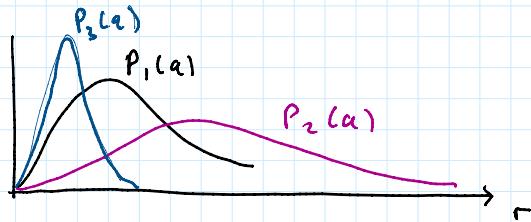


One-armed Bandits

choose $a \in \{1, \dots, A\} =: \mathcal{A}$



goal: what policy $\pi(a)$ over actions
maximises reward?

$$Q(a) := \mathbb{E}[R | A=a] = \mathbb{E}^{P_a}[R]$$

- greedy approach:

$$a^* = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(a)$$

$$\begin{matrix} R^{(1)}, & R^{(2)}, & R^{(3)}, & \dots, & R^{(n)} \\ a^{(1)}, & a^{(2)}, & a^{(3)}, & \dots, & a^{(n)} \end{matrix}$$

1.0	1.1	-5	+2	+1	10	-5	~1
1	1	3	2	2	3	3	2
							1

$$Q(1) \sim \frac{1}{2}(1.0 + 1.1) = 1.05 \leftarrow$$

$$Q(2) \sim \frac{1}{2}(2 + 1) = 1.5$$

$$\hookrightarrow Q(3) = \frac{1}{3}(2 + 1 - 1) = \frac{2}{3} \leftarrow$$

$$Q(3) \sim \frac{1}{3}(-5 + 10 - 5) = 0 \leftarrow$$

$$Q(a) = \frac{1}{n} \sum_{k=1}^n R^{(k)} \mathbb{1}_{a^{(k)}} \dots$$

$$Q(a) = \frac{1}{N(a)} \sum_{k=1}^N R^{(k)} \mathbb{1}_{a^{(k)}=a}$$

$$\hookrightarrow := \sum_{k=1}^N \mathbb{1}_{a^{(k)}=a}$$

$$a^* = \arg \max_a Q(a) \quad \text{apply action update}$$

- ϵ -greedy.

$H^{(n)}$ - Bernoulli r.v. success prob $\epsilon_k \in (0, 1)$

$$\epsilon_1 > \epsilon_2 > \epsilon_3 > \dots$$

$$\epsilon_n \xrightarrow[n \rightarrow \infty]{} 0$$

$$\text{e.g. } \epsilon_n = 1/n$$

$$\epsilon_n = \frac{C}{D+n}$$

\hookrightarrow r.v. uniform on $\{1, \dots, A\}$

$$a^{(n)} = \mathbb{1}(H^{(n)} = 1) A^{(n)}$$

$$+ \mathbb{1}(H^{(n)} = 0) \arg \max_a Q^{(n)}(a)$$

$$\text{suppose } Q^{(0)}(1) = Q^{(0)}(2)$$

$$a^{(n)} \rightarrow R^{(n)} \rightarrow \text{update } Q(a^{(n)})$$



$$M^{(n+1)} = \frac{1}{n+1} \sum_{k=1}^{n+1} Y^{(k)} = \frac{1}{n+1} \left(\sum_{k=1}^n Y^{(k)} + Y^{(n+1)} \right)$$

$\overbrace{Y^{(1)}, \dots, Y^{(n)}}^{\text{y}}, \overbrace{Y^{(n+1)}}^{\text{y}}$

$$= \frac{1}{n+1} (M^{(n)} n + Y^{(n+1)})$$

$$M^{(n)} = \frac{1}{n} \sum_{k=1}^n Y^{(k)}$$

$$= M^{(n)} \left(\frac{n}{n+1} \right) + \frac{Y^{(n+1)}}{n+1}$$

$\hookrightarrow 1 - \frac{1}{n+1}$



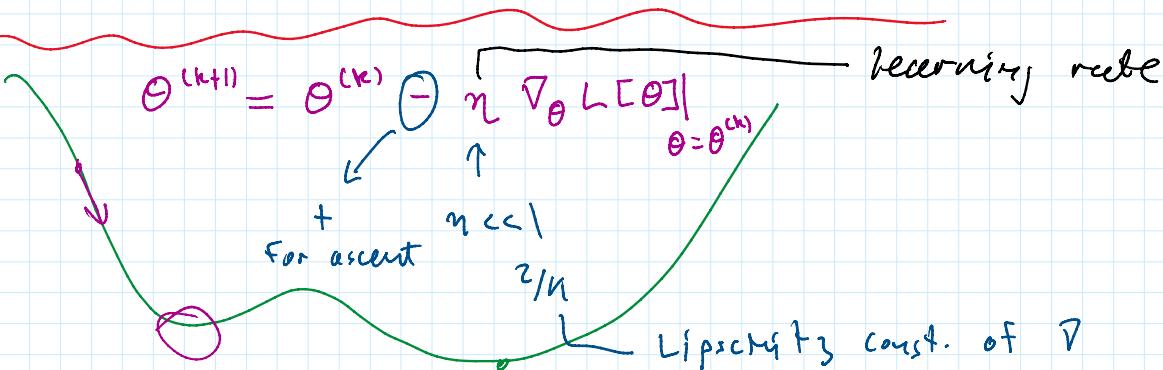
$$\mu^{(n+1)} = \mu^{(n)} + \frac{1}{n+1} (y^{(n+1)} - \mu^{(n)})$$

$$1 - \frac{1}{n+1}$$

$\pi_\theta(a)$ - distribution over actions parameterised by θ

$$Q(a) = \mathbb{E}_{P_a}^{\pi}[R]$$

$$Q[\theta] = \int Q(a) \pi_\theta(a) da =: \mathbb{E}_{P^\pi}^{\pi}[R]$$



$$Q[\theta] = \int Q(a) \pi_\theta(a) da = \mathbb{E}_{P^\pi}^{\pi}[R]$$

$$\text{e.g. } \pi_\theta(a) = p \phi\left(\frac{a-a_0}{\sigma_0}\right) + (1-p) \phi\left(\frac{a-a_1}{\sigma_1}\right)$$

$$a^{(1)} \sim \pi_\theta \rightarrow r^{(1)}$$

$$a^{(2)} \sim \pi_\theta \rightarrow r^{(2)}$$

$$\vdots \quad \vdots \quad \vdots$$

$$\frac{1}{N} \sum_{n=1}^N r^{(n)} \sim Q[\theta] \text{ unbiased estimator of}$$

REINFORCE

$$\nabla_\theta Q[\theta] = \int Q(a) \nabla_\theta \pi_\theta(a) da$$

$\dots P_\theta \dots$

$$\begin{aligned}
\nabla_{\theta} Q[\theta] &= \int Q(a) \nabla_{\theta} \pi_{\theta}(a) da \\
&= \int Q(a) \left(\frac{\nabla_{\theta} \pi_{\theta}(a)}{\pi_{\theta}(a)} \right) \cdot \pi_{\theta}(a) da \quad \text{likelihood ratio "trick"} \\
&\quad \xrightarrow{\text{IP}_{\theta}[R]} \nabla_{\theta} \log \pi_{\theta}(a) \\
&= \iint r \nabla_{\theta} \log \pi_{\theta}(a) \underbrace{p_a(r) \pi_{\theta}(a) dr da}_{\text{IP}(r|a) \text{IP}_{\theta}(a) = \text{IP}(r)} \\
&= \mathbb{E}^{\text{IP}_{\theta}}[R \nabla_{\theta} \log \pi_{\theta}(a)]_{a \sim A} \\
&\boxed{\nabla_{\theta} Q[\theta] = \mathbb{E}^{\text{IP}_{\theta}}[R \nabla_{\theta} \log \pi_{\theta}(A)]}
\end{aligned}$$

$$a^{(n)} \sim \pi_{\theta}(a) \rightarrow r^{(n)} \leftarrow$$

$$\begin{aligned}
\nabla_{\theta} Q[\theta] &\sim \frac{1}{N} \sum_{k=1}^N r^{(k)} \nabla_{\theta} \log \pi_{\theta}(a^{(k)}) \\
\pi_{\theta}(a) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2} \frac{(a-\bar{a})^2}{\sigma^2}} \quad a \sim N(\bar{a}; \sigma^2) \\
\log \pi_{\theta}(a) &= -\frac{1}{2} \left[\frac{(a-\bar{a})^2}{\sigma^2} + \log(2\pi\sigma) \right] \quad \eta := \sigma^{-2} \\
\partial_{\bar{a}} \log \pi_{\theta}(a) &= \frac{(a-\bar{a})}{\sigma} \\
\partial_{\eta} \log \pi_{\theta}(a) &= -\frac{1}{2} \left[(a-\bar{a})^2 - \frac{1}{2\eta} \right]
\end{aligned}$$

overall loss: $\text{IP}(Y=1 | X=x) = \rho^{\beta' x}$

recall logistic regression model

$$P(Y=1 | X=x) = \frac{e^{\beta' x}}{1 + e^{\beta' x}}$$

multi-class logistic regression

$$y, Y \in \{0, \dots, K-1\}$$

$$P(Y=y | X=x) = \frac{e^{\beta_y' x}}{\sum_{c=0}^{K-1} e^{\beta_c' x}}$$

$$\text{a canonical choice } \beta_0 = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\text{so take } \pi_\theta(a) = \frac{e^{H_\theta(a)}}{\sum_{a'} e^{H_\theta(a')}} \quad (\text{soft-max})$$

$$\nabla_\theta \log \pi_\theta(a) = \nabla_\theta H_\theta(a) - \nabla_\theta \log \sum_{a'} e^{H_\theta(a')}$$

$$= \nabla_\theta H_\theta(a) - \underbrace{\nabla_\theta \sum_{a'} e^{H_\theta(a')}}_{\sum_{a'} e^{H_\theta(a')}}$$

$$= \nabla_\theta H_\theta(a) - \underbrace{\sum_{a''} e^{H_\theta(a'')} \nabla_\theta H_\theta(a'')}_{\sum_{a'} e^{H_\theta(a')}}$$

$$H_\theta : A \rightarrow \mathbb{R}$$

$$H_\theta(a)$$

=

$$0 \longrightarrow H_\theta(0) = \gamma_0$$

$$1 \rightarrow H_\theta(1) = r_1$$

$$2 \rightarrow H_\theta(2) = r_2$$

$$\begin{aligned} \partial_{\gamma_k} \log \pi_\theta(a) &= \mathbb{1}_{a=k} - \frac{\sum_{a''} e^{H_\theta(a'')} \mathbb{1}_{a''=k}}{\sum_{a'} e^{H_\theta(a')}} \\ &= \mathbb{1}_{a=k} - \left(\frac{e^{H_\theta(k)}}{\sum_{a'} e^{H_\theta(a')}} \right) \end{aligned}$$

$$\boxed{\partial_{\gamma_k} \log \pi_\theta(a) = \mathbb{1}_{a=k} - \pi_\theta(k)}$$

$$\boxed{\partial_{\gamma_k} Q[\theta] = \mathbb{E}_{\pi}^{\theta} [R (\mathbb{1}_{A=k} - \pi_\theta(k))]}$$

$$\gamma_k^{(n)} \rightarrow \gamma_k^{(n)} + n \underbrace{\frac{1}{M} \sum_{m=1}^M r^{(m)} (\mathbb{1}_{a^{(m)}=k} - \pi_\theta^{(n)}(k))}_{\gamma_k^{(n)}}$$

$$\gamma_k^{(1)} = r^{(1)} (\mathbb{1}_{a^{(1)}=k} - \pi_\theta^{(1)}(k))$$

$$\gamma_k^{(n+1)} = \gamma_k^{(n)} + \frac{1}{n+1} \underbrace{(r^{(n+1)} (\mathbb{1}_{a^{(n+1)}=k} - \pi_\theta^{(n)}(k)) - \gamma_k^{(n)})}_{\gamma_k^{(n)}}$$

