

Глубокое обучение

Дмитрий Никулин

14 апреля 2021 г.

Лекция 4: 50 оттенков градиентного спуска

Agenda

- Разбираем способы улучшить SGD
- Экспериментируем с ними

50 оттенков градиентного спуска



Градиентный спуск (GD)

Проблема оптимизации:

$$L(w) = \frac{1}{n} \cdot \sum_{i=1}^n L(w, x_i, y_i) \rightarrow \min_w \quad \text{по датасету } \{(x_i, y_i)\}_{i=1}^n$$

Градиент указывает направление максимального роста

$$\nabla L(w) = \left(\frac{\partial L(w)}{\partial w_0}, \frac{\partial L(w)}{\partial w_2}, \dots, \frac{\partial L(w)}{\partial w_k} \right)$$

Идём в противоположную сторону:

$$w^1 = w^0 - \lambda \cdot \nabla L(w^0)$$

скорость обучения

Градиентный спуск (GD)

Проблема оптимизации:

$$L(w) = \sum_{i=1}^n L(w, x_i, y_i) \rightarrow \min_w \quad \text{по датасету } \{(x_i, y_i)\}_{i=1}^n$$

Инициализация w_0

while True:

$$g_t = \frac{1}{n} \sum_{i=1}^n \nabla L(w, x_i, y_i)$$

$$w_t = w_{t-1} - \lambda_t \cdot g_t$$

if $\|w_t - w_{t-1}\| < \varepsilon$:

break

(можно останавливаться и при превышении количества итераций / по метрикам на валидации / etc)

Стохастический градиентный спуск (SGD)

Проблема оптимизации:

$$L(w) = \sum_{i=1}^n L(w, x_i, y_i) \rightarrow \min_w \quad \text{по датасету } \{(x_i, y_i)\}_{i=1}^n$$

Инициализация w_0

while True:

 рандомно выбрали $b < n$ индексов $B = \{i_1, \dots, i_b\}$

$$g_t = \frac{1}{b} \sum_{j \in B} \nabla L(w, x_j, y_j)$$

$$w_t = w_{t-1} - \lambda_t \cdot g_t$$

if $\|w_t - w_{t-1}\| < \varepsilon$:

break

Momentum SGD

Мы считали на каждом шаге градиент по формуле

$$g_t = \frac{1}{b} \sum_{i=1}^b \nabla L(w_{t-1}, x_i, y_i).$$

После шага мы забывали его. **Давайте запоминать направление:**

$$m_t = \alpha \cdot m_{t-1} + g_t$$

$$w_t = w_{t-1} - \lambda \cdot m_t$$

- Движение поддерживается в том же направлении, что и на предыдущем шаге
- Нет резких изменений направления движения.
- Обычно $\alpha \approx 0.9$.

Крутой интерактив для моментума: <https://distill.pub/2017/momentum/>

Momentum SGD

- Бежим с горки и всё больше ускоряемся в том направлении, в котором были направлены сразу несколько предыдущих градиентов, но при этом движемся медленно там, где градиент постоянно меняется
- Хотелось бы не просто бежать с горы, но и хотя бы на полшага смотреть себе под ноги, чтобы внезапно не споткнуться \Rightarrow давайте смотреть на градиент в будущей точке
- В Momentum SGD $\alpha \cdot m_{t-1}$ точно будет использоваться при шаге, давайте искать $\nabla L(w_{t-1} - \lambda \cdot \alpha \cdot m_{t-1})$.

Nesterov Momentum SGD

- Мы теперь сначала прыгаем в том же направлении, в каком шли до этого, потом корректируем его
- Иллюстрация из статьи Хинтона 2013 года:

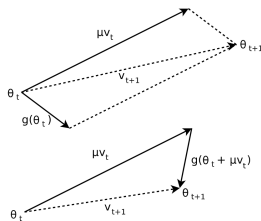


Figure 1. **(Top)** Classical Momentum **(Bottom)** Nesterov Accelerated Gradient

<http://www.cs.toronto.edu/~hinton/absps/momentum.pdf>

Nesterov Momentum SGD

- Мы теперь сначала прыгаем в том же направлении, в каком шли до этого, потом корректируем его

$$m_t = \alpha \cdot m_{t-1} + \nabla L(w_{t-1} - \alpha \cdot m_{t-1})$$

$$w_t = w_{t-1} - \lambda \cdot m_t$$

Разная скорость обучения

- Может сложиться, что некоторые веса уже близки к своим локальным минимумам, по этим координатам надо двигаться медленнее, а по другим быстрее \Rightarrow **адаптивные методы градиентного спуска**
- Шаг изменения должен быть меньше у тех параметров, которые в большей степени варьируются в данных, и больше у тех, которые менее изменчивы

AdaGrad

$$v_0^j = 0$$

$$v_t^j = v_{t-1}^j + g_{tj}^2$$

$$w_t^j = w_{t-1}^j - \frac{\lambda}{\sqrt{v_t^j + \varepsilon}} \cdot g_{tj}$$

- g_{tj} — градиент по j -ому параметру
- своя скорость обучения для каждого параметра
- обычно $\lambda = 0.01$, т.к. параметр не очень важен
- v_t^j всегда увеличивается, из-за этого обучение может рано останавливаться \Rightarrow RMSprop

RMSprop

$$v_0^j = 0$$

$$v_t^j = \alpha \cdot v_{t-1}^j + (1 - \alpha) \cdot g_{tj}^2$$

$$w_t^j = w_{t-1}^j - \frac{\lambda_t}{\sqrt{v_t^j + \varepsilon}} \cdot g_{tj}$$

- Обычно $\alpha = 0.99$
- Скорость обучения адаптируется к последнему сделанному шагу, бесконтрольного роста v_t^j больше не происходит
- RMSprop нигде не был опубликован, Хинтон просто привёл его в своей лекции, сказав, что это норм тема

Adam (Adaptive Moment Estimation)

$$m_t^j = \beta_1 \cdot m_{t-1}^j + (1 - \beta_1) \cdot g_{tj}$$

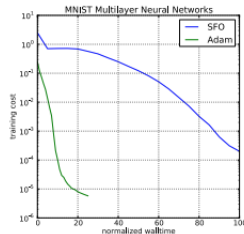
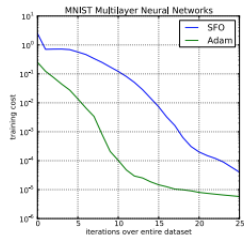
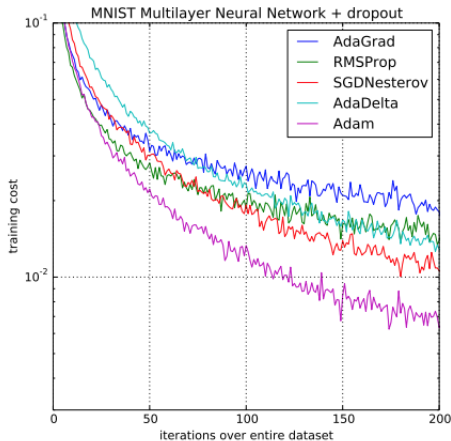
$$v_t^j = \beta_2 \cdot v_{t-1}^j + (1 - \beta_2) \cdot g_{tj}^2$$

$$\hat{m}_t^j = \frac{m_t^j}{1 - \beta_1^t} \quad \hat{v}_t^j = \frac{v_t^j}{1 - \beta_2^t}$$

$$w_t^j = w_{t-1}^j - \frac{\lambda_t}{\sqrt{\hat{v}_t^j} + \varepsilon} \cdot \hat{m}_t^j$$

- Комбинируем Momentum и индивидуальные скорости обучения
- Фактически, \hat{m}_t и \hat{v}_t — это оценки первого и второго моментов для стохастического градиента
- Деление на $1 - \beta_i^t$ нужно, чтобы оценки были несмещёнными

Сравнение на MNIST (из статьи про Adam)



<https://arxiv.org/pdf/1412.6980.pdf>

Лучший оптимизатор?..

STOCHASTIC GRADIENT DESCENT	
SGD WITH MOMENTUM	
RMSPROP	
ADAM	

imgflip.com

<https://towardsdatascience.com/optimization-algorithms-in-deep-learning-191bfc2737a4>

<https://arxiv.org/pdf/2007.01547.pdf>

A. List of optimizers and schedules considered

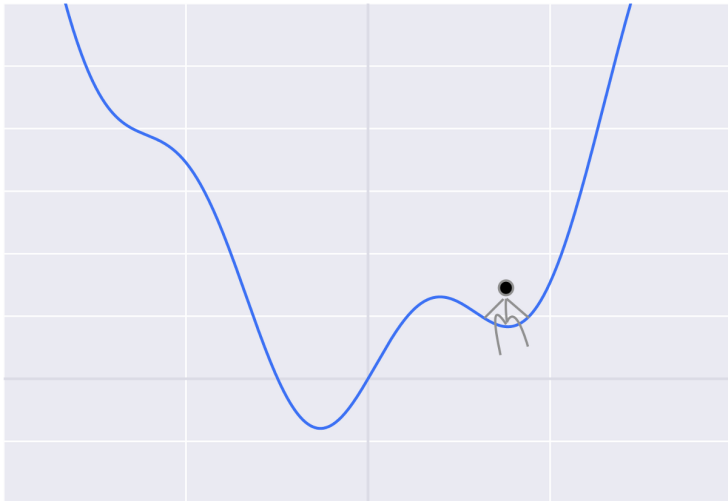
Table 2: List of optimizers we considered for our benchmark. Note, that this is far from being a complete list of all existing optimization methods applicable to deep learning, but only a subset, comprising of some of the most popular choices.

[illegible]

Резюме по методам градиентного спуска

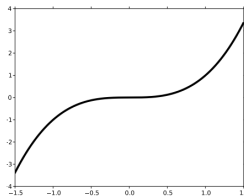
- Momentum SGD сохраняет направление шага и позволяет добиваться более быстрой сходимости
- Адаптивные методы позволяют находить индивидуальную скорость обучения для каждого параметра
- Adam комбинирует в себе оба подхода
- Давайте посмотрим [визуализацию 1](#) и [визуализацию 2](#)
- Но это же не все вызовы!

Боб чилит в локальном минимуме

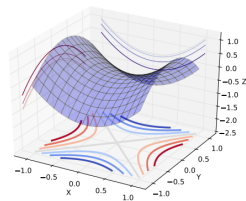


<https://hackernoon.com/life-is-gradient-descent-880c60ac1be8>

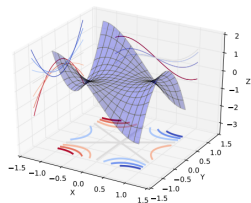
Седловые точки



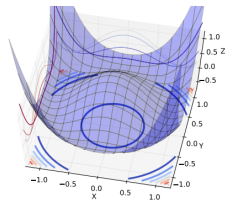
(a)



(b)



(c)

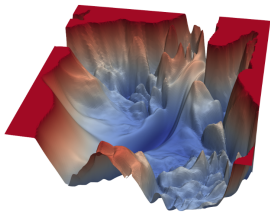


(d)

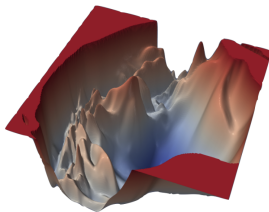
<https://arxiv.org/pdf/1406.2572.pdf>

Визуализация потерь

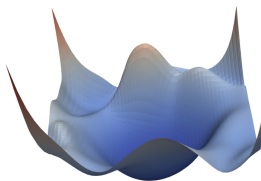
VGG-56



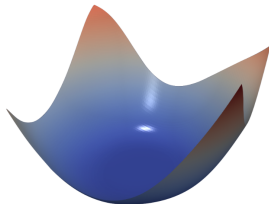
VGG-110



Renset-56



Densenet-121



<https://arxiv.org/pdf/1712.09913.pdf>

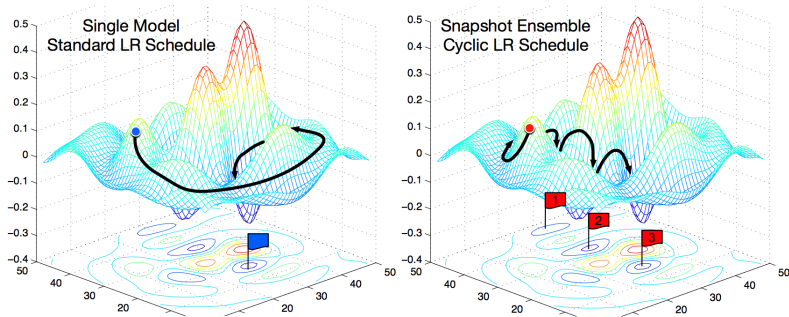
<https://github.com/tomgoldstein/loss-landscape>

Переменный learning rate

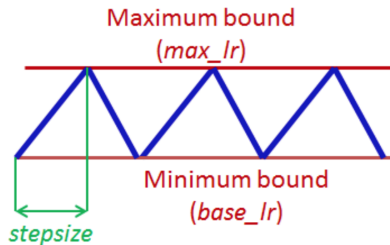
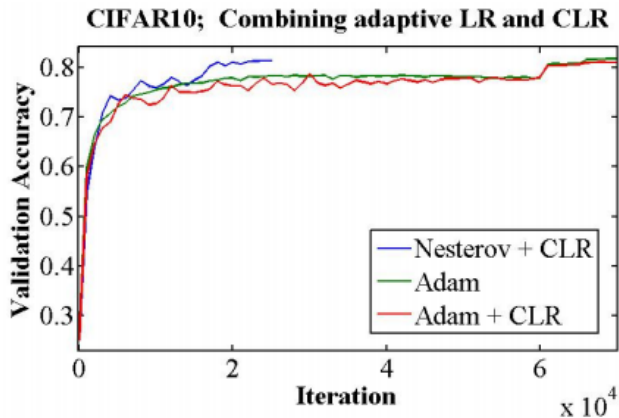
- Когда качество на валидации перестаёт улучшаться, уменьшаем learning rate в несколько раз

Циклическая скорость обучения (CLR)

- Хочется, чтобы был шанс вылезти из локального минимума, а также шанс сползти с седла \Rightarrow давайте менять глобальную скорость обучения циклически



Циклическая скорость обучения (CLR)



Нестеров с CLR отработал быстрее и лучше Adam

Нет одного правильного алгоритма на все случаи!

Всегда надо экспериментировать

<https://arxiv.org/pdf/1506.01186.pdf>

<https://openreview.net/pdf?id=BJYwwY9ll>

А что люди сейчас делают с оптимизаторами?

- Нейросеть пытается сама от слоя к слою выстраивать новые, более сложные фичи
- Самостоятельно фичи больше придумывать не надо
- При этом оптимизаторы какие-то слишком олдскульные и ручные
- Круто было бы, если бы оптимизаторы тоже сами обучались
- Такие работы в последнее время активно появляются, но какого-то суперпрогресса пока нет

<https://arxiv.org/pdf/2009.11243.pdf>

Экспериментируем!

