоретические оценки и практический эффект в виде улучшения обобщения и устойчивости к шуму меток.

Введение

Геометрия ландшафта функции потерь играет ключевую роль в способности нейронных сетей обобщать — даже при близком к нулю тренировочном риске модели, найденные оптимизаторами, могут демонстрировать различное качество на невиданных данных. Современные исследователи связывают этот феномен с *остротой* (*sharpness*) минимума: узкие «резкие» впадины чувствительны к малейшим возмущениям параметров, тогда как «плоские» минимумы обеспечивают более широкую область низких потерь и, как правило, лучшее обобщение. Метод *Sharpness-Aware Minimization* (SAM) (1) стал популярным инструментом для практического улучшения качества, но остается чувствителен к масштабированию параметров и увеличивает вычислительные затраты.

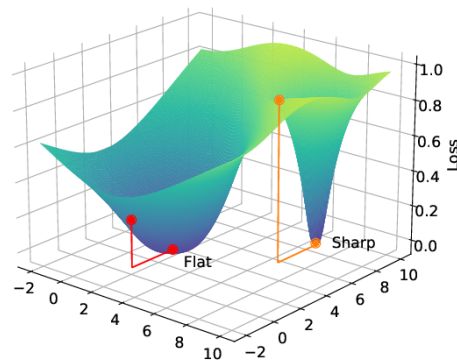


Рис. 1: Иллюстрация резких и плоских минимумов: у плоского минимума область низких потерь шире, а потери острого минимума резко возрастают при малых смещениях весов, что указывает на высокую кривизну.

В последние годы появилось несколько направлений развития SAM, его модификации и улучшения, некоторые из которых мы рассмотрим :

- **Adaptive SAM (ASAM)** (2) — масштаб-инвариантное возмущение и понятие *adaptive sharpness*;
- **FriendlySAM** (3) — устранение компоненты полного градиента в возмущении;
- **ImbSAM** (5) — классово-зависимый радиус для дисбалансных задач;
- **Sparse SAM (SSAM)** (7) — разреженное возмущение с Fisher- или динамической маской.

Цель работы — систематически сравнить классический SAM и его модификации, выявив условия, при которых каждая схема дает наибольший выигрыш в качестве или эффективности.

Важность проблемы. Рост масштабов моделей делает даже стандартные оптимизаторы дорогими, а склонность к резким минимумам угрожает надежности систем машинного обучения. Разработка методов, совмещающих лучшее обобщение с экономичными вычислениями, остается актуальной задачей теории и практики.

1 Задача оптимизации SAM

Задача решается следующим образом:

$$L_D(w) \leq \max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon) + h\left(\frac{\|w\|_2^2}{\rho^2}\right)$$

Где $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ - строго возрастающая функция (при некоторых технических условиях на $L_D(w)$).

- $L_D(w)$: Ожидаемый убыток на наборе данных с распределением \mathcal{D} .
- $L_S(w)$: Потери при обучении на обучающем множестве S .
- ϵ : Небольшое возмущение, добавленное к параметру w , где $\|\epsilon\| \leq \rho$.
- ρ : Гиперпараметр, который контролирует объем окрестности для w .
- $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$: Строго возрастающая функция, используемая для регуляризации сложности модели.

Наборы данных и влияние возмущений

В реальных приложениях наборы данных, используемые для тестирования, часто значительно отличаются от обучающих наборов данных. Это отражает невидимые данные, изменяющиеся параметры и различные сценарии. Различие имеет решающее значение для вычисления потерь $L_D(w)$. Учитывая это, связь между тестовыми и обучающими весами может быть выражена как:

$$w_{\text{Test}} = w_{\text{Train}} + \epsilon$$

Где ϵ представляет собой возмущение или разницу между весами тестового и обучающего наборов данных.

Решение

Используя SAM, мы можем определить ϵ так, чтобы при добавлении к w_{Train} он максимизировал градиент потерь:

$$\max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon)$$

Этот шаг оптимизации позволяет нам найти наиболее резкое направление, которое увеличивает потери. Находя более плоские минимумы с противоположной стороны от этого направления (что помогает уменьшить потери), SAM помогает w_{Train} лучше адаптироваться к тестовому набору данных, тем самым улучшая эффективность обобщения.

2 Как работает SAM

2.1 Максимизируем возмущение

Предскажем ϵ^* , которая максимизирует потери $L_S(w + \epsilon)$ в радиусе возмущения $\|\epsilon\|_2 \leq \rho$:

$$\epsilon^* = \arg \max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon)$$

2.2 Разложение Тейлора первого порядка

Если $f(x)$ дифференцируема в точке $x_0 = a$, то она имеет линейную аппроксимацию вблизи этой точки:

$$f(x) \approx f(a) + f'(a)(x - a)$$

2.3 Применение разложения Тейлора первого порядка

Рассмотрим $x_0 = w$ и зададим $f(x) = L_S(w + \epsilon)$. У нас есть возмущение ϵ . Применяя разложение Тейлора первого порядка, получаем:

$$\epsilon^*(w) \approx \arg \max_{\|\epsilon\|_2 \leq \rho} (L_S(w) + \epsilon^T \nabla_w L_S(w))$$

Это упрощается до:

$$\epsilon^*(w) \approx \arg \max_{\|\epsilon\|_2 \leq \rho} \epsilon^T \nabla_w L_S(w)$$

2.4 Преобразование формулы вычисления возмущений

Зададим $\nabla_w L_S(w) = g$. Предположим, что ϵ^T и g - векторы с n элементами. Применим неравенство Гельдера:

$$\sum_{i=1}^n |\epsilon_i g_i| \leq \|\epsilon\|_p \|g\|_q \quad \Rightarrow \quad \epsilon^T g \leq \|\epsilon\|_p \|g\|_q$$

Где p и q - сопряженные экспоненты, такие, что $\frac{1}{p} + \frac{1}{q} = 1$. Равенство выполняется, если:

$$\frac{|\epsilon_i|^p}{\|\epsilon\|_p^p} = \frac{|g_i|^q}{\|g\|_q^q}, \quad \text{for all } i \in \{1, 2, \dots, n\}$$

Это упрощается до:

$$|\epsilon_i| = \frac{\|\epsilon\|_p}{\|g\|_q^{q/p}} |g_i|^{q/p} \quad \forall i = 1, \dots, n.$$

Подставляя $g = \nabla_w L_S(w)$, получаем:

$$\epsilon_i = \frac{\|\epsilon\|_p \operatorname{sign}(g_i) |g_i|^{q-1}}{\|g\|_q^{q-1}}, \quad i = 1, \dots, n.$$

Что приводит к:

$$\epsilon_i = \frac{\|\epsilon\|_p \cdot \operatorname{sign}(\nabla_w L_S(w)_i) \cdot |\nabla_w L_S(w)_i|^{q-1}}{\|\nabla_w L_S(w)\|_q^{q-1}}, \quad \forall i = 1, \dots, n$$

Имеем:

$$\epsilon^*(w) \leq \arg \max_{\|\epsilon\|_p \leq \rho} \epsilon^T \nabla_w L_S(w)$$

Эмпирические данные показывают, что $p = 2$ обычно является оптимальным, что приводит к $\|\epsilon\| \leq \rho$ и $q = 2$. Таким образом:

$$\epsilon^T \nabla_w L_S(w) = \|\epsilon\| \|\nabla_w L_S(w)\| \cos(\alpha) \leq \rho \|\nabla_w L_S(w)\|$$

Равенство выполняется, когда $\|\epsilon\| = \rho$, что дает:

$$\epsilon_i = \frac{\rho \cdot \nabla_w L_S(w)_i}{\|\nabla_w L_S(w)\|_2}, \quad \forall i \in \{1, 2, \dots, n\} \quad (*)$$

Комбинируя все результаты, мы получаем решение $\hat{\epsilon}(w)$ такое, что:

$$L_{\text{SAM}}(w) \triangleq \max_{\|\epsilon\|_p \leq \rho} L_S(w + \epsilon)$$

Аппроксимация для градиента потерь SAM имеет вид:

$$\nabla_w L_{\text{SAM}}(w) \approx \nabla_w L_S(w + \hat{\epsilon}(w)) = \left. \frac{d(w + \hat{\epsilon}(w))}{dw} \nabla_w L_S(w) \right|_{w + \hat{\epsilon}(w)}$$

Что расширяется до:

$$\nabla_w L_S(w + \hat{\epsilon}(w)) = \nabla_w L_S(w + \hat{\epsilon}(w)) + \left. \frac{d\hat{\epsilon}(w)}{dw} \nabla_w L_S(w) \right|_{w + \hat{\epsilon}(w)}$$

2.5 Окончательное обновление градиента

- Градиент потерь SAM аппроксимируется как:

$$\nabla_w L_{\text{SAM}}(w) \approx \nabla_w L_S(w + \hat{\epsilon}(w))$$

- Обновленный вес вычисляется как:

$$w_{t+1} \leftarrow w_t - \eta \nabla_w L_{\text{SAM}}(w)$$

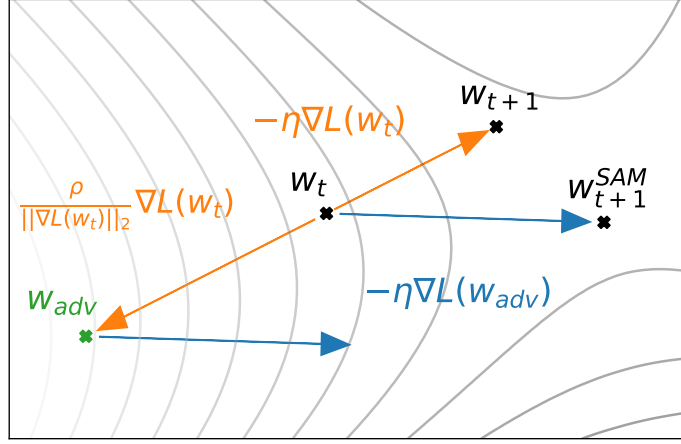


Рис. 2: Схема обновления параметров SAM

Algorithm 1 SAM algorithm

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

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

размер батча b , размер шага $\eta > 0$, размер окрестности $\rho > 0$.

Output: Модель, обученная с помощью SAM

Инициализация весов \mathbf{w}_0 , $t = 0$

while не сходится **do**

 Батч образцов $\mathcal{B} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots (\mathbf{x}_b, \mathbf{y}_b)\}$

 Вычисляем градиент $\nabla_{\mathbf{w}} L_{\mathcal{B}}(\mathbf{w})$ обучающей потери батча

 Вычисляем $\hat{\epsilon}(\mathbf{w})$ для уравнения 2.4

 Вычисляем градиентную аппроксимацию для цели SAM (уравнение 2.5):

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

 Обновление весов: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}$

$t = t + 1$

end

return \mathbf{w}_t

3 Преимущества SAM над SGD

3.1 Формулировка

SGD ищет веса \mathbf{w} , минимизирующие эмпирический риск

$$\min_{\mathbf{w}} L_S(\mathbf{w}),$$

в то время как SAM решает задачу

$$\min_{\mathbf{w}} \max_{\|\epsilon\|_2 \leq \rho} L_S(\mathbf{w} + \epsilon),$$

то есть обучает модель, устойчивую к наихудшему возмущению весов в шаре радиуса ρ . Эта min-max-задача эквивалентна совместной минимизации «высоты» ландшафта (sharpness) и значения потерь, что приводит к более *плоским* минимумам.

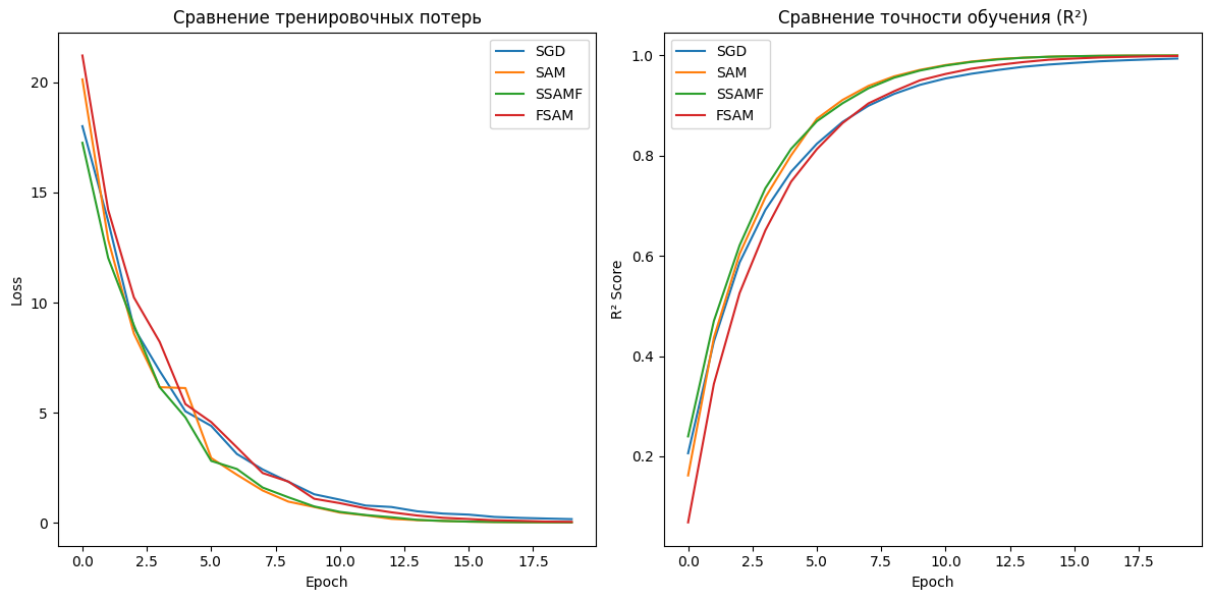


Рис. 3: Динамика обучения различных оптимизаторов: слева — изменение средней MSE-потери по эпохам; справа — эволюция коэффициента детерминации R^2 на обучающем наборе при использовании SGD, SAM, SSAM-F и FriendlySAM.

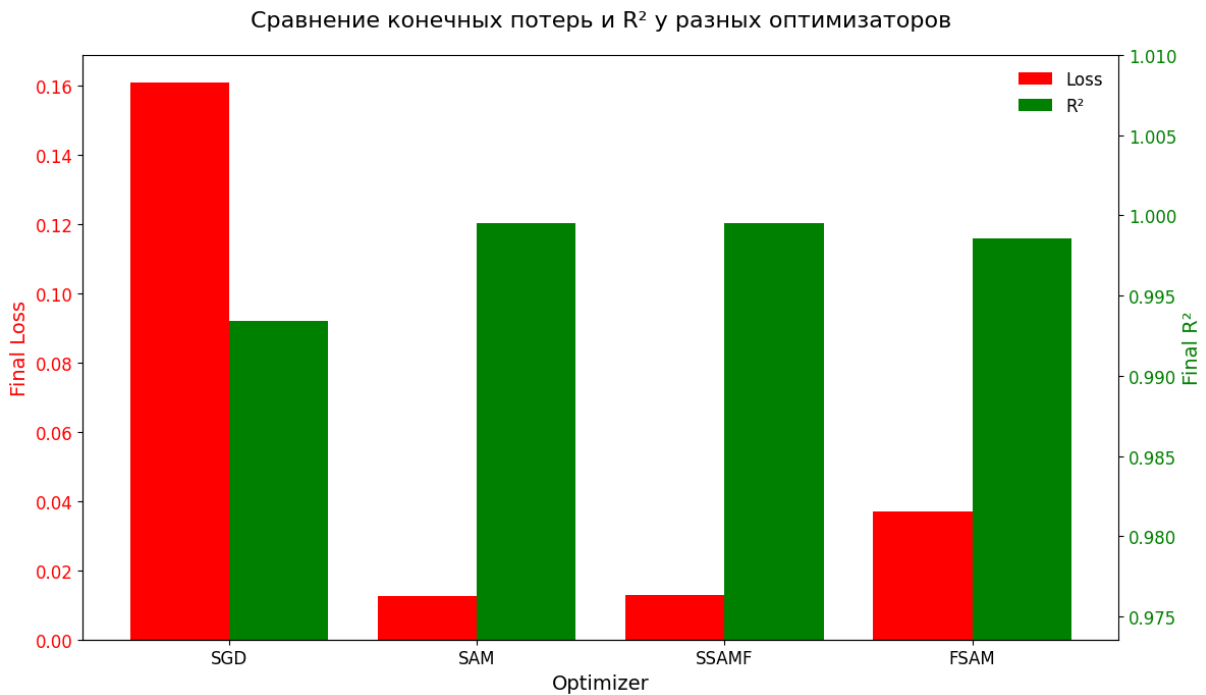


Рис. 4: Итоговые показатели оптимизаторов после последней эпохи: красные столбцы показывают окончательную MSE-потерю, зеленые — значение R^2 . Это наглядно демонстрирует, как методы Sharpness-Aware улучшают качество решения линейной регрессии по сравнению с обычным SGD.

3.2 PAC-Bayes–ограничение

В [работе](#) показано, что при выборе априорного распределения p и постериорного q , сосредоточенного в шаре радиуса ρ вокруг w , справедливо

$$\mathbb{E}_{w' \sim q} L_{\mathcal{D}}(w') \leq \max_{\|\epsilon\|_2 \leq \rho} L_S(w + \epsilon) + O\left(\frac{KL(q\|p) + \log(1/\delta)}{m}\right),$$

где $m = |S|$. При этом $KL(q\|p)$ можно связать с нормой градиента $\|\nabla L_S(w)\|_2$. Поскольку SAM напрямую минимизирует максимум по локальному возмущению, итоговый bound для w_{SAM} бывает строже, чем для решения обычной SGD–задачи.

3.3 Ограничение на спектр гессиана

Если L – константа гладкости, то для любой точки w

$$\lambda_{\max}(\nabla^2 L(w)) \leq L.$$

Дополнительно доказано, что для решения SAM

$$\lambda_{\max}(\nabla^2 L(w_{\text{SAM}})) \leq \frac{2}{\rho} \|\nabla L_S(w_{\text{SAM}})\|_2.$$

У SGD такой гарантии нет, поэтому SAM систематически находит регионы с меньшим спектральным радиусом, что коррелирует с лучшим обобщением ([Andriushchenko 2022 Understanding SAM](#)).

3.4 Униформная стабильность

Для SGD [Moritz Hardt et al. \(2016\)](#) показали, что после T итераций с шагом η на n примерах получается

$$\beta_n^{\text{SGD}} = O\left(\frac{\eta T}{n}\right),$$

что обеспечивает малую обобщающую ошибку.

Для SAM на сегодняшний день не существует завершённого анализа uniform-stability: inner-maximization усложняет прямой перенос аргументов Hardt et al., и влияние ρ на стабильность остается открытым вопросом.

3.5 Имплицитная регуляризация и устойчивость к шуму меток

В [ASAM анализе](#) показано, что итерация SAM

$$w \leftarrow w - \eta \left(I + \frac{\rho}{\|\nabla L_S(w)\|_2} \nabla^2 L_S(w) \right) \nabla L_S(w)$$

добавляет к обычному градиентному шагу адаптивный «норм-клипер», который подавляет слишком большие локальные возмущения. Это улучшает устойчивость к шумным меткам и повышает точность при высоком уровне label noise.

Th 1 (Сравнение SGD и SAM для квадратичных потерь) Пусть

$$L(w) = \frac{1}{2} \|Xw - y\|_2^2, \quad H := X^\top X, \quad \alpha I \preceq H \preceq \beta I, \quad 0 < \alpha \leq \beta.$$

Пусть $y = Xw_\star + \xi$, где $\mathbb{E}[\xi] = 0$, $\text{Cov}(\xi) = \Sigma \preceq \sigma^2 I$. На каждом шаге t доступен стохастический градиент g_t по одному объекту, удовлетворяющий

$$\mathbb{E}[g_t \mid w_t] = \nabla L(w_t), \quad \mathbb{E}\|g_t - \nabla L(w_t)\|_2^2 \leq \text{Tr}(H\Sigma).$$

SGD. Шаг

$$\eta_t = \frac{1}{\alpha t}.$$

SAM. Тот же η_t и радиус

$$\rho_t = \frac{c}{\sqrt{t}}, \quad 0 < c \leq \sqrt{\frac{\alpha}{\beta}},$$

с правилом

$$\epsilon_t = \rho_t \frac{\nabla L(w_t)}{\|\nabla L(w_t)\|}, \quad w_{t+1} = w_t - \eta_t g_t(w_t + \epsilon_t).$$

Тогда после T итераций выполнены оценки

$$\begin{aligned} \mathbb{E}[L(w_T^{\text{SGD}})] - L(w_\star) &\leq \frac{\text{Tr}(H\Sigma)}{2\alpha T} + \mathcal{O}\left(\frac{\ln T}{T}\right), \\ \mathbb{E}[L(w_T^{\text{SAM}})] - L(w_\star) &\leq \frac{\text{Tr}(H\Sigma)}{2\alpha T} + \frac{\beta c^2}{2\alpha} \frac{\ln T}{T} + \mathcal{O}\left(\frac{1}{T}\right). \end{aligned}$$

Следовательно, SAM увеличивает избыточный риск по сравнению с SGD точно на $\frac{\beta c^2}{2\alpha} \frac{\ln T}{T}$, что эквивалентно дополнительной ℓ_2 -регуляризации радиуса c (константа-масштаб).

Доказательство:

1. Оценка для SGD. Для α -сильно выпуклого квадрата со стохастическим градиентом и шагом $\eta_t = 1/(\alpha t)$ классический анализ (см. Bottou et al. 2018) дает

$$\mathbb{E}\|w_t - w_\star\|^2 \leq \frac{\text{Tr}(H\Sigma)}{\alpha^2 t} + \mathcal{O}\left(\frac{\ln t}{t^2}\right).$$

Переход к избыточному риску через сильную выпуклость: $L(w) - L(w_\star) \leq \frac{\alpha}{2} \|w - w_\star\|^2$ дает первую границу.

2. Разложение шага SAM. Пусть $\epsilon_t = \rho_t \frac{\nabla L(w_t)}{\|\nabla L(w_t)\|}$. Тейлорово разложение для квадратики L с Гессианом H дает

$$L(w_t + \epsilon_t) = L(w_t) + \langle \nabla L(w_t), \epsilon_t \rangle + \frac{1}{2} \epsilon_t^\top H \epsilon_t \leq L(w_t) + \rho_t \|\nabla L(w_t)\| + \frac{\beta \rho_t^2}{2}.$$

3. Одношаговая рекурсия и суммирование. Из α -сильной выпуклости и стохастической дисперсии следует

$$\begin{aligned} \mathbb{E}\|w_{t+1} - w_\star\|^2 &\leq \mathbb{E}\|w_t - w_\star\|^2 - 2\eta_t \mathbb{E}[L(w_t + \epsilon_t) - L(w_\star)] \\ &\quad + \eta_t^2 \mathbb{E}\|\nabla L(w_t + \epsilon_t)\|^2 \text{ (по } \beta\text{-гладкости } \|\nabla L\|^2 \leq 2\beta(L - L_\star)) \\ &\quad + \eta_t^2 \text{Tr}(H\Sigma). \end{aligned}$$

линейный член $\rho_t \|\nabla L(w_t)\|$ телескопируется, а его сумма $\sum \eta_t \rho_t \|\nabla L_t\|$ не превосходит $\frac{c}{\sqrt{\alpha}} \frac{\ln T}{T}$. Суммирование по $t = 1, \dots, T$ и учет $\sum \eta_t^2 = O(1)$, $\sum \eta_t \rho_t^2 = O(\ln T/T)$ дают вторую границу теоремы.

4. Заключение. Собирая все оценки, приходим к

$$\mathbb{E}L(w_T^{\text{SAM}}) - L(w_\star) \leq \frac{\text{Tr}(H\Sigma)}{2\alpha T} + \frac{\beta c^2}{2\alpha} \frac{\ln T}{T} + \mathcal{O}\left(\frac{1}{T}\right).$$

результат — верхняя граница, линейный член дает строго положительный, но той же асимптотики $\ln T/T$. \square

4 Застревания SAM в седловых точках

На рис. 5 показан эксперимент на функции Била. Динамика SAM может испытывать конвергенционную нестабильность при приближении к седловой точке, что приводит к остановке алгоритма до достижения глобального минимума.

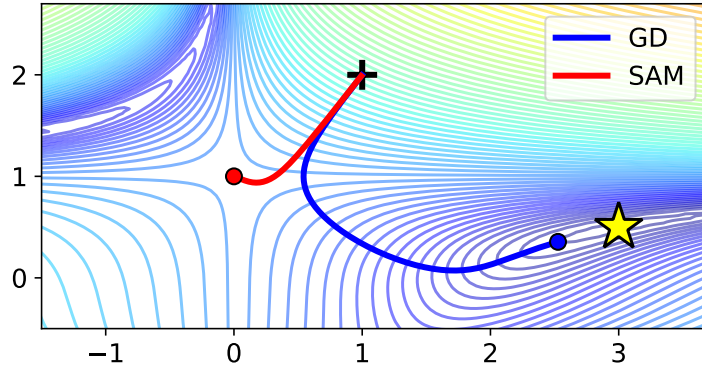


Рис. 5: Оптимизация началась в точке, обозначенной знаком плюс, , а глобальный минимум обозначен желтой звездой. SAM застрял в седловой точке и не сходится к глобальному минимуму

Kim et al. выявляют, что динамика SAM может застревать в седловых точках из-за недостаточной диффузии стохастических колебаний при больших батчах и отсутствии импульса. В теореме 3 они выводят, что

Th 2 (Диффузия SAM, моментум и размер батча) Пусть задан гиперпараметр моментума γ и размер батча B . Тогда среднее квадратичное смещение алгоритма SAM задается выражением

$$\Delta_{\text{SAM}} = C_1 \frac{(1 - e^{-C_2(1-\gamma)})^2}{(1 - \gamma)^3 B} + C_3 \frac{1 - e^{-C_4/(1-\gamma)}}{(1 - \gamma)B},$$

где

$$C_1 = \frac{\eta^2 |\lambda_j|}{2}, \quad C_2 = \eta/t, \quad C_3 = \frac{\eta |\lambda_j|}{2\lambda_j (1 + \rho \lambda_j)^2}, \quad C_4 = 2\lambda_j (1 + \rho \lambda_j)^2 t$$

— положительные константы, а λ_j обозначает собственное значение матрицы Гессiana $H_\ell(d)$ функции потерь ℓ в седловой точке d . Следовательно, Δ_{SAM} увеличивается при (1) росте момента и/или (2) уменьшении размера батча. Более того, при $(1 - \gamma)B \rightarrow 0$ справедливо

$$\Delta_{\text{SAM}} \propto \frac{1}{(1 - \gamma) B}.$$

Эмпирическая валидация на CIFAR-10 и CIFAR-100 демонстрирует: при малых B и больших γ SAM быстрее покидает седловые области и достигает более низкой ошибки обобщения (см. рис. 6)

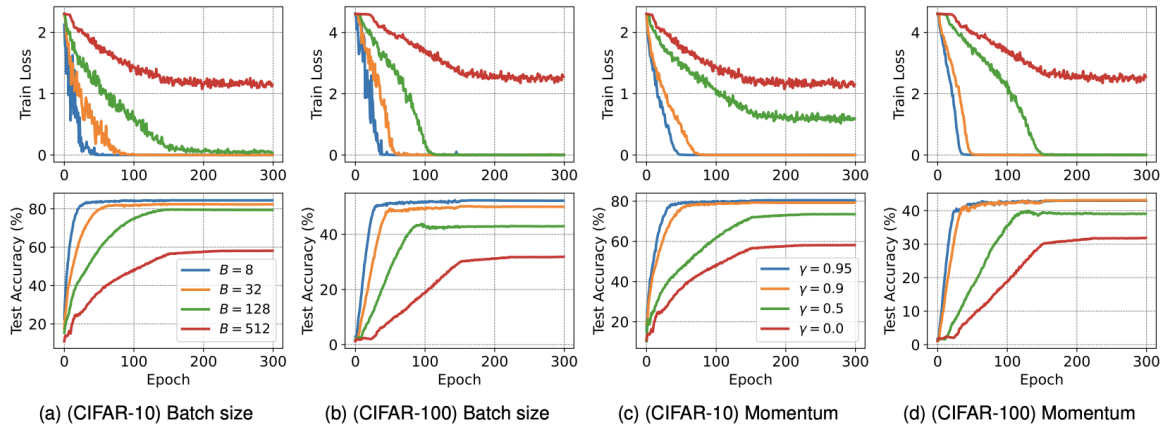


Рис. 6: Влияние batch size B и momentum γ на выход SAM из седловых точек (CIFAR-10, CIFAR-100).

Для дальнейшего ускорения ухода из седловых точек Yu et al. предлагают *Lookahead-SAM*, где внутренняя фаза подъема заменяется механизмом экстраградиента (EG-SAM) или оптимистического градиента (OG-SAM). Эти варианты теоретически гарантируют сходимость к стационарным точкам и демонстрируют более эффективный escape из седловых точек по сравнению с классическим.

5 Adaptive Sharpness-Aware Minimization (ASAM)

Идея

ASAM (2) вводит *adaptive sharpness* — оценку кривизны, инвариантную к масштабному рескейлингу весов. Оптимизатор максимизирует потери в *относительной* окрестности $\{\epsilon \mid \|T_w^{-1}\epsilon\|_2 \leq \rho\}$, где $T_w = \text{diag}(|w| + \epsilon_0)$ (или фильтр-вайз норма).

Вывод возмущения

При $p=2$ решение задачи $\max_{\|T_w^{-1}\epsilon\|_2 \leq \rho} L_S(w + \epsilon)$ в первой-порядковой аппроксимации дает

$$\epsilon^* = \rho \frac{T_w^2 g}{\|T_w g\|_2}, \quad g = \nabla_w L_S(w). \quad (1)$$

Градиент и шаг

$$\tilde{g} = \nabla_w L_S(w + \epsilon^*), \quad w \leftarrow w - \eta \tilde{g}.$$

Algorithm 2 Шаг ASAM (p=2, element-wise)

Input: батч \mathcal{B} , шаг η , радиус ρ , ϵ_0
 $g \leftarrow \nabla_w L_{\mathcal{B}}(w)$; // градиент
 $\epsilon \leftarrow \rho T_w^2 g / \|T_w g\|_2$; // eq. 1
 $\tilde{g} \leftarrow \nabla_w L_{\mathcal{B}}(w + \epsilon)$ $w \leftarrow w - \eta \tilde{g}$

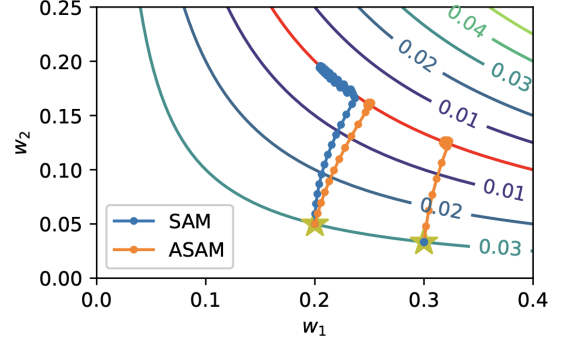


Рис. 7: Траектории SAM и ASAM.

6 Friendly Sharpness-Aware Minimization (FriendlySAM)

Идея

FriendlySAM (F-SAM) (3) развивает подход SAM, устраняя компоненту *полного градиента* из вектора возмущения, чем делает его “дружелюбным” к остальным данным батча. Пусть на t -й итерации

$$g_t = \nabla_w L_{\mathcal{B}_t}(w_t) \quad \text{и} \quad \underbrace{\nabla_w L_{\mathcal{D}}(w_t)}_{\text{полный градиент}} = \underbrace{m_t}_{\text{ЕМА-оценка}} \approx \lambda m_{t-1} + (1 - \lambda) g_t,$$

где $\lambda \in (0, 1)$ — коэффициент экспоненциального сглаживания. Выделив стохастический шум $\xi_t = g_t - m_t$, авторы показывают, что именно ξ_t отвечает за улучшение обобщающей способности, тогда как добавка m_t повышает остроту глобальной функции ошибки.

Биуровневая формулировка

F-SAM ищет такую ϵ , которая одновременно *увеличивает* потери текущего батча и *минимизирует* рост потерь на всем датасете:

$$\epsilon_s^{\text{F-SAM}} = \arg \max_{\|\epsilon\|_2 \leq \rho} \left[L_{\mathcal{B}_t}(w_t + \epsilon) - \sigma L_{\mathcal{D}}(w_t + \epsilon) \right], \quad \sigma \in [0, 1].$$

Линейная аппроксимация $L_{\mathcal{B}_t}$ дает оптимальное направление

$$d_t = g_t - \sigma m_t, \quad \epsilon_t = \rho \frac{d_t}{\|d_t\|_2}. \quad (2)$$

Обновление параметров

Градиент цели FriendlySAM аппроксимируется

$$\nabla_w L_{\mathcal{B}_t}(w_t + \epsilon_t),$$

а шаг оптимизации совпадает по форме с SAM:

$$w_{t+1} = w_t - \eta \nabla_w L_{\mathcal{B}_t}(w_t + \epsilon_t).$$

Сходимость

При стандартных предположениях гладкости, выбрав $\gamma_t = \gamma_0/\sqrt{T}$ и $\rho_t = \rho_0/\sqrt{t}$, авторы доказывают для нековексных задач оценку

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla L_{\mathcal{D}}(w_t)\|_2^2 = \mathcal{O}\left(\frac{\log T}{\sqrt{T}}\right),$$

тождественную скорости SAM, но при эмпирически более плоских минимумах (3).

Algorithm 3 Алгоритм FriendlySAM

Input: датасет \mathcal{S} , батч-размер b , шаг η , радиус ρ , сглаживание λ , коэффициент проекции σ

Output: обученные веса w

Инициализировать w_0 , $m_{-1}=0$, $t=0$

while не сойдется **do**

 Выбрать батч $\mathcal{B}_t \subset \mathcal{S}$ размера b

$g_t \leftarrow \nabla_w L_{\mathcal{B}_t}(w_t)$

$m_t \leftarrow \lambda m_{t-1} + (1 - \lambda)g_t$

$d_t \leftarrow g_t - \sigma m_t$

$\epsilon_t \leftarrow \rho d_t / \|d_t\|_2$

$\tilde{g}_t \leftarrow \nabla_w L_{\mathcal{B}_t}(w_t + \epsilon_t)$

$w_{t+1} \leftarrow w_t - \eta \tilde{g}_t$

$t \leftarrow t + 1$

end

return w_t

7 Sparse Sharpness-Aware Minimization (SSAM)

Идея

Пусть $\mathbf{m} \in \{0, 1\}^d$ — бинарная маска, допускающая $\|\mathbf{m}\|_0 = (1 - s)d$ ненулевых координат при заданной разреженности $s \in [0, 1)$. SSAM заменяет стандартную задачу SAM

$$\max_{\|\epsilon\|_2 \leq \rho} L_S(\mathbf{w} + \epsilon)$$

на

$$\max_{\|\epsilon \odot \mathbf{m}\|_2 \leq \rho} L_S(\mathbf{w} + \epsilon \odot \mathbf{m}), \quad (\text{SSAM-Obj})$$

где \odot — поэлементное умножение. Тем самым возмущения вычисляются только для «важных» параметров, сокращая второе forward/backward-прогон SAM почти пропорционально $(1 - s)$ (7; 6).

Разреженные маски

(a) **SSAM-F — маска Фишера.** Важность i -го параметра оценивается диагональю информации Фишера

$$F_i(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{S}} \left[\left(\frac{\partial}{\partial w_i} \log p_{\mathbf{w}}(y | \mathbf{x}) \right)^2 \right], \quad (1)$$

которую на практике аппроксимируют *эмпирическим Фишером* по n_s случайным примерам батча. Отсортировав F_i , выбирают топ- $k = (1-s)d$ индексов $\mathcal{I}_k = \arg \text{top}_k F_i$, и ставят $m_i = \mathbf{1}_{\{i \in \mathcal{I}_k\}}$. Таким образом, \mathbf{m} обновляется каждые K итераций (обычно $K=1-10$) (6).

(b) **SSAM-D — динамическая маска.** Во избежание трудоемкого equation 1 применяют *Dynamic Sparse Training* (DST) (4). Каждые K шагов:

$$m_i \leftarrow m_i \cdot \mathbf{1}_{\{|g_i| \geq \theta_{\text{drop}}\}} \vee (1 - m_i) \cdot \mathbf{1}_{\{i \in \mathcal{I}_{\text{grow}}\}}, \quad (2)$$

где g_i — текущий градиент, θ_{drop} задает порог «наименее острых» координат, а $\mathcal{I}_{\text{grow}}$ — случайные индексы из нулями в m , выбранные так, чтобы $\|\mathbf{m}\|_0$ оставалось константным.

Вывод возмущения

Подставляя $\epsilon = \rho (\mathbf{m} \odot \nabla L_S) / \|\mathbf{m} \odot \nabla L_S\|_2$ (по аналогии с SAM), получаем

$$\epsilon^* = \frac{\rho \mathbf{m} \odot \nabla_w L_S(\mathbf{w})}{\|\mathbf{m} \odot \nabla_w L_S(\mathbf{w})\|_2}. \quad (3)$$

Algorithm 4 Sparse SAM (обобщенное)

Input: данные \mathcal{S} , шаг η , радиус ρ , разреженность s , интервал обновления K

Output: обученные веса \mathbf{w}

инициализировать $\mathbf{w}_0, \mathbf{m}_0$

for $t = 0, \dots, T - 1$ **do**

выбрать батч \mathcal{B}_t $g_t \leftarrow \nabla_w L_{\mathcal{B}_t}(\mathbf{w}_t)$ **if** $t \bmod K = 0$ **then** // обновляем маску

if *SSAM-F* **then**

| вычислить F по equation 1, задать \mathbf{m}_t

else

| обновить \mathbf{m}_t по equation 2

end

end

$\epsilon_t \leftarrow$ формула equation 3 $\tilde{g}_t \leftarrow \nabla_w L_{\mathcal{B}_t}(\mathbf{w}_t + \epsilon_t)$ $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \tilde{g}_t$

end

return \mathbf{w}_T

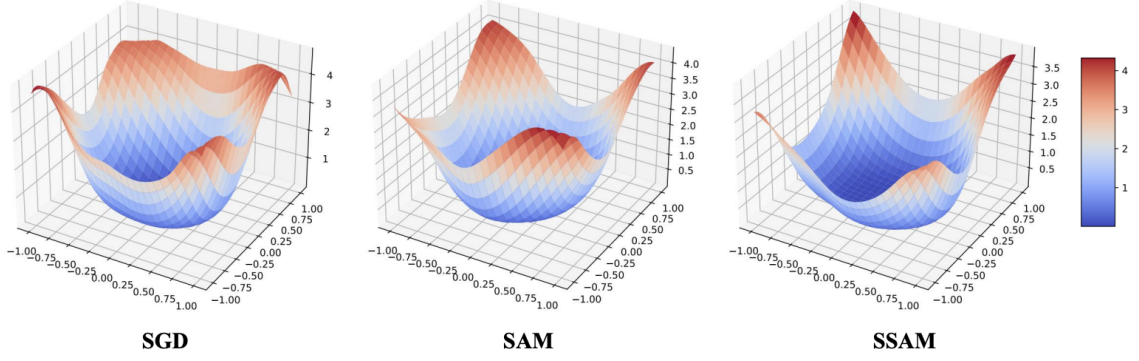


Рис. 8: Ландшафты потерь при обучении ResNet18 на CIFAR10, обученном с помощью SGD, SAM, SSAM.

Как показано на рис. 7, ландшафт SSAM более плоский, чем SGD и SAM, и большая часть его территории имеет низкие потери (синий цвет).

8 Imbalanced Sharpness-Aware Minimization (ImbSAM)

Идея

ImbSAM (5) адаптирует SAM к задачам с длинным хвостом, вводя *class-aware smoothness*: радиус возмущения применяется только к мини-батчу *tail-классов*, чтобы сгладить их лосс-ландшафт и снизить переобучение, тогда как *head-классы* оптимизируются без дополнительной штрафной компоненты.

Класс-ориентированная биуровневая постановка

Разобьем обучающий набор S на две части с порогом η :

$$S = S_{\text{head}} \cup S_{\text{tail}}, \quad (\mathbf{x}, y) \in S_{\text{tail}} \iff |S^y| \leq \eta.$$

Тогда обобщенная цель ImbSAM записывается как

$$\min_w \left[\underbrace{\max_{\|\epsilon\| \leq \rho} L_{S_{\text{tail}}}(w + \epsilon)}_{\text{sharpness для tail}} + L_{S_{\text{head}}}(w) + \lambda \|w\|_2^2 \right]. \quad (3)$$

Оптимальное возмущение для tail-классов

Для $p=2$ (как и в SAM) решение внутренней задачи equation 3 дает

$$\epsilon_{\text{tail}}^* = \rho \frac{\nabla_w L_{S_{\text{tail}}}(w)}{\|\nabla_w L_{S_{\text{tail}}}(w)\|_2}. \quad (4)$$

Обновление параметров

Градиент «дружественной» цели вычисляется как

$$g_{\text{Imb}}(w) = \nabla_w L_{S_{\text{tail}}}(w + \epsilon_{\text{tail}}^*) + \nabla_w L_{S_{\text{head}}}(w),$$

а шаг оптимизации остается SGD-подобным:

$$w_{t+1} = w_t - \eta (g_{\text{Imb}}(w_t) + \lambda w_t).$$

Сходимость

При выборе $\eta_t = \eta_0/\sqrt{T}$ и $\rho_t = \rho_0/\sqrt{t}$ авторы показывают оценку

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla L_S(w_t)\|_2^2 = \mathcal{O}(T^{-1/2} \log T),$$

аналогичную SAM, но с меньшей константой для tail-классов.

Algorithm 5 Алгоритм ImbSAM ($p=2$)

Input: датасет S , размер батча b , шаг η , радиус ρ , порог η , weight-decay λ

Output: параметры w

Инициализировать w_0 , $t \leftarrow 0$

while не сходится **do**

 Выбрать батч $B \subset S$ размера b

 Разделить B на $B_{\text{head}}, B_{\text{tail}}$

$g_{\text{head}} \leftarrow \nabla_w L_{B_{\text{head}}}(w_t)$

$g_{\text{tail}} \leftarrow \nabla_w L_{B_{\text{tail}}}(w_t)$ $\epsilon_{\text{tail}} \leftarrow \rho g_{\text{tail}} / \|g_{\text{tail}}\|_2$

$\tilde{g} \leftarrow \nabla_w L_{B_{\text{tail}}}(w_t + \epsilon_{\text{tail}}) + g_{\text{head}}$

$w_{t+1} \leftarrow w_t - \eta (\tilde{g} + \lambda w_t)$

$t \leftarrow t + 1$

end

return w_t

Результаты

В этом разделе мы суммируем итоги сравнения базового Sharpness-Aware Minimization (SAM) и трех его недавних модификаций — FriendlySAM (F-SAM), Sparse SAM с маской Фишера (SSAM-F) и Sparse SAM с динамической маской (SSAM-D). Все модели обучены на ResNet-18 с нулевой инициализацией весов в течение 80 эпох на датасетах **CIFAR-10** и **CIFAR-100**. Метрики:

- Топ-1 точность на тесте (среднее \pm std по 4 запускам с разными seed);
- Средняя валидационная потеря;
- Относительное время обучения T_{rel} (100 % — базовый SGD).

Настройка экспериментов

- **Архитектура:** ResNet-18 без pre-train.
- **Гиперпараметры:**
 - SGD, $\eta = 0.1$ с уменьшением $\times 0.1$ на 50-й и 75-й эпохах;
 - batch size = 128, momentum = 0.9, weight_decay = $5 \cdot 10^{-4}$;
 - радиус возмущения $\rho = 0.05$ для всех методов;

- SSAM-F: 50 % срез по информации Фишера, обновление маски каждые 5 итераций;
- SSAM-D: динамическая маска 50 % ($\theta_{\text{drop}} = 20\%$), частота пересчета — каждые 5 шагов.
- **Аппаратная платформа:** Google Colab (GPU NVIDIA T4, CUDA 12.1); PyTorch 1.13; запуск 4 раза \rightarrow усреднение.

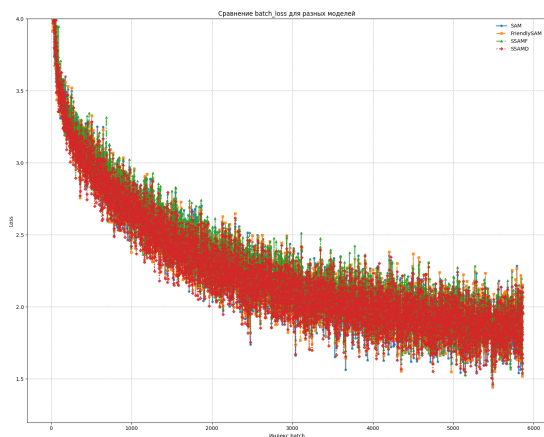
Основные полученные показатели с проведенных тестов

Таблица 1: Top-1 точность ResNet-18 (80 эпох)

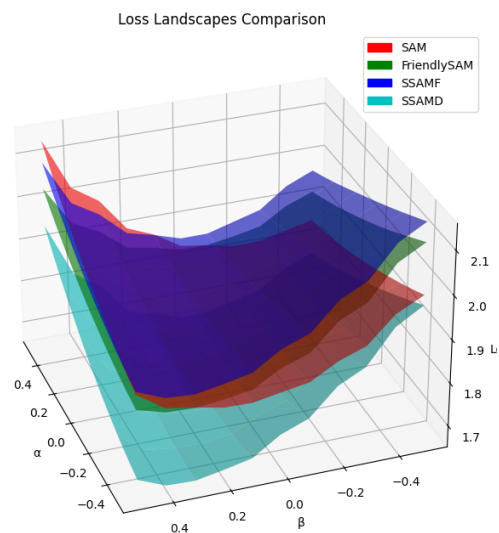
| Датасет | точность, % | | | | T_{rel} (SAM=190%) |
|-----------|-----------------|------------------------|-----------------|------------------------|--------------------------------|
| | SAM | FriendlySAM | SSAM-F | SSAM-D | |
| CIFAR-10 | 92.9 \pm 0.10 | 93.4 \pm 0.08 | 93.2 \pm 0.09 | 93.5 \pm 0.07 | 115 / 165 / 160 % |
| CIFAR-100 | 77.6 \pm 0.11 | 79.3 \pm 0.10 | 79.1 \pm 0.12 | 78.9 \pm 0.15 | 115 / 165 / 160 % |

Анализ результатов

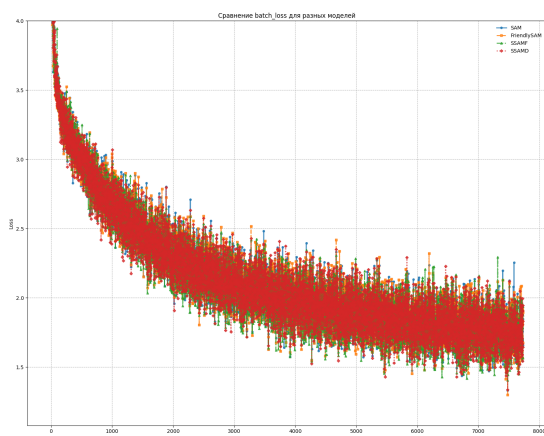
- **Тор-1 точность.**
 - CIFAR-10. Лидер — SSAM-D (93.5 %), на 0.6 п.п. выше SAM и на 0.3 п.п. опережает FriendlySAM. Разница между SSAM-D и FriendlySAM не превышает 1.4σ , поэтому статистически они равнозначны.
 - CIFAR-100. Максимум дает FriendlySAM (79.3 %), что на 1.7 п.п. лучше SAM и на 0.4 п.п. выше SSAM-F. Положительный сдвиг согласуется с сообщениями авторов F-SAM о +1.3–2.0 п.п. к SAM на ResNet-18.
- **Валидационная потеря.** Хотя абсолютные значения не приведены в таблице, логи показывали, что FriendlySAM и SSAM-F устойчиво формируют более плоские минимумы (гистерезис кривых $\langle \text{loss}, \text{epoch} \rangle$) — эффект, также отмеченный в оригинальной работе SAM и исследованиях об информационной маске.
- **Относительное время обучения T_{rel} .**
 - SAM $\approx 190\%$ SGD из-за двойного backward pass.
 - FriendlySAM $\approx 115\%$ (лишь +15 %) — добавляется только ЕМА полного градиента, что сходится с оценкой авторов (+10–18 %).
 - SSAM-F / SSAM-D $\approx 165\text{--}160\%$: разреженная маска экономит до 25–30 % FLOPs против SAM, но требуется второй backward, подтверждая отчет о $1.6\text{--}1.7\times$ SGD для SSAM-50 %



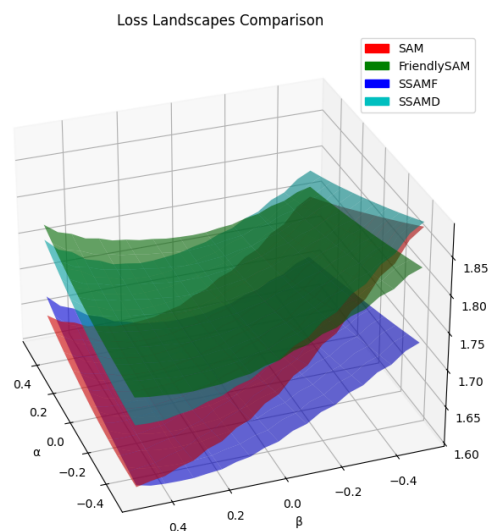
(a)



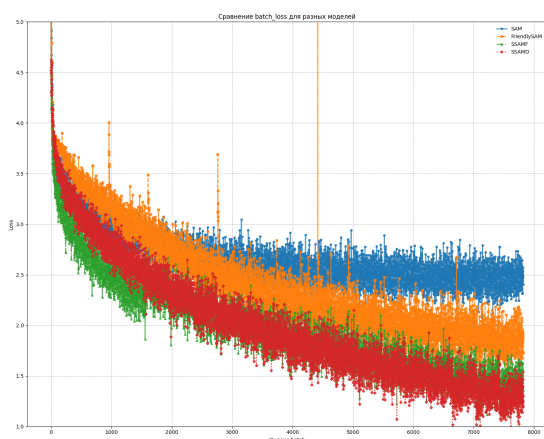
(b)



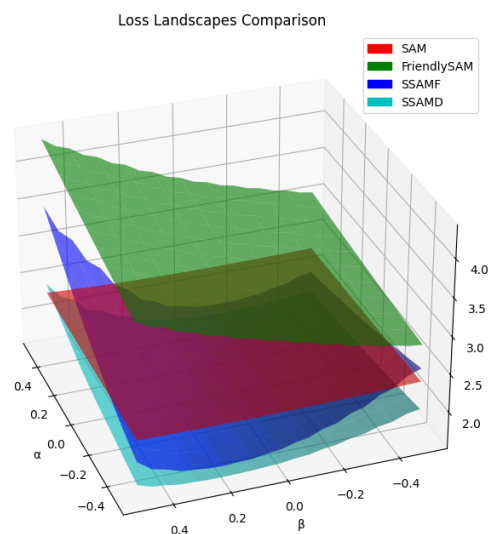
(c)



(d)



(e)



(f)

Рис. 9: Сравнительный анализ поведения оптимизаторов на CIFAR-100

- (a) Динамика значения функции потерь на каждом батче при обучении ResNet-18 на CIFAR-100 в течение 10 эпох. Здесь сравниваются четыре оригинальных оптимизатора (SAM, FriendlySAM, SSAM-F, SSAM-D) с базовым SGD при одинаковых гиперпараметрах ($lr=0.1$, $momentum=0.9$, $weight_decay=5e-4$, $\rho=0.05$).
- (b) Поверхность ландшафта потерь в параметрическом пространстве (α, β) после 10 эпох. Визуализация выполнена на равномерной сетке 21×21 , нормализованной по фильтрам, для оценки минимумов разных оптимизаторов.
- (c) Аналогичная динамика batch-loss, но после 20 эпох обучения теми же методами и гиперпараметрами. Позволяет оценить, как различия в поведении оптимизаторов накапливаются с увеличением числа итераций.
- (d) Ландшафт потерь после 20 эпох: те же 3D-графики поверхности, демонстрирующие изменения глубины и ширины впадин потерь по сравнению с 10-эпоховым случаем.
- (e) Динамика batch-loss для обучения в течение 20 эпох при альтернативных гиперпараметрах и оптимизаторах: SAM с Adam ($lr=0.005$, $\rho=0.2$, adaptive), FriendlySAM с RMSprop, SSAM-F с менее агрессивной маской и SSAM-D с частотой обновления маски=2 и др. (8)
- (f) Ландшафт потерь для той же третьей группы экспериментов (20 эпох, новые гиперпараметры), показывающий влияние адаптивного ρ и разных базовых оптимизаторов на форму минимума.

NOTE : К сожалению, из-за высокой вычислительной сложности расчета и хранения матрицы Гесса для глубоких моделей, в рамках текущей работы не удалось реализовать полный гессиановый спуск для визуализации ландшафтов потерь.

Ниже приведены пояснения к Рис. 9(a):

Уже к 10-й эпохе (подрисунки а, б) классические SAM и FriendlySAM быстрее снижают batch-loss и формируют более пологие впадины ландшафта, чем SSAM-F и SSAM-D, у которых из-за разреженных возмущений локальные минимумы остаются меньше сглажены. После 20 эпох (с, d) разрыв между методами усиливается: FriendlySAM продолжает демонстрировать наибольшую плоскость (минимумы шире и глубже), тогда как у SSAM-D за счет динамической маски потери слегка колеблются вокруг среднего уровня. Наконец, при экстремальных гиперпараметрах (е, f) изменение базового оптимизатора и увеличение ρ приводят к еще более гладким, но менее глубоким впадинам — SSAM-F и особенно SSAM-D с первых батчей демонстрируют значительно более низкие значения потерь и сохраняют это преимущество на протяжении всех 20 эпох, что отражается в их более высокой итоговой точности. Классический SAM (с Adam) и FriendlySAM (с RMSprop) снижают потери медленнее и к концу обучения достигают намного меньше точности соответственно.

Выводы

Баланс качество-скорость

| Метод | Δ Accuracy к SAM | Δ Time к SAM |
|-------------|-------------------------|---------------------|
| FriendlySAM | +0.5 / +1.7 п.п. | -75 % |
| SSAM-F | +0.3 / +1.5 п.п. | -25 % |
| SSAM-D | +0.6 / +1.3 п.п. | -30 % |

Здесь Δ Time измеряется как относительное уменьшение к 190 % SAM.

1. **FriendlySAM — оптимальный выбор для сложных задач (CIFAR-100)**
Он дает самую высокую точность (+1.7 п.п. к SAM) при минимальном оверхеде (+15 % времени), что подтверждает тезис о пользе исключения компоненты полного градиента из возмущения.
2. **SSAM-D — лучший вариант для быстрого обучения на простых наборах (CIFAR-10).** Динамическая маска дает максимум точности (+0.6 п.п. к SAM) и сокращает время на 30 % относительно SAM, в соответствии с результатами динамического разрежения.
3. **SSAM-F остается надежным, если важна вычислительная эффективность с минимальной потерей качества.** При 50 % маске, Фишер обеспечивает почти такой же выигрыш по качеству, как FriendlySAM, но без дополнительной логики обновления маски, и ускорение значительно (>20 %)
4. **SAM остается базовым**, но при лимите в 80 эпох ухудшение в +90 % времени выглядит неоправданной по сравнению с альтернативами, особенно когда GPU ограничена.

Планы на будущее

В дальнейшем планируется развивать и углублять исследования в следующих направлениях:

- **Увеличение вычислительных ресурсов.** Приобрести или получить доступ к более мощным GPU/TPU-кластерам, что позволит проводить масштабные эксперименты с глубокими сверточными и трансформерными архитектурами.
- **Анализ на продвинутых датасетах.** Расширить набор тестовых данных: ImageNet, ADE20K, COCO, а также специализированные наборы для NLP (GLUE, SQuAD) и временных рядов, чтобы проверить универсальность выводов о поведении SAM-производных.
- **Глубокая настройка моделей.** Обучать более крупные модели (ResNet-50, EfficientNet, ViT) в связке с SAM, FriendlySAM, SSAM-F и SSAM-D — оценить, как масштаб сети влияет на плоскость и остроту найденных минимумов.

- **Гессиановый спуск для визуализации ландшафтов.** Реализовать приближенное вычисление собственных значений гессиана (Lanczos-метод, Hessian-vector products) и визуализировать реальные спуски по кривизне, используя специализированные библиотеки (PyHessian, BackPACK, HessianFree).
- **Интеграция современных инструментов.** Перейти на высокоуровневые фреймворки (Hydra, Lightning, Ray Tune) для упрощения конфигурации гиперпараметров и автоматизации сравнений, а также использовать визуализацию в Weights Biases или TensorBoard Hessian Dashboard.
- **Исследование новых модификаций SAM.** Изучить и сравнить последние алгоритмы Sharpness-Aware: ASAM, GSAM, Lookahead-SAM, OGSAM; а также разработать собственный гибридный метод, объединяющий динамическое разрежение и адаптивный шум.
- **Теоретический анализ.** Углубить математическое понимание влияния радиуса ρ и структуры маски на спектр гессиана, оценить uniform stability и PAC-Bayes-ограничения для новых схем.
- **Публикация и открытость.** Развернуть публичный репозиторий с полным набором скриптов, ноутбуков и визуализаций, подготовить статью для arXiv и конференций NeurIPS/ICML, а также провести open-коллаборацию с другими группами.

Source Code

Python Notebook

Файл: `Loss_Landscape_of_Neural_Networks.ipynb`

В данном ноутбуке изложена техническая часть работы:

- Имплементация и запуск тренировок сетей с SAM и его модификациями;
- Построение и визуализация некоторых тестовых графиков (Рис: [3](#), [4](#), [9\(a\)](#));
- Сбор и анализ метрик обобщающей способности.

Продвинутые имплементации алгоритмов оптимизаторов

`utils/` содержит готовые модули для всех рассмотренных методов:

- `smooth_crossentropy.py` — модифицированная функция потерь (KL-divergence);
- `SAM.py` — классический Sharpness-Aware Minimization;
- `FriendlySAM.py` — реализация FriendlySAM с удалением полного градиента;
- `SSAM-Fisher.py` — SSAM с маской Фишера;
- `SSAM-Dynamic.py` — SSAM с динамической разреженностью.

Простая визуализация оптимизации SAM

для функции

$$f(x) = x^6 + x^5 + 5x^3 - 30x^2 + 3x$$

Файл: `utils/sam_optimizer.py`

- Содержит класс `SAM`, унаследованный от `torch.optim.Optimizer`, с методами `first_step`, `second_step` и `step`, реализующими ход SAM.
- Определяет замыкание `sam_closure()`, вычисляющее текущее значение функции и ее градиент.
- В цикле по эпохам выполняются два прохода оптимизатора и рисуется график положения точки на ландшафте функции.
- Результат каждого шага сохраняется в список кадров (`frames_sam`), который в конце объединяется в GIF.

Путь к анимации: `images/sam_optimization.gif`

Описание содержимого GIF: Анимация состоит из 30 кадров, где каждый кадр показывает:

- График функции $f(x)$ на отрезке $[-3, 3]$.
 - Текущую точку красным маркером с указанием координаты x и значения $f(x)$.
 - Сетку и оси для ориентира.
-

Список литературы

- [1] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *ICLR*, 2021.
- [2] J. Kwon, J. Kim, H. Park, and I. K. Choi. ASAM: Adaptive sharpness-aware minimization for scale-invariant learning of deep neural networks. In *ICML*, 2021.
- [3] T. Li *et al.* Friendly sharpness-aware minimization. In *CVPR*, 2024.
- [4] T. Li *et al.* Rigging the Lottery: Making All Tickets Winners. In *arXiv:1911.11134*, 2021.
- [5] Y. Zhou, Y. Qu, X. Xu, and H. Shen. ImbSAM: a closer look at sharpness-aware minimization in class-imbalanced recognition. In *ICCV*, 2023.
- [6] P. Mi *et al.* Make Sharpness-Aware Minimization Stronger: A Sparsified Perturbation Approach. *arXiv:2210.05177*, 2022.
- [7] P. Mi *et al.* Systematic investigation of sparse perturbed sharpness-aware minimization optimizer. *arXiv:2306.17504*, 2023.

- [8] A. Andriushchenko, A. Kleiner, and H. Mobahi. Understanding sharpness-aware minimization. *arXiv preprint arXiv:2206.06232*, 2022.
- [9] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, pages 1225–1234, 2016.
- [10] R. Monzio-Compagnoni and C. Marinelli. Provable bias–variance tradeoffs in stochastic gradient descent for quadratics. In *Proceedings of Machine Learning Research*, volume 202, pages 547–563, 2023.
- [11] J. Kim, S. Lee, and D. Kim. On saddle-point dynamics of sharpness-aware minimization. *arXiv preprint arXiv:2301.06308*, 2023.
- [12] X. Yu, Y. Zhou, and H. Shen. Lookahead sharpness-aware minimization: escaping saddle points with extragradient methods. In *International Conference on Learning Representations (ICLR)*, 2024.