**2 a)** Adam

$$V_t \leftarrow B_2 V_{t-1} + (1-B_2) g_t^2$$
$$\theta_t \leftarrow \theta_{t-1} - \alpha_A m_t / (\sqrt{V_t} + \epsilon_A)$$
$$\hookrightarrow m_t \leftarrow B_1 m_{t-1} + (1-B_1) g_t$$

RMSProp

$$V_t \leftarrow \gamma V_{t-1} + (1-\gamma) g_t^2$$
$$\theta_t \leftarrow \theta_{t-1} - \alpha_R g_t / (\sqrt{V_t} + \epsilon_R)$$

$\therefore$ Set $B_2 = \gamma$.

$$B_1 = 0 \Rightarrow m_t = g_t$$
$$\alpha_A = \alpha_R$$
$$\epsilon_A = \epsilon_R$$

**b)** SGD w/ momentum

$$P_t \leftarrow \mu P_{t-1} - (1-\mu) \nabla F(\theta_{t-1})$$
$$\theta_t \leftarrow \theta_{t-1} + \alpha_S P_t$$

Want: $\alpha_S P_t \approx \alpha_A m_t / (\sqrt{V_t} + \epsilon_A)$

Setting $B_2 = 1 \Rightarrow V_t \leftarrow V_{t-1} = 0 \quad \forall t$

Setting $\epsilon_A = 1, \alpha_A = -\alpha_S$

$$\alpha_S P_t = -\alpha_S m_t$$

Want: $P_t = -m_t$.

Setting $B_1 = \mu \Rightarrow m_t = \mu m_{t-1} + (1-\mu) \nabla F(\theta_{t-1})$

$\rightarrow$ update eq's almost the same. $m_t$ sums the gradient, $P_t$ sums the neg.

Noting that $P_0 = m_0 = 0$, then:

$$P_1 = \mu P_0^{\nearrow 0} - (1-\mu) \nabla F(\theta_0)$$
$$P_2 = -\mu(1-\mu) \nabla F(\theta_0) - (1-\mu) \nabla F(\theta_1)$$
$$P_3 = -\mu^2(1-\mu) \nabla F(\theta_0) - \mu(1-\mu) \nabla F(\theta_1) - (1-\mu) \nabla F(\theta_2)$$
$$P_T = \sum_{t=1}^{T} - \mu^{T-t}(1-\mu) \nabla F(\theta_{t-1})$$

Similarly: $m_T = \sum + \mu^{T-t}(1-\mu) \nabla F(\theta_{t-1})$

as required. $P_t = -m_t$

c) Suppose $\tilde{F} = C \cdot F$

$\tilde{g}_t \leftarrow \nabla C \cdot F = C \cdot \nabla F = C \cdot g_t$

$\tilde{m}_t \leftarrow B_1 \tilde{m}_{t-1} + (1 - B_1) C g_t$

$\tilde{v}_t \leftarrow B_2 \tilde{v}_{t-1} + (1 + B_2) C^2 g_t^2$

$\tilde{\theta}_t \leftarrow \tilde{\theta}_{t-1} - \alpha_A \tilde{m}_t / \sqrt{\tilde{v}_t}$

Recall: $m_T = \sum_{t=1}^{T} B_1^{T-t} (1 - B_1) \nabla F(\theta_{t-1})$

Clearly $\tilde{m}_T = C m_T$ since $\tilde{\nabla F} = C \nabla F$.

By a very similar argument (e.g. let $\hat{g}_t = g_t^2$)

$$v_T = \sum B_2^{T-t} (1 - B_2)(\nabla F(\theta_{t-1}))^2$$

Again, noting $\nabla \tilde{F} = C \nabla F$, $\tilde{v}_T = C^2 v_T$

Then $\tilde{\theta}_t \leftarrow \tilde{\theta}_{t-1} - \alpha_A C m_t / \sqrt{C^2 v_T}$

$\tilde{\theta}_t \leftarrow \tilde{\theta}_{t-1} - \alpha_A m_t / \sqrt{v_T}$

Assuming $\tilde{\theta}_{t-1} = \theta_{t-1}$, this is the same recurrence.

Base case: $\theta_0 = 0 = \tilde{\theta}_0$.

$\therefore$ If $\epsilon_A = 0$, then the trajectory is invariant to the scale of the loss fct.