

Goal: Maximise objective function

$$\begin{aligned} U(\theta) &= \mathbb{E}_{\pi_\theta}(\mathcal{R}(\tau)) \\ &= \sum_{\tau} P(\tau; \theta) \mathcal{R}(\tau) \end{aligned}$$

Use gradient ascent to optimise parameters.

$$\theta \leftarrow \theta + \alpha \cdot \nabla_{\theta}[U(\theta)]$$

$$\begin{aligned} \nabla_{\theta}[U(\theta)] &= \nabla_{\theta}[\mathbb{E}_{\pi_\theta}(\mathcal{R}(\tau))] \\ &= \nabla_{\theta}\left[\sum_{\tau} P(\tau; \theta) \mathcal{R}(\tau)\right] \\ &= \sum_{\tau} \nabla_{\theta}[P(\tau; \theta)] \mathcal{R}(\tau) \end{aligned}$$

Note we have introduced these as they cancel.

$$= \sum_{\tau} P(\tau; \theta) \cdot \frac{\nabla_{\theta}[P(\tau; \theta)]}{P(\tau; \theta)} \mathcal{R}(\tau)$$

$$= \sum_{\tau} P(\tau; \theta) \nabla_{\theta}[\log(P(\tau; \theta))] \mathcal{R}(\tau)$$

This is an expectation! i.e.

$$= \mathbb{E}_{\tau_\theta}(\nabla_{\theta}[\log(P(\tau; \theta))] \mathcal{R}(\tau))$$

We can use Monte Carlo sampling to approximate this:

$$\approx \frac{1}{m} \sum_{i=1}^m (\nabla_{\theta}[\log(P(\tau^i; \theta))] \mathcal{R}(\tau^i))$$

But for each trajectory $\tau^{(i)}$ we sample, how do we calculate

$$\nabla_{\theta} [\log(P(\tau^{(i)}; \theta))] R(\tau^{(i)})?$$

$R(\tau^{(i)})$ is easy - it's just the sum of rewards over trajectory i .

$$\begin{aligned} \nabla_{\theta} [\log(P(\tau; \theta))] &= \nabla_{\theta} [\log(\pi_{\theta}(a_0 | s_0) \cdot P(s_1 | s_0, a_0) \cdot \pi_{\theta}(a_1 | s_1) \cdot P(s_2 | s_1, a_1) \cdots)] \\ &= \nabla_{\theta} [\log(\prod_{t=0}^{T-1} \pi_{\theta}(a_t | s_t) \cdot P(s_{t+1} | s_t, a_t))] \\ &= \nabla_{\theta} \left[\sum_{t=0}^{T-1} (\log(\pi_{\theta}(a_t | s_t) \cdot P(s_{t+1} | s_t, a_t))) \right] \\ &= \nabla_{\theta} \left[\sum_{t=0}^{T-1} (\log(\pi_{\theta}(a_t | s_t))) + \sum_{t=0}^{T-1} (P(s_{t+1} | s_t, a_t)) \right] \\ &= \sum_{t=0}^{T-1} (\nabla_{\theta} [\log(\pi_{\theta}(a_t | s_t))]) + \sum_{t=0}^{T-1} (\underbrace{\nabla_{\theta} [P(s_{t+1} | s_t, a_t)]}_{\text{no dependence on } \theta \text{ so derivative is 0}}) \\ &= \sum_{t=0}^{T-1} (\nabla_{\theta} [\log(\pi_{\theta}(a_t | s_t))]) \end{aligned}$$

So

$$\nabla_{\theta} [u(\theta)] \approx \frac{1}{m} \sum_{i=1}^m \left(\sum_{t=0}^{T-1} (\nabla_{\theta} [\log(\pi_{\theta}(a_t^{(i)} | s_t^{(i)})]) R(\tau^{(i)}) \right)$$