

Оптимизаторы - алгоритмы, используемые для поиска минимума пар-тов с целью ускорения работы модели.

AdaGrad - оптимизатор, использующий скорости обучения  $\eta$  для каждого пар-та на каждом шаге. Работает по правилу ф-ции ошибки.

$$g_t = \nabla_{\theta} J(\theta_t) \text{ - правило ф-ции}$$

$$G_t = G_t + g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot g_t \text{ - обновление пар-тов}$$

$\eta$  - скорости обучения, как итерируется для зад. пар-ра  $\theta$  в данный момент времени на основе предыдущих градиентов, рассчит. для данного пар-та.

Adadelta - расширение AdaGrad. Огранич. окно памяти. прошлых  $\eta$  градиентов до нек. фикс. размера, вместо того, чтобы



хранить их все. Испытуется эксп.

скользящая средняя, а не сумма градиентов.

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1-\gamma) g_t^2$$

$$RMS[g]_t = \sqrt{E[g^2]_t} \cdot \epsilon$$

$$\theta_{t+1} = \theta_t - \frac{RMS[\Delta \theta]_{t-1}}{RMS[g]_t} \cdot g_t$$

Adam - работает с импульсами 1 и 2 порядков. Уже замечается в уменьшении скорости во избежание проскакивания минимума. В дополнение к квадратичной экспоненц. скользящей средним (как в Adelta) сокращ. экспоненц. скользящую среднюю.

$$m_t^1 = \frac{m_t^1}{1 - \beta_1^t} \quad v_t^1 = \frac{v_t^1}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t^1} \cdot \epsilon} \cdot m_t^1$$



Название  
Gradient  
Descent

Формула  
$$\theta = \theta - \alpha \nabla J(\theta)$$

Характеристики  
Можно использовать для всех данных сразу.  
Величина шага должна быть выбрана для всего набора данных.  
Требуется большой объем памяти.

Stochastic  
Gradient Descent

Формула  
$$\theta = \theta - \alpha \nabla J(\theta; x_i, y_i)$$
  
 $x_i, y_i$  - одна пара

Характеристики  
Высокая дисперсия параметров.  
Чтобы получить хороший результат, необходимо многократно снижать значение скорости обучения.

Mini-Batch  
Gradient Descent

Формула  
$$\theta = \theta - \alpha \nabla J(\theta; B_j)$$
  
 $B_j$  - batch одна пар.

Характеристики  
Можно использовать для всех данных сразу.  
При выборе скорости обучения, необходимо многократно снижать значение скорости обучения.

Adagrad

Формула  
$$g_t = \nabla J(\theta_t)$$
  
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t}} \cdot g_t$$

Характеристики  
Вычисл. дорого (очень много операций на пар-х)  
Медленное обучение

Adadelta

Формула  
$$RMS[g]_t = \sqrt{E[g]_t^2}$$
  
$$\theta_{t+1} = \theta_t - \frac{RMS[g]_t}{RMS[g]_t} \cdot g_t$$

Характеристики  
Вычисл. дорого

RMSProp

Формула  
$$\theta_{t+1} = \theta_t - \frac{\eta}{RMS[g]_t} \cdot g_t$$

Характеристики  
Выс. дорого

Adam

Формула  
$$m_t = \frac{\eta}{1 - \beta_1^t} \cdot \hat{g}_t = \frac{\eta}{1 - \beta_1^t} \cdot \frac{1}{\sqrt{1 - \beta_2^t}} \cdot \hat{g}_t$$
  
$$\theta_{t+1} = \theta_t + m_t$$

Характеристики  
Выс. дорого

Nesterov  
Accelerated  
Gradient

Формула  
$$\tilde{\theta} = \theta - \gamma \nabla J(\theta - \gamma \nabla J(\theta))$$
  
$$\theta \leftarrow \tilde{\theta}$$

Характеристики  
Нужна скорость обучения