
Algorithm 1: The Minibatch Stochastic Gradient Descent Algorithm.

input : Initial value $\theta^{(0)}$, Batch size m , Number of epochs T ,
Learning rate schedule $\eta_1, \eta_2, \dots, \eta_T$
output: minimizer $\theta^{(T)}$

```

1 for  $t = 1$  to  $T$  do
2   Split the data into  $n/m$  minibatches  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{n/m}\}$ 
3   for  $b = 0$  to  $n/m - 1$  do
4      $\theta^{(t,b+1)} := \theta^{(t,b)} - \eta_t \cdot \frac{1}{m} \sum_{i \in \mathcal{B}_b} \nabla \ell(y_i, f(x_i; \theta^{(t,b)}))$ 
5   end
6    $\theta^{(t+1,0)} := \theta^{(t,n/m-1)}$ 
7 end

```

Algorithm 2: The AdaGrad Algorithm.

input : Initial value $\theta^{(0)}$, Batch size m , Number of epochs T ,
Global learning rate η
output: minimizer $\theta^{(T)}$

```

1 for  $t = 1$  to  $T$  do
2   Split the data into  $n/m$  minibatches  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{n/m}\}$ 
3   for  $b = 0$  to  $n/m - 1$  do
4     Let  $\mathbf{g} = \frac{1}{m} \sum_{i \in \mathcal{B}_b} \nabla \ell(y_i, f(x_i; \theta^{(t,b)}))$ .
5     Accumulate squared gradient:  $\mathbf{r} \leftarrow \mathbf{r} + \mathbf{g} \odot \mathbf{g}$ 
6     Compute update:  $\theta^{(t,b+1)} := \theta^{(t,b)} - \frac{\eta}{\sqrt{\mathbf{r} + \epsilon}} \odot \mathbf{g}$  (Division and
      square root applied element-wise;  $\epsilon$  is a small number for
      numerical stability)
7   end
8    $\theta^{(t+1,0)} := \theta^{(t,n/m-1)}$ 
9 end

```

Algorithm 3: The RMSProp Algorithm.

input : Initial value $\theta^{(0)}$, Batch size m , Number of epochs T ,
Global learning rate η , **Decay rate** ρ
output: minimizer $\theta^{(T)}$

```
1 for  $t = 1$  to  $T$  do
2   Split the data into  $n/m$  minibatches  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{n/m}\}$ 
3   for  $b = 0$  to  $n/m - 1$  do
4     Let  $\mathbf{g} = \frac{1}{m} \sum_{i \in \mathcal{B}_b} \nabla \ell(y_i, f(x_i; \theta^{(t,b)}))$ .
5     Accumulate squared gradient:  $\mathbf{r} \leftarrow \rho \mathbf{r} + (1 - \rho) \mathbf{g} \odot \mathbf{g}$ 
6     Compute update:  $\theta^{(t,b+1)} := \theta^{(t,b)} - \frac{\eta}{\sqrt{\mathbf{r} + \epsilon}} \odot \mathbf{g}$  (Division and
       square root applied element-wise;  $\epsilon$  is a small number for
       numerical stability)
7   end
8    $\theta^{(t+1,0)} := \theta^{(t,n/m-1)}$ 
9 end
```

Algorithm 4: The Adam Algorithm.

input : Initial value $\theta^{(0)}$, Batch size m , Number of epochs T ,
Global learning rate η , **Decay rate** ρ_1, ρ_2
output: minimizer $\theta^{(T)}$

```
1 for  $t = 1$  to  $T$  do
2   Split the data into  $n/m$  minibatches  $\{\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_{n/m}\}$ 
3   for  $b = 0$  to  $n/m - 1$  do
4     Let  $\mathbf{g} = \frac{1}{m} \sum_{i \in \mathcal{B}_b} \nabla \ell(y_i, f(x_i; \theta^{(t,b)}))$ .
5     Update biased 1st moment estimate:  $\mathbf{v} \leftarrow \rho_1 \mathbf{v} + (1 - \rho_1) \mathbf{g}$ 
6     Update biased 2nd moment estimate:  $\mathbf{r} \leftarrow \rho_2 \mathbf{r} + (1 - \rho_2) \mathbf{g} \odot \mathbf{g}$ 
7     Correct bias in 1st moment:  $\hat{\mathbf{v}} \leftarrow \frac{\mathbf{v}}{1 - \rho_1^t}$ 
8     Correct bias in 2nd moment:  $\hat{\mathbf{r}} \leftarrow \frac{\mathbf{r}}{1 - \rho_2^t}$ 
9     Compute update:  $\theta^{(t,b+1)} := \theta^{(t,b)} - \frac{\eta}{\sqrt{\hat{\mathbf{r}} + \epsilon}} \odot \hat{\mathbf{v}}$ 
10  end
11   $\theta^{(t+1,0)} := \theta^{(t,n/m-1)}$ 
12 end
```
