

# Week 10 - RL3 - Blackboard 1

statistical weight  $P(y, \vec{x}) =$

$$(1) \langle R \rangle = \sum_{\vec{x} \text{ input}} \sum_{y \in \{0,1\} \text{ output}} \overset{\text{reward}}{R(y, \vec{x})} \cdot \underbrace{\pi_{\omega}(y | \vec{x}) \cdot P(\vec{x})}_{\substack{\text{policy} \\ \text{depends on } \omega}} \cdot \underbrace{1}_{y=0}$$
$$= \sum_{\vec{x}} P(\vec{x}) \left[ R(y=1, \vec{x}) \cdot g(\vec{\omega} \cdot \vec{x}) + R(y=0, \vec{x}) \cdot (1 - g(\vec{\omega} \cdot \vec{x})) \right]$$

take derivative and update (batch)

$$(2) \Delta \omega_j = d \cdot \frac{\partial}{\partial \omega_j} \langle R \rangle = d \sum_{\vec{x}} P(\vec{x}) \left[ \underbrace{R(1, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x})}_{y=1} - \underbrace{R(0, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x})}_{y=0} \right] \cdot x_j$$

online rule? need to make statistical weight explicit! need  $\pi(y | \vec{x}) \cdot P(\vec{x})!$

use  $\pi_{\omega}(y | \vec{x}) = g(\vec{\omega} \cdot \vec{x})$  for  $y=1$  and  $\pi_{\omega}(y | \vec{x}) = (1 - g(\vec{\omega} \cdot \vec{x}))$  for  $y=0$

$$\Delta \omega_j = d \cdot \frac{\partial}{\partial \omega_j} \langle R \rangle = d \sum_{\vec{x}} P(\vec{x}) \left[ \underbrace{\frac{\pi(y=1 | \vec{x})}{g(\vec{\omega} \cdot \vec{x})} \cdot R(1, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x})}_{y=1} - \underbrace{\frac{\pi(y=0 | \vec{x})}{1 - g(\vec{\omega} \cdot \vec{x})} \cdot R(0, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x})}_{y=0} \right] \cdot x_j$$

$$(3) \Delta \omega_j = " = d \sum_{\vec{x}} \sum_{y \in \{0,1\}} \underbrace{P(\vec{x}) \cdot \pi(y | \vec{x})}_{\text{statistical weight}} \cdot R(y, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x}) \left[ \frac{y}{g(\vec{\omega} \cdot \vec{x})} - \frac{1-y}{1-g(\vec{\omega} \cdot \vec{x})} \right] \cdot x_j$$

if condition

online: cut statistical weight  $\Rightarrow$  self-averaging over samples

$$(4) \Delta \omega_j = d R(y, \vec{x}) \cdot g'(\vec{\omega} \cdot \vec{x}) \cdot \left[ \frac{y}{g} - \frac{1-y}{1-g} \right] \cdot x_j$$

Wed 10

Blackboard 2

log-likelihood trick

copy (1)

$$\langle R \rangle = \sum_{\vec{x}} \sum_y R(y, \vec{x}) \pi_{\omega}(y | \vec{x}) \cdot P(\vec{x})$$

$$\Delta \omega_j = \alpha \frac{\partial}{\partial \omega_j} \langle R \rangle = \alpha \sum_{\vec{x}} \sum_y R(y, \vec{x}) P(\vec{x}) \frac{\pi_{\omega}(y | \vec{x})}{\pi_{\omega}(y | \vec{x})} \frac{\partial}{\partial \omega_j} \pi_{\omega}(y | \vec{x})$$

$$= \alpha \sum_{\vec{x}} \sum_y \underbrace{P(\vec{x}) \pi_{\omega}(y | \vec{x})}_{\text{statistical weight}} \cdot R(y, \vec{x}) \frac{\partial}{\partial \omega_j} \ln \pi_{\omega}(y | \vec{x})$$

outline

$$(5) \parallel \Delta \omega_j = \alpha \cdot \begin{matrix} R(y, \vec{x}) \\ \uparrow \\ \text{reward} \end{matrix} \frac{\partial}{\partial \omega_j} \ln \begin{matrix} \pi_{\omega}(y | \vec{x}) \\ \uparrow \\ \text{policy} \end{matrix} \parallel \text{"log-likelihood trick"}$$

evaluate for ow case: (Blackboard 2b/continued)

"if-condition"

$$\text{if } y=1 \quad \Pi_{\omega}(y=1|\vec{x}) = g(\vec{\omega} \cdot \vec{x})$$

$$\text{if } y=0 \quad \Pi_{\omega}(y=0|\vec{x}) = (1-g(\vec{\omega} \cdot \vec{x}))$$

$$\Pi_{\omega}(y|\vec{x}) = g^y \cdot (1-g)^{(1-y)}$$

$$\Rightarrow \ln \Pi_{\omega}(y|\vec{x}) = y \cdot \ln g + (1-y) \cdot \ln(1-g)$$

compare: log-likelihood

$$\Rightarrow \frac{\partial}{\partial \omega_j} \ln_{\omega} \Pi(y|\vec{x}) = \frac{y}{g} \cdot g' \cdot x_j - \frac{(1-y)}{(1-g)} \cdot g' \cdot x_j$$

with (5)

$$\Delta \omega_j = d \cdot R(y, \vec{x}) \cdot g' \left[ \frac{y}{g} - \frac{(1-y)}{1-g} \right] \cdot x_j \quad \text{compare (4)}$$

evaluate further:

$$\Delta \omega_j = d \cdot R(y, \vec{x}) \frac{g'}{g \cdot (1-g)} \left[ \cancel{(1-g)} \cdot y - g \cdot \cancel{(1-y)} \right] \cdot x_j$$

$$\Delta \omega_j = d \cdot R(y, \vec{x}) \frac{g'}{g \cdot (1-g)} [y - g] \cdot x_j$$

$$\uparrow \quad g = \langle y \rangle = 1 \cdot \text{Prob}(y=1) + 0 \cdot \text{Prob}(y=0)$$

$$\Delta \omega_j = d \cdot R(y, \vec{x}) \frac{g'}{g \cdot (1-g)} [y - \langle y \rangle] \cdot x_j$$

Week 10 - Blackboard 3: multi-step policy gradient

estimated return (total discounted future reward)

$$V_{\pi}^{\text{est}}(s_t) = \langle r_t + \gamma V_{\pi}^{\text{est}}(s_{t+1}) + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots \rangle$$

depends on policy

Bellman

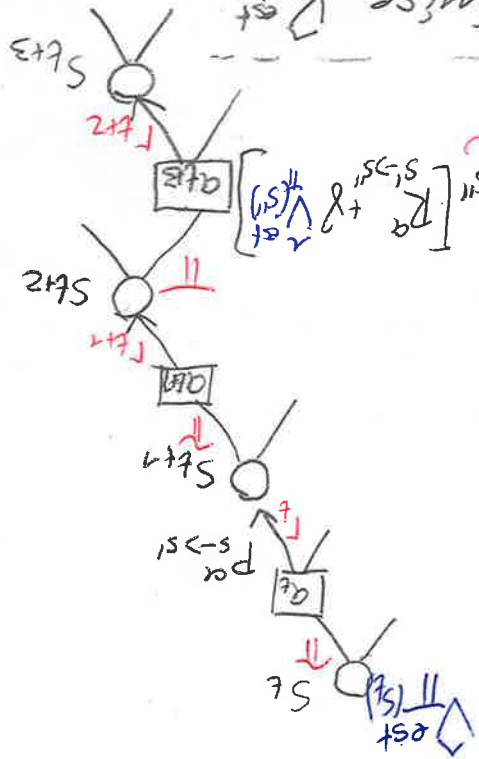
$$= \sum_{a_t} \pi(a_t | s_t) \cdot \sum_{s_{t+1}} p_{a_t}^{s_t \rightarrow s_{t+1}} [r_{a_t} + \gamma \cdot V_{\pi}^{\text{est}}(s_{t+1})]$$

natural statistical weight

expand

$$\sum_{a_{t+1}} \pi(a_{t+1} | s_t) \sum_{s_{t+1}} p_{a_{t+1}}^{s_t \rightarrow s_{t+1}} [r_{a_{t+1}} + \gamma \cdot V_{\pi}^{\text{est}}(s_{t+1})]$$

all paths from  $s_t$



Change parameters  $\theta$  of policy  $\pi_{\theta}(a|s)$  so as to maximise  $V_{\pi}^{\text{est}}(s_t)$  so as to maximise  $V_{\pi}^{\text{est}}(s_t)$

contains natural statistical weight

$$\Delta \theta = \alpha \cdot \frac{\partial}{\partial \theta} V_{\pi}^{\text{est}}(s_t) = \alpha \sum_{a_t} \pi(a_t | s_t) \cdot \frac{\partial}{\partial \theta} \ln \pi(a_t | s_t) \sum_{s_{t+1}} p_{a_t}^{s_t \rightarrow s_{t+1}} [r_{a_t} + \gamma \cdot V_{\pi}^{\text{est}}(s_{t+1})]$$

(product rule)

$$+ \alpha \sum_{a_t} \pi(a_t | s_t) - \sum_{s_{t+1}} p_{a_t}^{s_t \rightarrow s_{t+1}} \cdot \gamma \cdot \frac{\partial}{\partial \theta} V_{\pi}^{\text{est}}(s_{t+1})$$

online rule, drop statistical weight

expand iteratively

$$\Delta \theta = \alpha \cdot \frac{\partial}{\partial \theta} \ln \pi(a_t | s_t) [r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots]$$

$$+ \alpha \cdot \frac{\partial}{\partial \theta} \ln \pi(a_{t+1} | s_{t+1}) [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots]$$