ⓐ Show that $\sum_{t}^{T} \mathbb{E}_{\tau \sim P_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) \right] = 0$

$$P_\theta(\tau) = P_\theta(s_t, a_t) P_\theta \left( \frac{\tau}{s_t, a_t} \Big| s_t, a_t \right)$$

$\mathbb{E}_{\tau \sim P_\theta} \left[ \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) \right]$

$= \int_\tau \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) p(\tau) \, \partial\tau$

$= \iint \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) P_\theta(s_t, a_t) P_\theta \left( \frac{\tau}{s_t, a_t} \Big| s_t, a_t \right) \partial(s_t, a_t) \partial\left( \frac{\tau}{s_t, a_t} \right)$

$= \int \nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t) P_\theta(s_t, a_t) \partial(s_t, a_t) \underbrace{\int P_\theta \left( \frac{\tau}{s_t, a_t} \right) \partial\left( \frac{\tau}{s_t, a_t} \right)}_{1}$

$\nabla p(x) = p(x) \nabla \log p(x)$

$p(x, y) = p(x) p(y | x)$
$P_\theta(a_t | s_t) = \pi_\theta(a_t | s_t)$

$= \iint \frac{\nabla_\theta \pi_\theta(a_t | s_t)}{\pi_\theta(a_t | s_t)} b(s_t) \underbrace{P_\theta(s_t, a_t)}_{P_\theta(a_t | s_t) p(s_t) = \pi_\theta(a_t | s_t) p(s_t)} \partial s_t \, \partial a_t$

$= \iint \nabla_\theta \pi_\theta(a_t | s_t) b(s_t) p(s_t) \, \partial s_t \, \partial a_t$

$= \int p(s_t) b(s_t) \partial s_t \, \nabla_\theta \int \pi_\theta(a_t | s_t) \, \partial a_t$

$= \int p(s_t) b(s_t) \partial s_t \, \nabla_\theta 1$

$= 0$

ⓑ   $p_\theta(\tau) = p_\theta(s_{0:t}, a_{0:t-1}) \, p_\theta(s_{t+1:T}, a_{t:T} \mid s_{0:t}, a_{0:t-1})$

show   $\sum\limits_{t=0}^{T} \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[ \nabla_\theta \log \underbrace{p_\theta(a_t \mid s_t)}_{= \pi_\theta(a_t \mid s_t)} \, b(s_t) \right] = 0$

$P(s_x, a_t \mid \pi) = \sum\limits_{s_0} \sum\limits_{a_0} \sum\limits_{s_1} \cdots \sum\limits_{s_{t+1}} \sum\limits_{a_{t+1}} P(s_0, a_0, s_1, a_1 \cdots s_{t+1}, a_{t+1}, \cdots \mid \pi)$

$\sum\limits_{t} \sum\limits_{\tau} p_\theta(\tau) \nabla_\theta \log p_\theta(a_t \mid s_t) \, b(s_t)$

$= \sum\limits_{t} \sum\limits_{s_{0:t}} \sum\limits_{a_{0:t-1}} \sum\limits_{s_{t+1:T}} \sum\limits_{a_{t:T}} P(s_{0:t}, a_{0:t-1}) \, P(s_{t+1:T}, a_{t:T} \mid s_{0:t}, a_{0:t-1}) \cdot$

$\nabla_\theta \log p_\theta(a_t \mid s_t) \, b(s_t)$

$= \sum\limits_{t} \mathbb{E}_{s_{0:t}, a_{0:t-1}} \left[ b(s_t) \, \mathbb{E}_{s_{t+1:T}, a_{t:T}} \left[ \nabla_\theta \log p_\theta(a_t \mid s_t) \mid s_{0:t}, a_{0:t-1} \right] \right]$

<span style="color:blue">Markov property means $s_{t+1:T}, a_{t:T}$ only depends on $s_t$</span>

$\mathbb{E}_{s_{t+1:T}, a_{t:T}} \left[ \nabla_\theta \log p_\theta(a_t \mid s_t) \mid s_t \right]$

$= \sum\limits_{s_{t+1}} \sum\limits_{s_{t+2}} \cdots \sum\limits_{s_T} \cdots \sum\limits_{a_t} \cdots \sum\limits_{a_T} p_\theta(s_{t+1} \mid s_t, a_t) p_\theta(s_{t+2} \mid s_{t+1}, a_{t+1}) \cdots p_\theta(a_t \mid s_t) p_\theta(a_{t+1} \mid s_{t+1})$

$\cdots p_\theta(a_T \mid s_T) \, p_\theta(s_T \mid s_{T-1}, a_{T-1}) \cdot \nabla_\theta \log p_\theta(a_t \mid s_t)$

$= \sum\limits_{a_t} p_\theta(a_t \mid s_t) \nabla_\theta \log p_\theta(a_t \mid s_t) \underbrace{\sum\limits_{s_{t+1}} p(s_{t+1} \mid s_t, a_t)}_{} \sum\limits_{s_{t+2}} \cdots$ <span style="color:blue">↗ 1</span>

<span style="color:blue">all sum to 1 since main term doesn't use any variables besides $s_t, a_t$</span>

$$= \sum_{s_t} \sum_{a_t} P_\theta(a_t | s_t) \nabla_\theta \log P_\theta(a_t | s_t) P_\theta(s_t) b(s_t)$$

$$\nabla P(x) = P(x) \nabla \log P(x)$$

$$\underbrace{\frac{\nabla P_\theta(a_t | s_t)}{P_\theta(a_t | s_t)}}$$

$$= \sum_{s_t} \sum_{a_t} \nabla_\theta P_\theta(a_t | s_t) P_\theta(s_t) b(s_t)$$

$$= \sum_{s_t} b(s_t) P_\theta(s_t) \nabla_\theta \underbrace{\sum_{a_t} P_\theta(a_t | s_t)}_{1} \quad 1$$

$$\to 0$$

$$= \sum_{s_t} b(s_t) P_\theta(s_t) \nabla_\theta 1$$

$$= 0$$

$$\nabla P(x) = P(x) \nabla \log P(x)$$