1 **Input:** $\epsilon_0$, initial learning rate
2 **Input:** $\alpha$, decay rate of learning rate
3 **Input:** $\beta_1$, 1st order moment for plain gradients
4 **Input:** $\beta_2$, 2nd order moment for squared gradients
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:** $\mathbf{X}$, training dataset inputs
10 **Input:** $\mathbf{y}$, training dataset targets
11 **Initialize:** $s \leftarrow 0$, accumulation variable for historical gradients
12 **Initialize:** $r \leftarrow 0$, accumulation variable for historical squared gradients
13 **Initialize:** $t \leftarrow 0$, counter of each gradient update
14 **Initialize:** $j \leftarrow 1$, current epoch
15 **while** $j \leq k$ **do**
16     update learning rate $\epsilon_j \leftarrow \epsilon_0 + \alpha(\epsilon_{j-1} - \epsilon_0)$
17     **while** *stopping criteria is not satisfied* **do**
18         $\{\mathbf{x}^1...\mathbf{x}^m\}, \{\mathbf{y}^1...\mathbf{y}^m\} \leftarrow$ get a sample from $\mathbf{X}$ and $\mathbf{y}$ randomly
19         calculate estimation of gradient $\hat{g} \leftarrow \frac{1}{m} \sum_{i=1}^{m} L(f(\mathbf{x}^i; \theta), \mathbf{y}^i)$
20         accumulate historical graidents $s \leftarrow \beta_1 s + (1 - \beta_1)\hat{g}$
21         accumulate historical squared graidents $r \leftarrow \beta_2 r + (1 - \beta_2)\hat{g} \odot \hat{g}$
22         $t \leftarrow t + 1$
23         apply bias correction to 1st order momentum $\hat{s} = \frac{s}{1-\beta_1^t}$
24         apply bias correction to 2nd order momentum $\hat{r} = \frac{r}{1-\beta_2^t}$
25         calculate step size $\Delta\theta \leftarrow -\epsilon_j \frac{\hat{s}}{\sqrt{\zeta+\hat{r}}}$
26         update weights $\theta \leftarrow \theta + \Delta\theta$
27     $j \leftarrow j + 1$ go to next epoch