**Q(5)** Adam is suitable for large models with complex datasets, where there are noisy gradients or vanishing gradients /exploding gradients.

Mathematical Formulation

$f_t \rightarrow$ ~~hypoth~~ at $t^{th}$ step (hypothesis)

Repeat till convergence :

$\theta \rightarrow$ parameters

$$g_t = \nabla_\theta f_t(\theta_{t-1})$$

$$m_t \leftarrow \beta_1 m_{t-1} + (1-\beta_1) \cdot g_t$$

[update momentum (weighted avg of gradients)]

$$v_t \leftarrow \beta_2 v_{t-1} + (1-\beta_2) \cdot g_t^2$$

[update second moment (variance of gradients)]

$$\hat{m}^t \leftarrow \frac{m_t}{1-\beta_1^t}$$

$$\hat{v}^t \leftarrow \frac{v_t}{1-\beta_2^t}$$

(Bias-correction)

$$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \varepsilon}$$