Part2

a) let $V_t$ of Adam equal to $V_t$ of RMSprop
   that is
   $$B_2 V_{t-1} + (1-B_2) g_t^2 = \gamma V_{t-1} + (1-\gamma) g_t^2$$
   So $B_2 = \gamma$
   let $\theta_t$ of Adam equal to $\theta_t$ of RMSprop
   then that is
   $$\theta_{t-1} - d_A m_t / (\sqrt{V_t} + \epsilon_A) = \theta_{t-1} - d_R g_t / (\sqrt{V_t} + \epsilon_R)$$
   So $m_t = g_t$, $\epsilon_A = \epsilon_R$, $d_A = d_R$,
   By Adam, we know
   $$m_t \leftarrow B_1 m_{t-1} + (1-B_1) g_t$$
   to make $m_t = g_t$, we need
   $$B_1 = 0,$$
   the we get hyperparameters $(d_R, 0, \gamma, \epsilon_R)$
   that matches $(d_A, B_1, B_2, \epsilon_A)$

B) We want to make Adam approximately equivalent to momentum SGD, we need to find a set of $(\alpha_A, \beta_1, \beta_2, \epsilon_A)$ such that

$$\theta_{t+1} = \alpha_A m_t / (\sqrt{v_t} + \epsilon_A) \approx \theta_{t+1} + \alpha_s P_t$$

then let $\alpha_A = \alpha_s$, we get

$$m_t / (\sqrt{v_t} + \epsilon_A) \approx -P_t$$

According to the algorithms; we can get

$$(\beta_1 m_{t-1} + (1-\beta_1) g_t / \sqrt{\beta_2 v_{t-1} + (1-\beta_2) g_t^2} + \epsilon_A \approx -(\mu P_{t-1} - (1-\mu) \nabla (\theta_{t-1}))$$

let $\beta_2 = 1$, $\epsilon_A = 1$ we can get

$$\frac{\beta_1 m_{t-1} + (1-\beta_1) g_t}{\sqrt{v_{t-1}} + 1} \qquad \text{where } v_{t-1} = v_t$$

By Adam, $v_0 = 0$, $v_t$ will be $0$ on all iterations. then $\sqrt{v_{t-1}} + 1 = 1$

then we got

$$\underset{m_t}{\underbrace{\beta_1 m_{t-1} + (1-\beta_1) g_t}} \approx \underset{-P_t}{\underbrace{-\mu P_{t-1} + (1-\mu) g_t}}$$

let $\beta_1 = \mu$, if $\mu$ is small enough, then they will be closer

Then we can find hyperparameters $(\alpha_s, \mu, 1, 1)$ where $\mu \to 0$.

C) According to the question, we can denote the quantities as $\tilde{g}_t$, $\tilde{m}_t$, $\tilde{v}_t$, $\tilde{\theta}_t$

WTS for $\epsilon_A = 0$, Adam is invariant to rescaling. that is $\tilde{\theta}_t = \theta_t$ for $\forall t \in N$, by hint we can use induction.

Base case:
  let $m_0 = \tilde{m}_0 = v_0 = \tilde{v}_0 = 0$, $\tilde{\theta}_0 = \theta_0$.
  WTS $\tilde{\theta}_1 = \theta_1$.
  We know $\tilde{g}_1 = (c \nabla J)(\tilde{\theta}_0)$
  then $= (c \nabla J)(\partial_0)$ since $\tilde{\theta}_0 = \theta_0$.
  $= cg_1$

  $\tilde{m}_1 = \beta_1 \tilde{m}_0 + (1-\beta_1)\tilde{g}_1$  $(\tilde{m}_0 = 0)$
  $= (1-\beta_1)(cg_1$

  $m_1 = \beta_1 m_0 + (1-\beta_1)g_1$
  $= (1-\beta_1)g_1$
  then $\tilde{m}_1 = cm_1$

  $\tilde{v}_1 = \beta_2 \tilde{v}_0 + (1-\beta_2)\tilde{g}_1^2$
  $= (1-\beta_2)c^2 g_1^2$   $(\tilde{v}_0 = 0)$

  $v_1 = \beta_2 v_0 + (1-\beta_2)g_1^2$
  $= (1-\beta_2)g_1^2$

  then $\tilde{v}_1 = c^2 v_1$

  $\tilde{\theta}_1 = \tilde{\theta}_0 - \alpha_A \tilde{m}_1 / (\sqrt{\tilde{v}_1} + \epsilon_A)$  $(\epsilon_A = 0)$
  $= \theta_0 - \alpha_A cm_1 / \sqrt{c^2 v_1}$   (By above)
  $= \theta_0 - \alpha_A m_1 / \sqrt{v_1}$
  $= \theta_0 - \alpha_A m_1 / (\sqrt{v_1} + \epsilon_A) = \theta_1$

I.S.

I.H $\tilde{\theta}_{t-1} = \theta_{t-1}$, $\tilde{m}_{t-1} = c \, m_{t-1}$, $\tilde{v}_{t-1} = c^2 v_{t-1}$

WTS. $\tilde{\theta}_t = \theta_t$

we get $\tilde{g}_t = c \, \nabla j(\tilde{\theta}_{t-1})$ $\quad (\tilde{\theta}_{t-1} = \theta_{t-1})$

$\qquad = c \cdot \nabla j(\theta_{t-1})$

$\qquad = c \, g_t$

$\tilde{m}_t = \beta_1 \tilde{m}_{t-1} + (1-\beta_1) \tilde{g}_t$

$\qquad = \beta_1 c \cdot m_{t-1} + (1-\beta_1) \cdot c \cdot g_t \quad$ (By I.H)

$\qquad = c \, (\beta_1 m_{t-1} + (1-\beta_1) g_t)$

$\qquad = c \, m_t$

$\tilde{v}_t = \beta_2 \tilde{v}_{t-1} + (1-\beta_2) \tilde{g}_t^2$

$\qquad = \beta_2 \, c^2 v_{t-1} + (1-\beta_2) c^2 \cdot g_t^2 \quad$ (By I.H)

$\qquad = c^2 (\beta_2 v_{t-1} + (1-\beta_2) g_t^2)$

$\qquad = c^2 \, v_t$

then $\tilde{\theta}_t = \tilde{\theta}_{t-1} - \alpha_A \tilde{m}_t / (c\sqrt{\tilde{v}_t} + \epsilon_A) \quad (\epsilon_A = 0)$

$\qquad = \theta_{t-1} - \alpha_A \, c \, m_t / \sqrt{c^2 v_t}$

$\qquad = \theta_{t-1} - \alpha_A \, m_t / (\sqrt{v_t} + \epsilon_A)$

$\qquad = \theta_t$

then we get Adam is invariant
to this rescaling