

Momentum: 结合当前梯度与上一次更新信息

$$\begin{aligned} \text{指数加权平均: } v_t &= \beta \cdot v_{t-1} + (1-\beta) \cdot \theta_t \\ &\quad \uparrow \quad \quad \uparrow \quad \quad \uparrow \\ &\quad \text{权重} \quad \text{之前平均值} \quad \text{当前参数值} \\ &= \sum_{i=1}^t (1-\beta) \cdot \beta^{t-i} \cdot \theta_{t-i} \end{aligned}$$

β 表示对 $\frac{1}{1-\beta}$ 次数据的记忆

不考虑 Momentum(m) 时的梯度下降

$$w_{i+1} = w_i - lr \cdot g(w_i)$$

考虑 m 时 (振荡收敛)

$$v_i = m \cdot v_{i-1} + g(w_i)$$

v_i : 更新量

$$w_{i+1} = w_i - lr \cdot v_i$$

optim. SGD

params: 参数组

lr: 初始学习率

momentum: 还是

weight_decay: L2 正则化系数

nesterov: 是否采用 NAG

其它优化器

optim. Adagrad 自适应学习率梯度下降法

RMSprop: Adagrad 的改进

Adadelta: Adagrad 的改进

Adam: RMSprop 结合 Momentum

Adamax: Adam 增加学习率上限.

SparseAdam: 稀疏版Adam

ASGD: 随机平均梯度下降.

Rprop: 弹性反向传播

LBFGS: BFGS 的改进