Author: Liu Jian

Time: 2021-03-17

1. SGD

- 1.1 动量法
 - 1.1.1 Momentum
 - 1.1.2 Nesterov Momentum
- 1.2 权重自适应 (Adaptive Learning Rates)
 - 1.2.1 AdaGrad
 - 1.2.2 RMSProp
- 1.3 动量法 + 权重自适应
 - 1.3.1 RMSProp with Nesterov momentum
 - 1.3.2 Adam

2. 二阶优化算法

2.1 牛顿法

1. SGD

记:

$$g(heta) riangleq
abla_{ heta} \left(rac{1}{m} \sum_{i=1}^{m} \mathcal{L}(f(x^{(i)}; heta), y^{(i)})
ight)$$

随机梯度下降 ($t = 1, 2, \dots,$ 下同):

$$heta^{t+1} = heta^t - arepsilon g(heta^t)$$

1.1 动量法

1.1.1 Momentum

Require: Learning rate ε , momentum parameter α

Require: Initial parameter θ^1 , initial velocity v^0

$$v^t = \alpha v^{t-1} - \varepsilon g(\theta^t)$$

 $\theta^{t+1} = \theta^t + v^t$

Momentum 的另一种表达形式:

$$egin{aligned} v^t &= lpha v^{t-1} - arepsilon g(heta^t) = arepsilon \left(rac{lpha}{arepsilon} v^{t-1} - g(heta^t)
ight) riangleq arepsilon u^t \ & heta^{t+1} = heta^t + arepsilon u^t \end{aligned}$$

接着寻找 u^t 和 u^{t-1} 之间的递推关系:

$$\begin{split} u^t &= \frac{\alpha}{\varepsilon} v^{t-1} - g(\theta^t) \\ v^t &= \varepsilon u^t \Rightarrow v^{t-1} = \varepsilon u^{t-1} \\ & \qquad \qquad \downarrow \\ u^t &= \frac{\alpha}{\varepsilon} \varepsilon u^{t-1} - g(\theta^t) = \alpha u^{t-1} - g(\theta^t) \end{split}$$

将u替换成v,整理得如下的等价迭代格式:

$$v^{t} = \alpha v^{t-1} - g(\theta^{t})$$
$$\theta^{t+1} = \theta^{t} + \varepsilon v^{t}$$

这就是 CS231n 给出的格式。

1.1.2 Nesterov Momentum

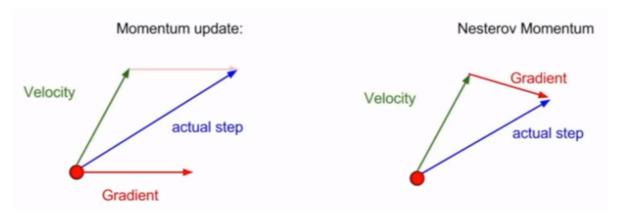
Require: Learning rate ε , momentum parameter α

Require: Initial parameter θ^1 , initial velocity v^0

$$v^t = \alpha v^{t-1} - \varepsilon g(\theta^t + \alpha v^{t-1})$$

 $\theta^{t+1} = \theta^t + v^t$

Nesterov Momentum 和普通的Momentum的区别:



1.2 权重自适应 (Adaptive Learning Rates)

1.2.1 AdaGrad

Require: Global learning rate arepsilon, small contant δ (perhaps 10^{-7}) used for numerical stabilization

Require: Initial parameter $heta^1$, initial gradient accumulation $r^0=\mathbf{0}$

$$r^t = r^{t-1} + g(\theta^t) \odot g(\theta^t) \text{ (\odot: Multiplication applied element-wise)}$$

$$\theta^{t+1} = \theta^t - \frac{\varepsilon}{\delta + \sqrt{r^t}} \odot g(\theta^t) \text{ (Division and square root applied element-wise)}$$

适用于凸优化问题。

1.2.2 RMSProp

The RMSProp algorithm modifies AdaGrad to perfrom better in the non-convex setting by changing the gradient accumulation r into an exponentially weighted moving average.

AdaGrad shrinks the learning rate according to the entire history of the squared gradient and may have made the learning rate too small before arriving at such a convex structure. RMSProp uses an exponentially decaying average to discard history from the extreme past so that it can converge rapidly after finding a convex bowl, as if it were an instance of the AdaGrad algorithm initialized within that bowl.

Require: Global learning rate ε , decay rate ρ controls the length scale of the moving average, small contant δ (usually 10^{-6}) used for numerical stabilization

Require: Initial parameter $heta^1$, initial gradient accumulation $r^0=\mathbf{0}$

$$r^t =
ho r^{t-1} + (1-
ho)g(heta^t) \odot g(heta^t)$$
 $heta^{t+1} = heta^t - rac{arepsilon}{\sqrt{\delta + r^t}} \odot g(heta^t) \left(rac{1}{\sqrt{\delta + r^t}} ext{applied element-wise}
ight)$

1.3 动量法 + 权重自适应

1.3.1 RMSProp with Nesterov momentum

Require: Global learning rate ε , decay rate ρ , momentum cofficient α

Require: Initial parameter $heta^1$, initial velocity v^0 , initial gradient accumulation $r^0=\mathbf{0}$

$$egin{aligned} r^t &=
ho r^{t-1} + (1-
ho) g(heta^t + lpha v^{t-1}) \odot g(heta^t + lpha v^{t-1}) \ v^t &= lpha v^{t-1} - rac{arepsilon}{\sqrt{r}} \odot g(heta^t + lpha v^{t-1}) \ (rac{1}{\sqrt{r^t}} ext{applied element-wise}) \ heta^{t+1} &= heta^t + v^t \end{aligned}$$

1.3.2 Adam

Adam 算法十分健壮,是首选。

Require: Global learning rate ε (Suggested default: 0.001), decay rates ρ_1 and ρ_2 for moment estimates (Suggested default: 0.9 and 0.999 respectively), small contant δ used for numerical stabilization (Suggested default: 10^{-8})

Require: Initial parameter $heta^1$, initial 1st and 2nd moment variables $s^0=\mathbf{0}, r^0=\mathbf{0}$

Biased 1st moment estimate: $s^t = \rho_1 s^{t-1} + (1 - \rho_1) g(\theta^t)$

Biased 2nd moment estimate: $r^t = \rho_2 r^{t-1} + (1 - \rho_2) g(\theta^t) \odot g(\theta^t)$

Corrected 1st moment estimate: $\hat{s}^t = \frac{s^t}{1 - (
ho_1)^t}$

Corrected 2nd moment estimate: $\hat{r}^t = \frac{r^t}{1 - (\rho_2)^t}$

$$heta^{t+1} = heta^t - rac{arepsilon}{\delta + \sqrt{\hat{r}^t}} \hat{s}^t$$

其中, $(\rho_1)^t, (\rho_2)^t$ 表示 t 次方而不是第 t 次迭代。

Adam 名称来源于 Adaptive moments (自适应矩),可视为 RMSProp 和 momentum 相结合的变体:

- 1. 上式中动量 (momentum) v 对应一阶矩 (moment) s , 只不过前面我们没有明确地初始化动量 v^0 , 而这里 $s^0=0$;
- 2. 考虑到初始化时令 $s^0=\mathbf{0}, r^0=\mathbf{0}$,Adam 还包含对一阶矩和二阶矩的偏差修正,比如 t=1 时:

$$s^1 = (1 -
ho_1)g(heta^1)$$
 $r^1 = (1 -
ho_2)g(heta^1)\odot g(heta^1)$

修正过后:

$$\hat{s}^1=rac{s^1}{1-
ho_1}=g(heta^1) \ \hat{r}^1=rac{r^1}{1-
ho_2}=g(heta^1)\odot g(heta^1)$$

而 RMSProp 虽然包含形式相同的二阶矩,且 $r^0=\mathbf{0}$,但 RMSProp 未对其进行修正,因此,相比 Adam,训练初始阶段 RMSProp 对二阶矩的估计可能存在较大的偏差。

2. 二阶优化算法

2.1 牛顿法