

FEW VARIANTS OF SGD

Method	Formula
Learning Rate	$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \cdot \nabla \ell(\mathbf{w}^{(t)}, \mathbf{z}) = \mathbf{w}^{(t)} - \eta \cdot \nabla \mathbf{w}^{(t)}$
Adaptive Learning Rate	$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \cdot \nabla \mathbf{w}^{(t)}$
Momentum [Qian 1999]	$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \mu \cdot (\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}) - \eta \cdot \nabla \mathbf{w}^{(t)}$
Nesterov Momentum [Nesterov 1983]	$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \mathbf{v}_t; \quad \mathbf{v}_{t+1} = \mu \cdot \mathbf{v}_t - \eta \cdot \nabla \ell(\mathbf{w}^{(t)} - \mu \cdot \mathbf{v}_t, \mathbf{z})$
AdaGrad [Duchi et al. 2011]	$\mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} - \frac{\eta \cdot \nabla \mathbf{w}_i^{(t)}}{\sqrt{A_{i,t} + \epsilon}}; \quad A_{i,t} = \sum_{\tau=0}^t \left(\nabla \mathbf{w}_i^{(\tau)} \right)^2$
RMSProp [Hinton 2012]	$\mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} - \frac{\eta \cdot \nabla \mathbf{w}_i^{(t)}}{\sqrt{A'_{i,t} + \epsilon}}; \quad A'_{i,t} = \beta \cdot A'_{i,t-1} + (1 - \beta) \left(\nabla \mathbf{w}_i^{(t)} \right)^2$
Adam [Kingma and Ba 2015]	$\mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} - \frac{\eta \cdot M_{i,t}^{(1)}}{\sqrt{M_{i,t}^{(2)} + \epsilon}}; \quad M_{i,t}^{(m)} = \frac{\beta_m \cdot M_{i,t-1}^{(m)} + (1 - \beta_m) \left(\nabla \mathbf{w}_i^{(t)} \right)^m}{1 - \beta_m^t}$