



ТЕХНОСФЕРА

Лекция 7 Методы оптимизации, часть 2

Полыковский Даниил

20 марта 2017 г.

Постановка задачи

- ▶ $\theta_* = \min_{\theta} J(\theta)$
- ▶ В любой точке можем вычислить $\nabla_{\theta} J(\theta)$

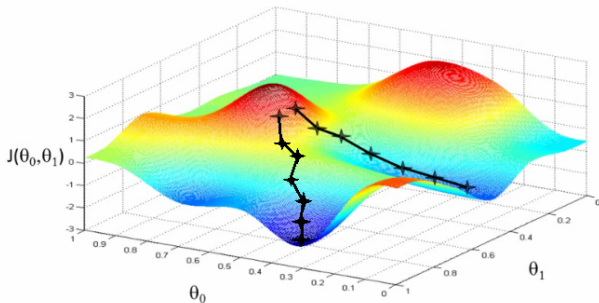


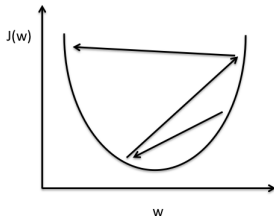
Рис.: Пример функции для оптимизации

Batch Gradient Descend

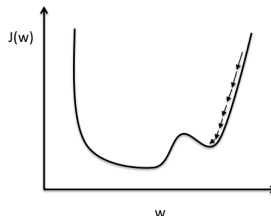
Формула пересчета:

$$\theta_t = \theta_{t-1} - \eta \nabla_{\theta} J(\theta_{t-1})$$

- Требуется обработать все объекты для одного шага
- Нет режима online обучения
- + Гарантируется сходимость к (локальному) минимуму



Large learning rate: Overshooting.



Small learning rate: Many iterations until convergence and trapping in local minima.

Рис.: Выбор темпа обучения

SGD / Mini-batch SGD

- ▶ Какие функции оптимизируем?

SGD / Mini-batch SGD

- ▶ Какие функции оптимизируем?
- ▶ Большие суммы функций: $J(\theta) = \sum_{i=1}^N J_i(\theta)$
- ▶ Формула пересчета: $\theta_t = \theta_{t-1} - \eta \nabla_{\theta} J_i(\theta_{t-1})$
- ▶ Mini-batch SGD: $\theta_t = \theta_{t-1} - \eta \sum_{i \in \{i_1, i_2, \dots, i_k\}} \nabla_{\theta} J_i(\theta_{t-1})$
- ▶ Легко попасть в регион согласованности, тяжело найти общий оптимум

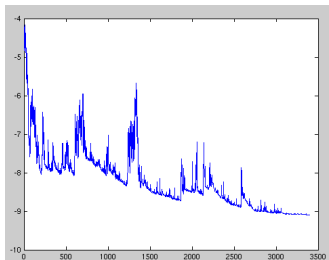


Рис.: Изменение значения J во время обучения

Momentum

- ▶ $\nu_t = \gamma \nu_{t-1} + \eta \nabla_{\theta} J(\theta) \leftarrow$ “скорость”
- ▶ $\theta = \theta - \nu_t$
- ▶ Рекомендовано брать $\gamma = 0.9$
- ▶ Проблема: метод приводит к перескокам через локальный минимум

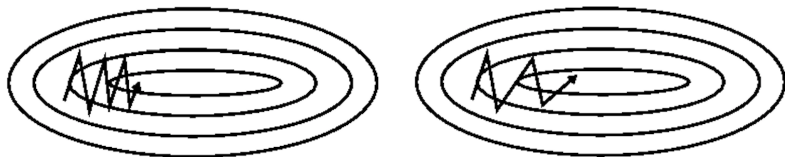


Рис.: Слева: без импульса, справа: с импульсом

Nesterov accelerated gradient

- ▶ Следующая позиция приближенно равна $\theta - \gamma\nu_{t-1}$
- ▶ Вычисление градиента дает возможность узнать будущее направление градиента
- ▶ $\nu_t = \gamma\nu_{t-1} + \eta\nabla_{\theta}J(\theta - \gamma\nu_{t-1})$
- ▶ $\theta = \theta - \nu_t$

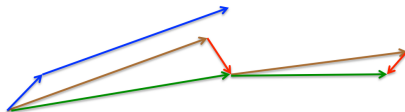


Рис.: NAG

- ▶ Сначала делаем шаг в направлении накопленного градиента
- ▶ Затем вычисляем градиент там и делаем поправку

Методы

- ▶ SGD $\nu_t = \eta \nabla_{\theta} J_i(\theta_{t-1})$
Momentum $\nu_t = \gamma \nu_{t-1} + \eta \nabla_{\theta} J(\theta)$
NAG $\nu_t = \gamma \nu_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma \nu_{t-1})$
- ▶ $\theta = \theta - \nu_t$
- ▶ Общая проблема: одинаковый шаг для всех параметров
- ▶ Трудно подобрать η или η_t
- ▶ Примеры расписаний: $\eta_t = \gamma^t \eta$, $\eta_t = \begin{cases} \alpha_1 & t \leq A \\ \alpha_2 & t > A \end{cases}$

Adagrad

- ▶ $\mathbf{g}_{t,i} = \nabla_{\theta} J(\theta_i)$
- ▶ $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot \mathbf{g}_{t,i}$
- ▶ $G_{t,ii}$ – сумма квадратов значений $\mathbf{g}_{t,i}$ вплоть до текущего
- ▶ Векторно: $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot \mathbf{g}_t$
- ▶ Стандартные значения: $\eta = 0.01$, $\epsilon = 10^{-8}$
- ▶ Мотивация: маленькие обновления для часто встречающихся параметров, большие для редких
- ? Какова проблема этого метода?

Adagrad

- ▶ $\mathbf{g}_{t,i} = \nabla_{\theta} J(\theta_i)$
- ▶ $\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot \mathbf{g}_{t,i}$
- ▶ $G_{t,ii}$ – сумма квадратов значений $\mathbf{g}_{t,i}$ вплоть до текущего
- ▶ Векторно: $\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot \mathbf{g}_t$
- ▶ Стандартные значения: $\eta = 0.01$, $\epsilon = 10^{-8}$
- ▶ Мотивация: маленькие обновления для часто встречающихся параметров, большие для редких
- ? Какова проблема этого метода?
 $G_{t,ii}$ неубывает \Rightarrow затухание обновлений

RMSProp / Adadelta

- ▶ Будем использовать последние несколько значений g_t^2 для подсчета G_t
- ▶ Экспоненциальное среднее: $E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2$,
 $\gamma = 0.9$
- ▶ $\theta_t = \theta_{t-1} - \Delta\theta_t$
- ▶ $\Delta\theta_t = \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t = \frac{\eta}{RMS[g]_t} g_t \leftarrow \text{RMSprop}$
- ▶ Adadelta: избавимся от η
- ▶ $\Delta\theta_t = \frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$

Adam

- ▶
$$\begin{cases} m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ \nu_t = \beta_2 \nu_{t-1} + (1 - \beta_2) g_t^2 \end{cases}$$
- ▶ m_t, ν_t инициализируются нулями, поэтому долгий “разгон” \Rightarrow надо уменьшить инерцию в начале обучения
- ▶
$$\begin{cases} \hat{m}_t = \frac{m_t}{1 - \beta_1^t} \\ \hat{\nu}_t = \frac{\nu_t}{1 - \beta_2^t} \end{cases}$$
- ▶
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{\nu}_t} + \epsilon} \hat{m}_t$$
- ▶ $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

Критерий остановки

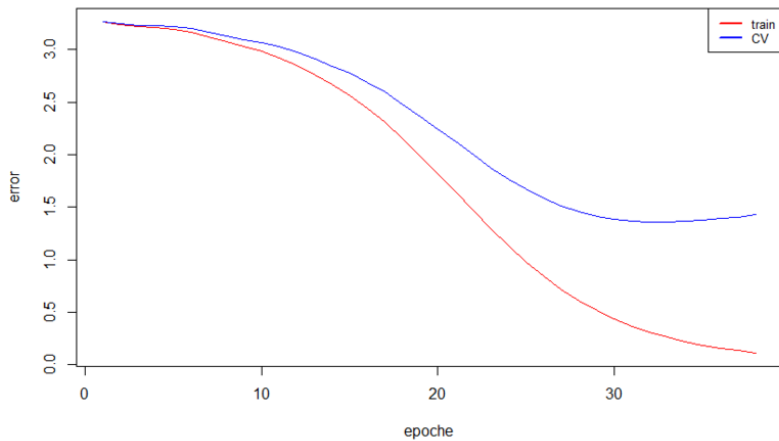


Рис.: Кроссвалидация

Визуализация

- ▶ 2D визуализация (gif)
- ▶ Седловая точка (gif)

Вопросы

