

CB Regret

$$\max_{\pi} \sum_t \mathbb{E}_{a \sim \pi(q)} r(a) - \sum_t r(a_t)$$

RL Regret:

$R_t = \sum_h r(a_h)$ in t^{th} episode

$$\text{Regret} = \max_{\pi} \sum_t \mathbb{E}_{\tau \sim \pi} \left[\sum_h r_h \right] - \sum_t R_t$$

\uparrow
 $\langle q_1, a_1, r_1, q_2, \dots, q_{H-1}, r_H \rangle$

$$\mathbb{E}_{q_1 \sim \text{world}} \mathbb{E}_{a_1 \sim \pi(q_1)} \mathbb{E}_{r_1, q_2 \sim \text{world}} \mathbb{E}_{a_2 \sim \pi(q_2)} \dots \mathbb{E}_{r_H} \left[\sum_h r_h \right]$$

Policy gradient (for $\pi(B)$)

$$\pi(a|o) = \frac{\exp[g(o,a)]}{\sum_{a'} \exp[g(o,a')]}$$

Goal: $\max_{\pi} \underbrace{\mathbb{E}_{o,r} \mathbb{E}_{a \sim \pi(o)} r(a)}_{R(\pi)}$

attempt: $\nabla_{\pi} R(\pi)$

$$R(\pi) = \mathbb{E}_{o,r} \sum_a \pi(a|o) r(a)$$

∇_{π}

$$\frac{\partial \log f(x)}{\partial x} = \frac{\partial f}{\partial x} \cdot \frac{1}{f(x)}$$

$$\Leftrightarrow \frac{\partial f}{\partial x} = f(x) \frac{\partial \log f(x)}{\partial x}$$

$$\begin{aligned}
 \nabla_{\pi} R(\pi) &= \mathbb{E}_{o,r} \sum_a [\nabla_{\pi} \pi(a|o)] r(a) \\
 &= \mathbb{E}_{o,r} \sum_a [\pi(a|o) \nabla_{\pi} \log \pi(a|o)] \times r(a) \\
 &= \mathbb{E}_{o,r} \mathbb{E}_{a \sim \pi} \underline{r(a)} \nabla_{\pi} \log \pi(a|o)
 \end{aligned}$$

Monte-Carlo estimate

~~$\mathbb{E}_{(o,r)} \sum_a \pi(a|o) \nabla_{\pi} \log \pi(a|o) r(a)$~~

$$\approx \sum_{(o,r)} \underbrace{r(\pi(o))}_{\text{sample from } \pi(a)} \nabla_{\pi} \log \pi(o|o)$$