**Gradient Descent**

For each epoch

      For each batch:

        Calculate dw, db

        w = w - alpha * dw

        b = b - alpha * db

**Exponential Moving Averages**

$V_0 = 0$

$V1 = \beta*V0 + (1-\beta)*\Theta1$ [Current observation at time 1] [$\beta$=0.9]

$V2 = \beta*V1 + (1-\beta)*\Theta2$ [Current observation at time 2]

$Vt = \beta*Vt\text{-}1 + (1-\beta)*\Theta t$[Current observation at time t]

Generalised form

$V = \beta*V + (1-\beta)*\Theta$

where  $\Theta1, \Theta2, \Theta3$ … are observations at time t=1,2,3

**Gradient Descent (mini-batch) with Momentum (SGD - Minibatch)**

For each batch:

Compute dw, db

$Vdw = \beta*Vdw + (1-\beta)*dw$

$Vdb = \beta*Vdb + (1-\beta)*db$

w = w - alpha * Vdw [ alpha = Learning rate]

b = b - alpha * Vdb

**Root Mean Square Propagation (RMSProp)**

$Sdw = \beta2*Sdw + (1-\beta2)*dw\text{^}2$

$Sdb = \beta2*Sdb + (1-\beta2)*db\text{^}2$

w = w - alpha * dw/**sqrt**(Sdw + epsilon)

b = b - alpha * db/**sqrt**(Sdb + epsilon)

Epsilon = small value to prevent division by 0. Normally 10^(-8)

**Adaptive Moment with Estimation (Adam)**

For each batch:

$Vdw = \beta1*Vdw + (1-\beta1)*dw$

$Vdw = Vdw/(1-\beta1^t)$

$Vdb = \beta1*Vdb + (1-\beta1)*db$

$Vdb = Vdb/(1-\beta1^t)$

$Sdw = \beta2*Sdw + (1-\beta2)*dw\text{^}2$

$Sdw = Sdw/(1-\beta2^t)$

$Sdb = \beta2*Sdb + (1-\beta2)*db\text{^}2$

$Sdb = Sdb/(1-\beta2^t)$

$$w = w - alpha * Vdw/[\textbf{sqrt}(Sdw)+Epsilon]$$
$$b = b - alpha * Vdb/[\textbf{sqrt}(Sdb)+Epsilon]$$

**β1 = 0.9**
**β2 = 0.999**
**Epsilon = $10^{-8}$**