

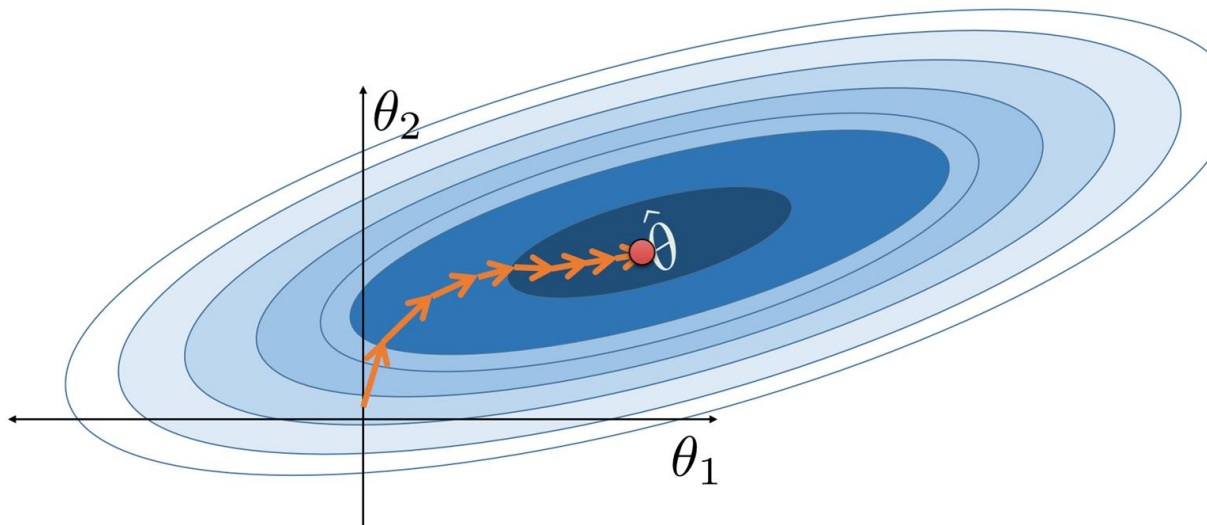
Оптимизаторы

Тетерин Михаил Александрович

ОБРАЗОВАТЕЛЬНЫЕ ПРОГРАММЫ В ОБЛАСТИ
ТЕХНОЛОГИЙ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА

- легко реализуется и приспособлен для динамического обучения, когда объекты поступают потоком
- способен обучаться на больших выборках
- может не сходиться
- может застрять в одном из локальных минимумов
- усреднение по маленьким батчам может приводить к шумным

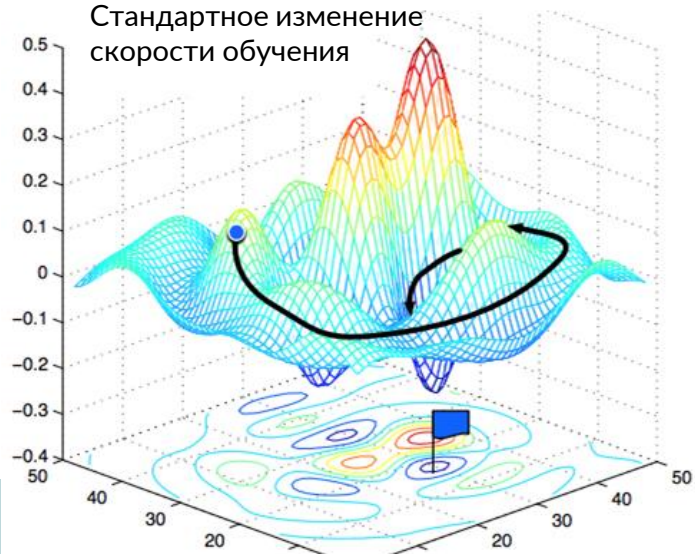
$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$



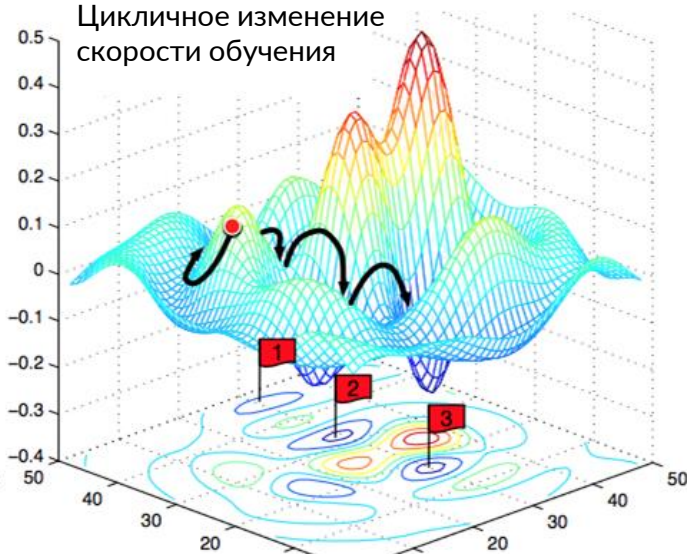
- легко реализуется и приспособлен для динамического обучения, когда объекты поступают потоком
- способен обучаться на больших выборках
- **может не сходиться**
- **может застрять в одном из локальных минимумов**
- **усреднение по маленьким батчам может приводить к шумным**

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

Стандартное изменение скорости обучения

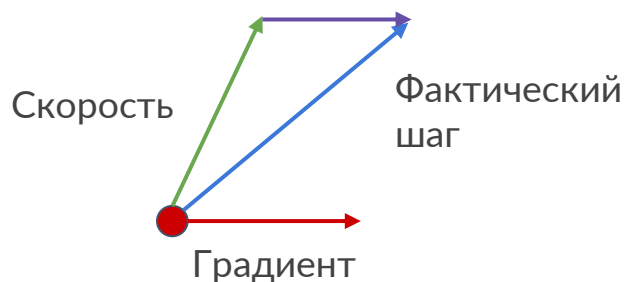


Циклическое изменение скорости обучения

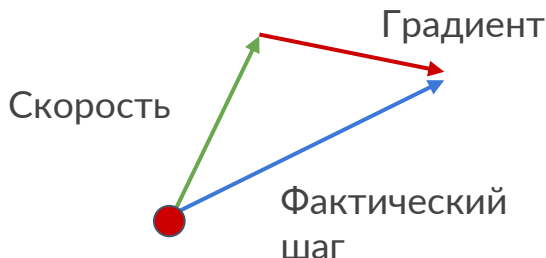


- легко реализуется и приспособлен для динамического обучения, когда объекты поступают потоком
- способен обучаться на больших выборках
- выше вероятность что сойдется
- **может не сходиться**
- **может застрять в одном из локальных минимумов (седловые точки)**

Обновление импульса:



Импульс Нестерова:



$$v_{t+1} = \rho v_t + \nabla f(x_t) \quad v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1} \quad x_{t+1} = x_t + v_{t+1}$$

- стохастический градиент спуск + кеширование
- уменьшаем обновления для элементов, которые мы часто используем
- скорость обучения уменьшается слишком быстро
- глобальную скорость обучения надо подбирать, и она может быть хороша для одних размерностей, но плоха для других

$$cache_{t+1} = cache_t + (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{cache_{t+1} + \epsilon}$$

- стохастический градиент спуск + кеширование + экспоненциальное сглаживание
- уменьшаем обновления для элементов, которые мы часто используем
- не нужно выбирать скорость обучения
- нет паралича алгоритма
- высокая вероятность преодоления локального оптимума

$$cache_{t+1} = \beta cache_t + (1 - \beta)(\nabla f(x_t))^2$$

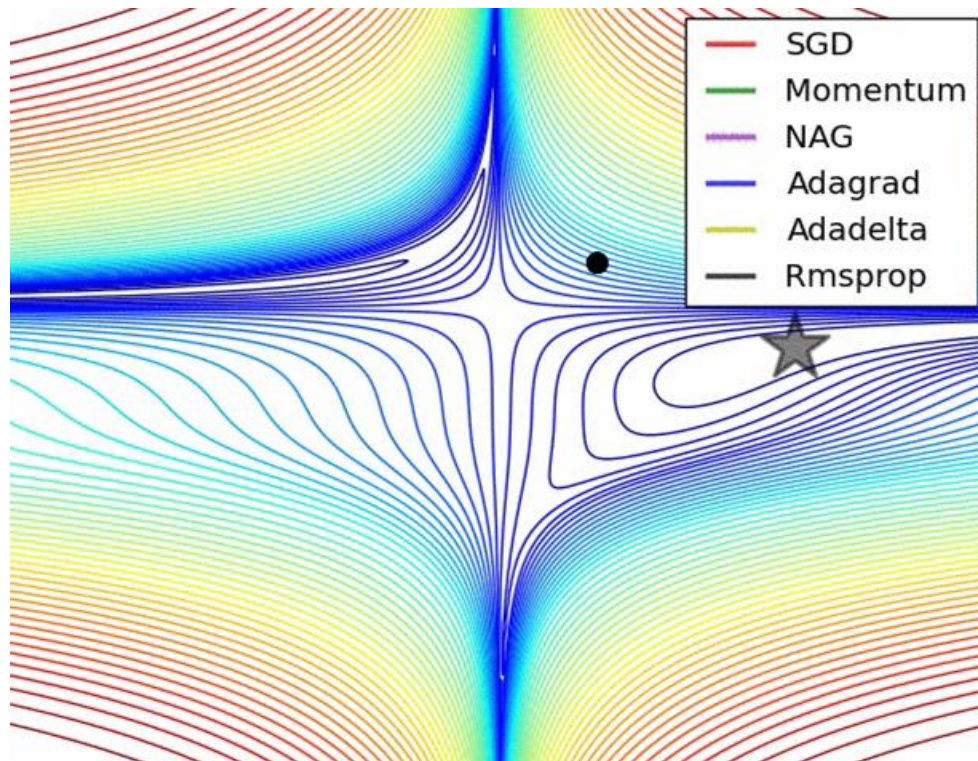
$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{cache_{t+1} + \epsilon}$$

- такие же преимущества как и у RMSProp
- быстрее сходится

$$\nu_{t+1} = \gamma \nu_t + (1 - \gamma) \nabla f(x_t)$$

$$cache_{t+1} = \beta cache_t + (1 - \beta) (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nu_{t+1}}{cache_{t+1} + \epsilon}$$



Спасибо за внимание

ОБРАЗОВАТЕЛЬНЫЕ ПРОГРАММЫ В ОБЛАСТИ
ТЕХНОЛОГИЙ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА