## 1. **Optimizer**

1.1. **Momentum.** Momentum is a method that helps accelerate SGD in the relevant direction and dampens oscillations. It does this by adding a fraction $\gamma$ of the update vector of the past time step to the current update vector:

$$\text{(1)} \qquad v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta)$$

$$\text{(2)} \qquad \theta = \theta - v_t,$$

where the momentum term is usually set to 0.9 or a similar value.
The momentum term <mark>increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions</mark>. As a result, we gain faster convergence and reduced oscillation.

1.2. **Nesterov.** We know that we will use our momentum term $\gamma v_{t-1}$ to move the parameters $\theta$. Computing $\theta - \gamma v_{t-1}$ thus gives us an *approximation* of the next position of the parameters. We can now effectively look ahead by calculating the gradient not w.r.t. to our current parameters $\theta$ but w.r.t. the approximate future position of our parameters:

$$\text{(3)} \qquad v_t = \gamma v_{t-1} + \eta \nabla_\theta J(\theta - \gamma v_{t-1})$$

$$\text{(4)} \qquad \theta = \theta - v_t,$$

Momentum first computes the current gradient ($J(\theta)$), and then takes a big jump in the direction of the updated accumulated gradient ($\theta = \theta - v_t$). <mark>Nesterov accelerated gradient (NAG) first makes a big jump in the direction of the previous accumulated gradient</mark> ($\theta - \gamma v_{t-1}$), measures the gradient ($J(\theta - \gamma v_{t-1})$) and then makes a correction ($\theta = \theta - v_t$) which results in the complete NAG update.

1.3. **Adagrad.** Previously, we performed an update for all parameters $\theta$ at once as every parameter $\theta_i$ used the same learning rate $\eta$. Adagrad uses <mark>a different learning rate for every parameter $\theta_i$ at every time step $t$</mark>. For brevity, we use $g_t$ to denote the gradient at time step $t$. $g_{t,i}$ is then the partial derivative of the objective function w.r.t. to the parameter $\theta_i$ at time step $t$:

$$\text{(5)} \qquad g_{t,i} = \nabla J(\theta_{t,i}).$$

Adagrad sets learning rate $\eta$ at each time step $t$ for every parameter $\theta_i$ based on the past gradients that have been computed for $\theta_i$

$$\text{(6)} \qquad \theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} g_{t,i},$$

where $G_t \in \mathbb{R}^{d \times d}$ is a diagonal matrix. Each diagonal element in $G_t$ is the sum of the squares of the gradients w.r.t. $\theta_i$ up to time step $t$. $\epsilon$ is a smoothing term that avoids division by zero (e.g., $10^{-8}$).
As $G_t$ contains the sum of the squares of the past gradients w.r.t. to all parameters $\theta$ along its diagonal, we can now vectorize our implementation by performing a matrix-vector product $\odot$ between $G_t$ and $g_t$:

$$\text{(7)} \qquad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t,$$

One of Adagrad's main benefits is that it eliminates the need to manually tune the learning rate. Most implementations use a default value of 0.01 and leave it at that. Adagrad's main weakness is its accumulation of the squared gradients in the denominator: Since every added term is positive, the accumulated sum keeps growing during training. This in turn causes the

learning rate to shrink and eventually become infinitesimally small, at which point the algorithm is no longer able to acquire additional knowledge.

**1.4. Adadelta.** Adadelta is an extension of Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size $\omega$. Instead of inefficiently storing $\omega$ previous squared gradients, the sum of gradients is recursively defined as a decaying average of all past squared gradients. The running average $E[g^2]_t$ at time step $t$ then depends (as a fraction $\gamma$ similarly to the Momentum term) only on the previous average and the current gradient:

$$(8) \qquad E[g^2]_t = \gamma E[g^2]_{t-1} + (1-\gamma)g_t^2,$$

where $\gamma$ can be set as a similar value as the momentum term, around 0.9.

Based on Eq. (7), we simply replace the diagonal matrix $G_t$ with the decaying average over past squared gradients $E[g^2]_t$, and get

$$(9) \qquad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}g_t.$$

Note that the denominator is just the root mean squared (RMS) of the gradient, we can rewrite it as

$$(10) \qquad \theta_{t+1} = \theta_t - \frac{\eta}{RMS[g]_t}g_t.$$

Note that the units in this update (as well as in SGD, Momentum, or Adagrad) do not match, i.e. the update should have the same hypothetical units as the parameter. To realize this, we define another exponentially decaying average, this time not of squared gradients but of squared parameter updates:

$$(11) \qquad E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1-\gamma)\theta_t^2,$$

where

$$(12) \qquad \Delta\theta_t = -\frac{\eta}{RMS[g]_t}g_t.$$

The root mean squared error of parameter updates is thus:

$$(13) \qquad RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}.$$

Since $RMS[\Delta\theta]_t$ is unknown, we approximate it with the RMS of parameter updates until the previous time step. Replacing the learning rate $\eta$ in the previous update rule with $RMS[\Delta\theta]_{t-1}$ finally yields the Adadelta update rule:

$$(14) \qquad \Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t}g_t.$$

$$(15) \qquad \theta_{t+1} = \theta_t + \Delta\theta_t.$$

With Adadelta, we do not even need to set a default learning rate, as it has been eliminated from the update rule.

**1.5. RMSprop.** RMSprop and Adadelta have both been developed independently around the same time stemming from the need to resolve Adagrad's radically diminishing learning rates. RMSprop in fact is identical to the first update vector of Adadelta that we derived above:

$$(16) \qquad E[g^2]_t = \gamma E[g^2]_{t-1} + (1-\gamma)g_t^2,$$

$$(17) \qquad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}}\ g_t.$$

1.6. **Adam.** Adaptive Moment Estimation (Adam) is another method that computes adaptive learning rates for each parameter. In addition to storing an exponentially decaying average of past squared gradients $v_t$ ike Adadelta and RMSprop, Adam also keeps an exponentially decaying average of past gradients $m_t$, similar to momentum. We compute the decaying averages of past and past squared gradients $m_t$ and $v_t$ respectively as follows:

$$(18) \qquad\qquad m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$(19) \qquad\qquad v_t = \beta_1 v_{t-1} + (1 - \beta_1) g_t^2$$

$m_t$ and $v_t$ are are estimates of the first moment (the mean) and the second moment (the uncentered variance) of the gradients respectively, hence the name of the method. As $m_t$ and $v_t$ are initilizaed as vectors of 0's, the authors of Adam observe that they are biased towards zero, especially during the initial time steps, and especially when the decay rates are small. They counteract these biases by computing bias-corrected first and second moment estimates:

$$(20) \qquad\qquad \hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$(21) \qquad\qquad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

They then use these to update the parameters just as we have seen in Adadelta and RMSprop, which yields the Adam update rule:

$$(22) \qquad\qquad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t.$$