SGD:

$$g_t = \hat{g}_t^{\text{SGD}} = \left. \frac{\partial}{\partial \theta} \mathbb{E}[L(\theta)] \right|_{\theta=\theta_t} \tag{1}$$

SGD with momentum:

$$\hat{g}_t^{\text{Mom}} = \beta \hat{g}_{t-1}^{\text{Mom}} + (1-\beta)g_t \tag{2}$$

$$= \sum_{t=1}^{T} (1-\beta)\beta^{T-t} g_t \tag{3}$$

$$\left. \frac{\partial \hat{L}^{\text{Mom}}}{\partial \theta} \right|_{\theta=\theta_t} = \hat{g}_t^{\text{Mom}} \tag{4}$$

$$\hat{L}^{\text{Mom}} = \int \sum_{t=1}^{T} (1-\beta)\beta^{T-t} g_t d\theta \tag{5}$$

$$= (1-\beta) \sum_{t=1}^{T} \beta^{T-t} \int g_t d\theta \tag{6}$$

$$= (1-\beta) \sum_{t=1}^{T} \beta^{T-t} \mathbb{E}[L(\theta_t)] \tag{7}$$

SGD with momentum is minimizing an exponential average of the loss of the model. So it is stabilizing learning.