$$\langle R(\vec{\omega}) \rangle = \underbrace{\sum_{\vec{x}}}_{\substack{input}} \underbrace{\sum_{y \in \{0,1\}}}_{\substack{output}} \overset{reward}{R(y,\vec{x})} \cdot \overbrace{\underset{\vec{\omega}}{\Pi} (y|\vec{x}) \cdot P(\vec{x})}^{\text{Statistical weight } P(y,\vec{x}) =}$$

$$\underbrace{}_{\text{sum all possibilities}} \qquad \underset{\substack{policy \ depends \\ on \ \vec{\omega}}}{\nearrow}$$

$$= \sum_{\vec{x}} P(\vec{x}) \left[ R(y=1,\vec{x}) \cdot g(\vec{\omega}\cdot\vec{x}) + R(y=0,\vec{x}) \cdot (1 - g(\vec{\omega}\cdot\vec{x})) \right]$$

derivative / batch update

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

This is the correct batch update!
But it is not easy to go from here to online,
because we do not have the correct statistical
weight :   we miss   $\sum_{y} P_{\vec{\omega}}(y|\vec{x})$

aim:   make statistical weight

$$\sum_{\vec{x}} \sum_{y} P_{\vec{\omega}}(y|\vec{x}) \cdot P(\vec{x})$$

visible!   Must show up explicitly!

approach ①   " pedestrian "

approach ②   " log-likelihood trick "

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

make "if-condition" explicit

$$(\text{x}) \quad \langle R \rangle = \sum_{\substack{\vec{x} \\ \text{input}}} \sum_{\substack{y \in \{0,1\} \\ \text{output}}} \overbrace{R(y,\vec{x})}^{\text{reward}} \cdot \underbrace{\Pi_\omega(y|\vec{x})}_{\substack{\text{policy} \\ \text{depends on } \omega}} \cdot P(\vec{x})$$

$y=1$     if $y=0$

$$= \sum_{\vec{x}} \sum_y P(\vec{x}) \left[ R(y=1,\vec{x}) \cdot g(\vec{\omega}\cdot\vec{x}) \cdot y + R(y=0,\vec{x}) \cdot (1-g(\vec{\omega}\cdot\vec{x})) \cdot (1-y) \right]$$

if I add the sum over $y$, I need to add the "if-condition"!

take derivative **and update (batch)**

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

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

use (1) $\Pi_\omega(y|\vec{x}) = g(\vec{\omega}\cdot\vec{x})$ for $y=1$ and (2) $\Pi_\omega(y|\vec{x}) = (1-g(\vec{\omega}\cdot\vec{x}))$ for $y=0$

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

sum all possibilities       if condition

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

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

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

Week 4: — Blackboard 2    log - likelihood trick

copy (1) $\qquad \langle R \rangle = \sum_{\vec{x}} \sum_{y} R(y, \vec{x}) \, \Pi_\omega(y \mid \vec{x}) \cdot P(\vec{x})$

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

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

online

$$(5) \left\| \; \Delta \omega_j = \alpha \cdot \underset{\underset{\text{reward}}{\uparrow}}{R(y, \vec{x})} \frac{\partial}{\partial \omega_j} \ln \underset{\underset{\text{policy}}{\uparrow}}{\Pi_\omega(y \mid \vec{x})} \; \right\|$$

"log - likelihood trick"

Week 4 – Blackboard 3 :   multi-step policy gradient
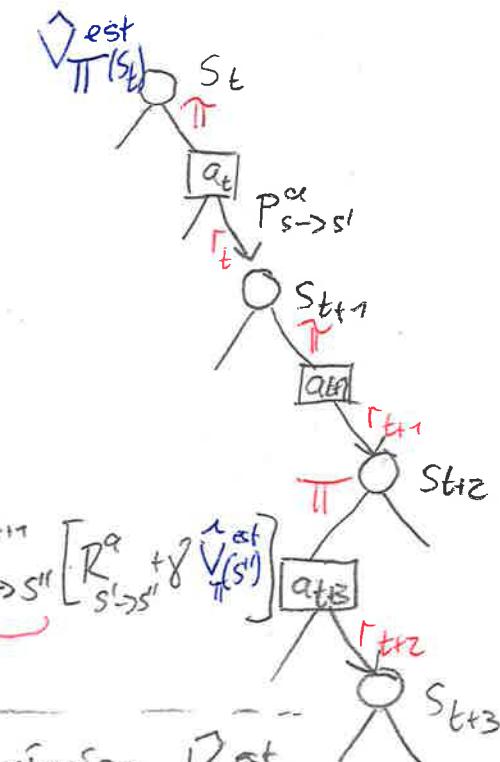
estimated return ( total discounted future reward )

$$\hat{V}_{\Pi}^{est}(s_t) = \left\langle r_t + \gamma\, r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \ldots \right\rangle_{\substack{\text{all paths} \\ \text{from } s_t}}$$

→ depends on policy

Bellman

$$= \sum_{a_t} \Pi_\theta(a_t|s_t) \cdot \sum_{s'} P_{s_t \to s'}^{a_t} \left[ R_{s_t \to s'}^{a} + \gamma \cdot \hat{V}_{\Pi}^{est}(s') \right]$$

natural statistical weight

expand

$$\sum_{a_{t+1}} \Pi_\theta(a_{t+1}|s') \sum_{s''} P_{s' \to s''}^{a_{t+1}} \left[ R_{s' \to s''}^{a} + \gamma\, \hat{V}_{\Pi}^{a_t}(s'') \right]$$



---

Change parameters $\theta$ of policy $\Pi_\theta(a|s)$  so as to maximise $\hat{V}_{\Pi}^{est}(s_t)$

$$\Delta\theta = \alpha \frac{\partial}{\partial\theta} \hat{V}_{\Pi}^{est}(s_t) = \alpha \sum_{a_t} \Pi_\theta(a_t|s_t) \cdot \frac{\partial}{\partial\theta} \ln\Pi(a_t|s_t) \sum_{s'} P_{s_t \to s'}^{a_t} \left[ R_{s_t \to s'}^{a} + \gamma\, \hat{V}_{\Pi}^{est}(s') \right]$$

(product rule)

contains natural statistical weight

↳ also depends on $\Pi$

$$+ \alpha \sum_{a_t} \Pi_\theta(a_t|s_t) \cdot \sum_{s'} P_{s_t \to s'}^{a_t} \cdot \gamma \cdot \frac{\partial}{\partial\theta} \hat{V}_{\Pi}^{est}(s')$$

online rule, drop statistical weight

expand iteratively

$$\Delta\theta = \alpha \cdot \frac{\partial}{\partial\theta} \ln\Pi(a_t|s_t) \underbrace{\left[ r_t + \gamma\, r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \ldots \right]}_{} \quad R_{s_t \to s_{end}}^{a_t} \text{ "Return"}$$

$$+ \alpha \cdot \gamma \cdot \frac{\partial}{\partial\theta} \ln\Pi(a_{t+1}|s_{t+1}) \underbrace{\left[ r_{t+1} + \gamma\, r_{t+2} + \gamma^2 r_{t+3} + \ldots \right]}_{} \quad R_{s_{t+1} \to s_{end}}^{a_{t+1}}$$

$$+ \alpha \cdot \gamma^2 \ldots$$