



(a)

$$\text{MAP}(R_\theta \mid o^1, o^2, \dots, o^N)$$

$$\arg \max_{R_\theta} \Pr(R_\theta \mid o^1, o^2 \dots o^N) \quad (1)$$

$$= \arg \max_{R_\theta} \Pr(o^1, o^2 \dots o^N \mid R_\theta) \cdot \Pr(R_\theta) \quad (2)$$

$$= \arg \max_{R_\theta} (\log \Pr(o^1, o^2 \dots o^N \mid R_\theta) \cdot \Pr(R_\theta)) \quad (3)$$

$$= \arg \max_{R_\theta} (\log \Pr(o^1, o^2 \dots o^N \mid R_\theta) + \log \Pr(R_\theta)) \quad (4)$$

$$= \arg \max_{R_\theta} (\log \sum_{s^1, a^1, \dots, s^N, a^N} \Pr(o^1, o^2 \dots, o^N, s^1, a^1 \dots, s^N, a^N \mid R_\theta) + \log \Pr(R_\theta)) \quad (5)$$

**Denote:**  $\tau, L_\theta^{lh}, L_\theta^{pr}$

$$\tau = \{(s^1, a^1), (s^2, a^2), \dots, (s^N, a^N)\} \quad (6)$$

$$L_\theta^{pr} = \log \Pr(R_\theta) \quad (7)$$

$$L_\theta^{lh} = \log \sum_{\tau} \Pr(o^1, o^2 \dots, \tau \mid R_\theta) \quad (8)$$

**Goal:**

$$L_\theta = L_\theta^{pr} + L_\theta^{lh} \quad (9)$$

$$\frac{\partial L_\theta}{\partial \theta} = \frac{\partial L_\theta^{pr}}{\partial \theta} + \frac{\partial L_\theta^{lh}}{\partial \theta} \quad (10)$$

**(b)**

**Derivative of the log prior:**  $L_\theta^{pr}$

$$L_\theta^{pr} = \log Pr(R_\theta) \quad (11)$$

$$Pr(R_\theta) = \frac{1}{(2\pi)^{d/2} |\Sigma_{R_\theta}|^{1/2}} \exp\left(-\frac{1}{2}(R_\theta - \mu_{R_\theta})^T \Sigma_{R_\theta}^{-1} (R_\theta - \mu_{R_\theta})\right) \quad (12)$$

$$\frac{\partial \log Pr(R_\theta)}{\partial \theta} = \frac{1}{Pr(R_\theta)} \cdot \frac{\partial Pr(R_\theta)}{\partial R_\theta} \cdot \frac{\partial R_\theta}{\partial \theta} \quad (13)$$

$$= \frac{1}{Pr(R_\theta)} \cdot (-Pr(R_\theta) \cdot \Sigma_{R_\theta}^{-1} (R_\theta - \mu_{R_\theta})) \cdot \frac{\partial R_\theta}{\partial \theta} \quad (14)$$

$$= -\Sigma_{R_\theta}^{-1} (R_\theta - \mu_{R_\theta}) \cdot \frac{\partial R_\theta}{\partial \theta} \quad (15)$$

**(c)**

**Derivative of the log likelihood:**  $L_\theta^{lh}$

$$L_\theta^{lh} = \log \sum_{\tau} Pr(o^1, o^2, \dots, \tau \mid R_\theta) \quad (16)$$

$$= \log \sum_{\tau} Pr(s^1) \prod_{i=1}^N Pr(o^i \mid s^i, a^i) \prod_{j=2}^N Pr(s^j \mid s^{j-1}, a^{j-1}) \prod_{z=1}^N Pr(a^z \mid s^z, R_\theta) \quad (17)$$

**Denote:**  $h_\theta(\tau), \pi_\theta(a^z \mid s^z)$

$$\pi_\theta(a^z | s^z) = Pr(a^z | s^z, R_\theta) \quad (18)$$

$$h_\theta(\tau) = Pr(s^1) \prod_{i=1}^N Pr(o^i | s^i, a^i) \prod_{j=2}^N Pr(s^j | s^{j-1}, a^{j-1}) \prod_{z=1}^N \pi_\theta(a^z | s^z) \quad (19)$$

$$L_\theta^{lh} = \log \sum_{\tau} h_\theta(\tau) \quad (20)$$

$$\frac{\partial L_\theta^{lh}}{\partial \theta} = \frac{1}{\sum_{\tau} h_\theta(\tau)} \cdot \frac{\partial \sum_{\tau} h_\theta(\tau)}{\partial \theta} \quad (21)$$

$$= \frac{1}{\sum_{\tau} h_\theta(\tau)} \cdot \sum_{\tau} \frac{\partial h_\theta(\tau)}{\partial \theta} \quad (22)$$

**(c)**

**Derivative of  $h_\theta(\tau)$**

$$h_\theta(\tau) = Pr(s^1) \prod_{i=1}^N Pr(o^i | s^i, a^i) \prod_{j=2}^N Pr(s^j | s^{j-1}, a^{j-1}) \prod_{z=1}^N \pi_\theta(a^z | s^z) \quad (23)$$

$$= c(\tau) \prod_{z=1}^N \pi_\theta(a^z | s^z) \quad (24)$$

$$\frac{\partial h_\theta(\tau)}{\partial \theta} = c(\tau) \cdot \left( \sum_{z=1}^N \left( \frac{\partial \pi_\theta(a^z | s^z)}{\partial \theta} \prod_{j \neq z}^N \pi_\theta(a^j | s^j) \right) \right) \quad (25)$$

$$\pi_\theta(a^z | s^z) = \frac{\exp(Q_\theta^*(s^z, a^z))}{\sum_{a' \in A} \exp(Q_\theta^*(s^z, a'))} \quad (26)$$

$$\frac{\partial \pi_\theta(a^z | s^z)}{\partial \theta} = \pi_\theta(a^z | s^z) \left( \frac{\partial Q_\theta^*(s^z, a^z)}{\partial \theta} - \sum_{a' \in A} \pi_\theta(a' | s^z) \frac{\partial Q_\theta^*(s^z, a')}{\partial \theta} \right) \quad (27)$$

**(d)**

**Derivative of  $Q_\theta^*(s, a)$**

$$\frac{\partial Q_{\theta}^*(s, a)}{\partial \theta} = \phi_{\theta}(s, a) \quad (28)$$

**Solving fixed point equation:**

$$\phi_{\theta}(s, a) = \frac{\partial R_{\theta}(s, a)}{\partial \theta} + \gamma \sum_{s' \in S} T(s, a, s') \sum_{a' \in A} \pi_{\theta}(a' \mid s') \phi_{\theta}(s', a') \quad (29)$$

**(iter)**

initialize  $\theta_0$

For  $k = 1$  to  $MaxIter$ :

Compute  $\frac{\partial L_{\theta}}{\partial \theta}$

$\theta_k = \theta_{k-1} + \sigma \frac{\partial L_{\theta}}{\partial \theta}$