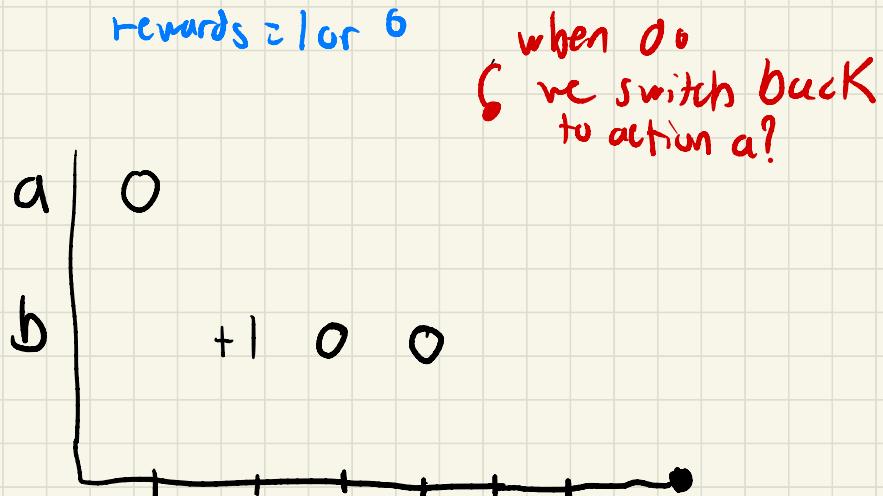


Settings:

- Env. has a single state
 - ↳ actions no longer have long-term consequences
 - ↳ actions gives immediate reward
 - ↳ can ignore other observations

How to learn policy in this setting?

Example:



- Learning agents must balance between 2 things

Exploitation := Maximise performance based on current knowledge

Exploration := Increase Knowledge

Gather info. to make better decisions

- long term strategy may involve short term sacrifice

Formalizing: Multi-Armed Bandit

- Multi-Armed Bandit is a set of distributions $\{P_a | a \in A\}$

P_a distribution of rewards
known actions

- Each time step t , agent selects $A_t \in A$

- Env. generates $R_t \sim P_{A_t}$

- goal = maximize $\sum_{i=1}^t R_i$ → do this by learning policy

Action value: $q(a) = E[R_t | A_t = a]$

Optimal value: $V_\star = \max_{a \in A} q(a) = \max_a E[R_t | A_t = a]$

regret : $\Delta_a = v_a - q(a)$

Want to minimise total regret

$$L_t = \sum_{n=1}^t v_{A_n} - q(A_n) = \sum_{n=1}^t \Delta_{A_n}$$

Maximise cumulative reward \equiv minimise total regret

Algorithms

several:

- 1) Greedy
- 2) ϵ -Greedy
- 3) UCB
- 4) Thompson Sampling
- 5) Policy gradients

→ use action value estimates
 $Q_t(a) \approx q(a)$

↳ estimate

Simple one:

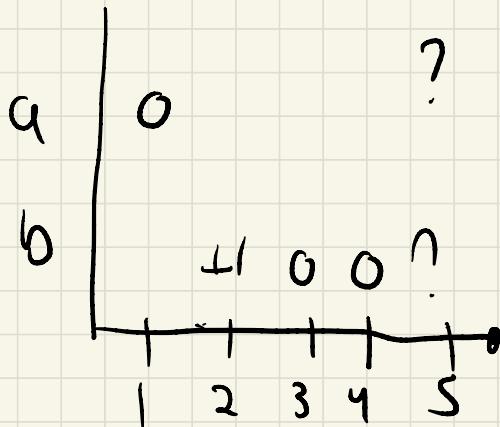
$$Q_t(a) = \frac{\sum_{n=1}^t I(A_n=a) R_n}{\sum_{n=1}^t I(A_n=a)}$$

$I(\text{true})=1$
 $I(\text{false})=0$ Leverage

1) Greedy

Select action with highest value: $A_t = \arg \max_a Q_t(a)$

regret:



$$Q_4(b) = \frac{1}{3}$$

Greedy would keep choosing action b

$$P(+1|a) = 0.8 = q(a)$$

$$P(+1|b) = 0.2 = q(b)$$

$$\Delta a = V_{\text{opt}} - q(a)$$

$$\Delta a = 0$$

$$\Delta b = 0.6$$

?

Q) ϵ -greedy

-Greedy gets stuck on suboptimal action forever
⇒ linear expected total regret.

Comes ϵ -greedy!

- with probability $1-\epsilon$ select optimal
- (1) ϵ select random action

$$\pi_+(u) = \begin{cases} (1-\epsilon) + \epsilon / |A| & \text{if } Q_+(u) = \max_b Q_b(u) \\ \epsilon / |A| & \text{otherwise} \end{cases}$$

Continues to explore, still linear expected total regret

Policy Search \rightarrow (can we learn $\pi(a)$ directly, instead of learning values)

Action preferences $H_t(a)$ and policy

$$\pi(a) = \frac{e^{H_t(a)}}{\sum_b e^{H_t(b)}}$$

Softmax

\hookrightarrow goal: Optimize preferences

Preferences := values; learnable policy parameters

Use gradient ascent

$$\Theta_{t+1} = \Theta_t + \alpha \nabla_\Theta E[R_t | \pi_{\Theta_t}]$$

vector

$\Theta_t \rightarrow$ current policy parameters

How to compute gradient

Log-likelihood trick (see slides)

\hookrightarrow can sample this

$$\nabla_\Theta E[R_t | \Theta] = E[R_t \nabla_\Theta \log \pi_\Theta(A_t)]$$

$$\theta = \theta + \alpha R_t \nabla_{\theta} \log \pi_{\theta}(A_t)$$

↳ stochastic gradient ascent

Note:

$$\sum_a \pi_{\theta}(a) = 1 \Rightarrow b$$

$$\sum_a b \nabla_{\theta} \pi_{\theta}(a) = b \nabla_{\theta} \sum_a \pi_{\theta}(a) = 0$$

as long as b does not depend on θ or a ,

We can subtract a baseline and use:

$$\theta = \theta + \alpha (R_t - b) \nabla_{\theta} \log \pi_{\theta}(A_t)$$

↳ changes variance.

Gradient can be stuck on local optimum

Theory! What is possible?

Theorem

$$\lim_{t \rightarrow \infty} L_t \geq \sum_{a \in A} \frac{D_a}{KL(Q_{\text{all}} \| R_a)} \propto D_a^2$$

So regret grows at least logarithmically.

Can we hit that lower bound?
Logarithmic upper bound!

$$D_a = \sqrt{Q_a - q(a)}$$

Total regret depends on D_a and action count

$$L_t = \sum_{n=1}^t D_{A_n} = \sum_{a \in A} N_t(a) D_a$$

Good algorithms ensure small counts for large action
regrets

UCB

- Estimate upper confidence $U_+(a)$ & action value
 - | $a(a) \leq Q_+(a) + U_+(a)$ with high probability
- Select $\underset{a \in A}{\operatorname{argmax}}^{\text{action}}$ upper confidence bound (ucb)
$$a_t = \underset{a \in A}{\operatorname{argmax}} Q_+(a) + U_+(a)$$

↳ uncertainty depends on $N(a)$

Large = high $U_+(a)$
Large = small $U_+(a)$

a only selected if
 - or $Q_+(a)$ large
 - $U_+(a)$ large

Hoeffding's inequality → See slides.

UCB derivation

Note: $\Delta a \geq 0$

$$L_t = \sum_a N_t(a) \Delta a$$

only random acty

$$x_a \neq a^* \quad N_t(a) \Delta a \leq x_a \log t$$

$$m \leq t \quad N_m(a) \Delta a \leq x_a \log m \leq x_a \log t$$

Let us possible time step which this holds.

$$n \in \{m+1, \dots, t\} : N_n(a) \Delta a \geq x_a \log n$$

$$\mathbb{E}[N_t(a)] = \mathbb{E}\left[\sum_{n=1}^t \mathbb{I}(A_n=a)\right]$$

$$= \mathbb{E}\left[N_m(a) + \sum_{n=m+1}^t \mathbb{I}(A_n=a)\right]$$

$$= \mathbb{E}\left[x_a \cdot \frac{\log t}{\Delta a} + \sum_{n=m+1}^t \mathbb{I}(A_n=a)\right]$$

$$= x_a \frac{\log t}{\Delta a} + \sum_{n=m+1}^t \mathbb{E}[\mathbb{I}(A_n=a) | N_n(a) \Delta a \geq x_a \log n]$$

$$= x_a \frac{\log t}{\Delta a} + \sum_{n=m+1}^t P(A_n=a, \dots)$$

$P(A_n=a)$? Given $N_m(a) \Delta a \geq x_a \log n$ ①

$$A_t = \arg \max_a Q(a) + \sqrt{\frac{2 \log t}{N_t(a)}}$$

$$\leq P(Q+U) + U_f(u) \geq Q^* + u^* + U_f(u^*)$$

$$= P(Q+U \geq Q^* + u^*) \stackrel{?}{\leq} 1$$

$$= P(Q+U \geq Q^* + u^* \mid Q^* + u^* \leq y) \cdot P(Q^* + u^* \leq y) \stackrel{st}{\leq}$$

$$+ P(Q+U \geq Q^* + u^* \mid Q^* + u^* > y) \cdot P(Q^* + u^* > y)$$

$$P(Q^* + u^* \leq q^*) \leq e^{-N_n(q^*) U_n(q^*)^2}$$

" $U(q^*)$ " "missing σ^2 " pick $y = q(q^*)$

$$e^{-N_n^2} = \frac{1}{n} \Rightarrow n \geq \sqrt{\frac{\log n}{N}}$$

$\sum_{n=1}^{\infty} \frac{1}{n} < \log t + 1$

Now bound:

$$P(Q+U \geq Q^* + u^* \mid Q^* + u^* > q^*) \leq \frac{1}{n}?$$

$$\leq P(Q+U \geq q^*) = P(-Q-U \leq -q^* + q - q)$$

$= -\Delta q$

$$= P(-Q + \underline{\Delta q} - U \leq q)$$

bound this

(1): $N \cdot \Delta q \geq \log n$

$\Delta q \geq \frac{\log n}{N_n(u)}$

$\underline{= x_u(u+\Delta q)^2}$

if we pick $x_a \geq \frac{1}{\Delta a}$ $\Rightarrow \Delta a^2 / (a + \Delta a)^2$
 $\Delta a / (a + \Delta a)$!

PICK $x_0 = \frac{4}{\partial u} \Rightarrow \Delta u > 2u$ now

$$\Rightarrow P(-Q + Du - u \leq -a) \leq P(-Q + 2u - u \leq -a)$$

$$= P(-Q + u \leq -a) \leq e^{-N \cdot u^2} \quad | \quad u = \sqrt{\frac{\log n}{N}}$$

$$\Rightarrow D(A_n - u) \leq \frac{2}{n}$$

$$\sum_{n=m}^{+\infty} P(A_n \cap A) \leq 2(\log t + 1)$$