
```

1 Input:  $\epsilon_0$ , initial learning rate
2 Input:  $\alpha$ , decay rate of learning rate
3 Input:  $\beta$ , momentum rate
4 Input:  $\rho$ , discount factor for historical gradient
5 Input:  $\zeta$ , small constant to avoid zero division
6 Input:  $m$ , minibatch size
7 Input:  $k$ , epoch size
8 Input:  $\theta$ , initial weights
9 Input:  $v$ , initial velocity
10 Input:  $\mathbf{X}$ , training dataset inputs
11 Input:  $\mathbf{y}$ , training dataset targets
12 Initialize:  $r \leftarrow 0$ , accumulation of historical gradient
13 Initialize:  $j \leftarrow 1$ , current epoch
14 while  $j \leq k$  do
15     update learning rate  $\epsilon_j \leftarrow \epsilon_0 + \alpha(\epsilon_{j-1} - \epsilon_0)$ 
16     while stopping criteria is not satisfied do
17          $\{\mathbf{x}^1 \dots \mathbf{x}^m\}, \{\mathbf{y}^1 \dots \mathbf{y}^m\} \leftarrow$  get a sample from  $\mathbf{X}$  and  $\mathbf{y}$  randomly
18         calculate estimation of gradient  $\hat{g} \leftarrow \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^i; \theta), \mathbf{y}^i)$ 
19         accumulate historical gradients  $r \leftarrow \rho r + (1 - \rho)\hat{g} \odot \hat{g}$ 
20         calculate step size  $v \leftarrow \beta v - \frac{\epsilon_j}{\sqrt{\zeta + r}} \odot \hat{g}$ 
21     update weights  $\theta \leftarrow \theta + v$ 
22  $j \leftarrow j + 1$  go to next epoch

```
