

Stochastic Multi-Armed Bandit

$m \rightarrow$ No. of arms

$A = \{a_1, a_2, \dots, a_m\}$ Arms

Bernoulli (0,1) rewards

Mean reward of arm a_i $\theta = p_i \in [0, 1]$

T: Horizon/budget

Highest mean $= p_*$; optimal arm a_*

Algorithm

a^0, r^0 $r^0 \in \text{Bernoulli}(\theta^0)$

a^1, r^1 $r^1 \in \text{Bernoulli}(p_a)$

a^{t-1}, r^{t-1} | Which arm to pull at t ?

An algorithm is a mapping from the set of histories to the set of arms,

Probability distribution over arms (A)
(randomized)

Deterministic.

① ϵ -first $\epsilon \in [0, 1]$

Pull arms uniformly at random (or round robin) for ET steps.

- Identify a_{best} an arm with the highest empirical mean.
- For the remaining pulls, sample a_{best}

ϵ -greedy

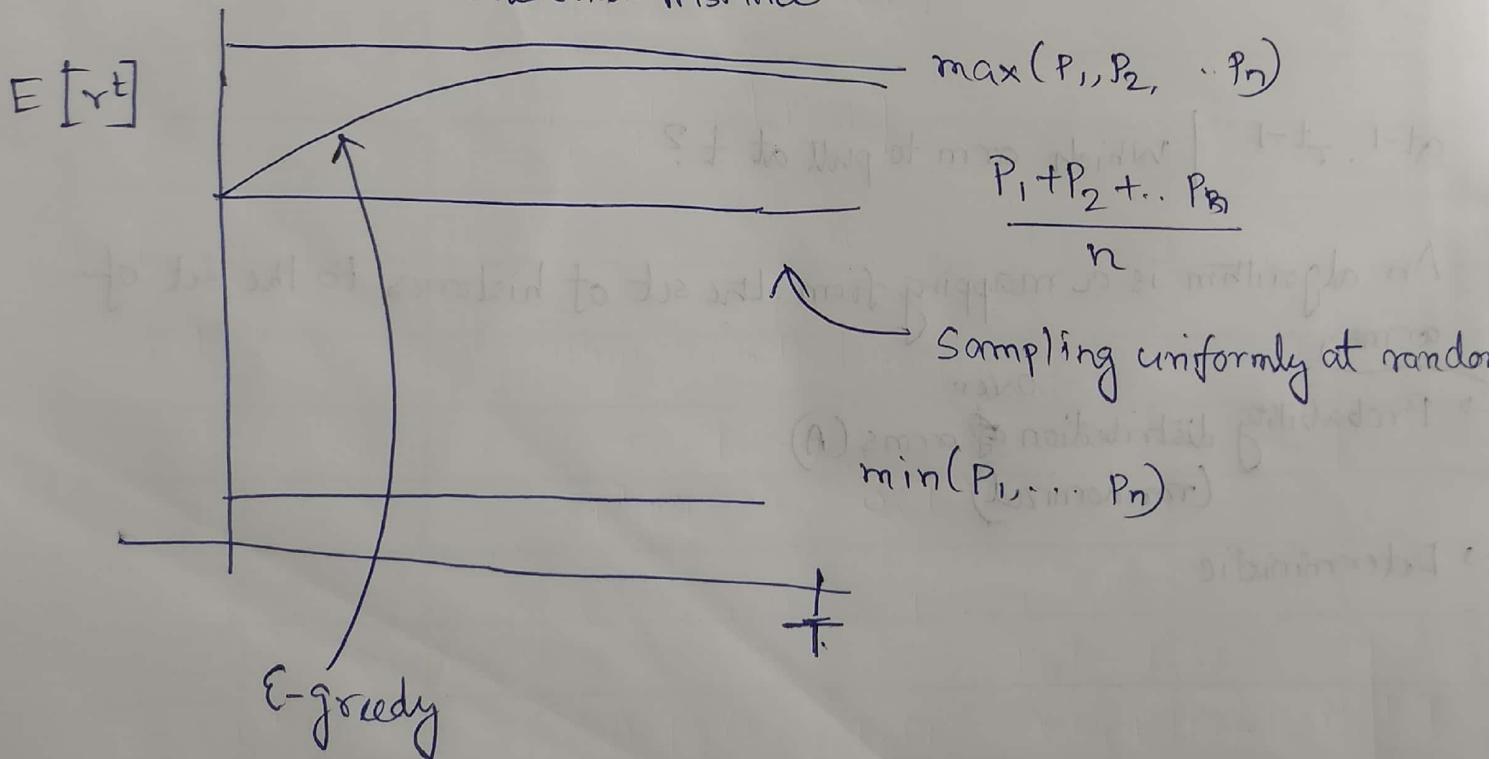
For $t \in 0, 1, 2, \dots$

With probability ϵ , pick an arm uniformly at random to pull.

With probability $1-\epsilon$, pull the arm with the highest empirical mean.

Graph

Some Given instance

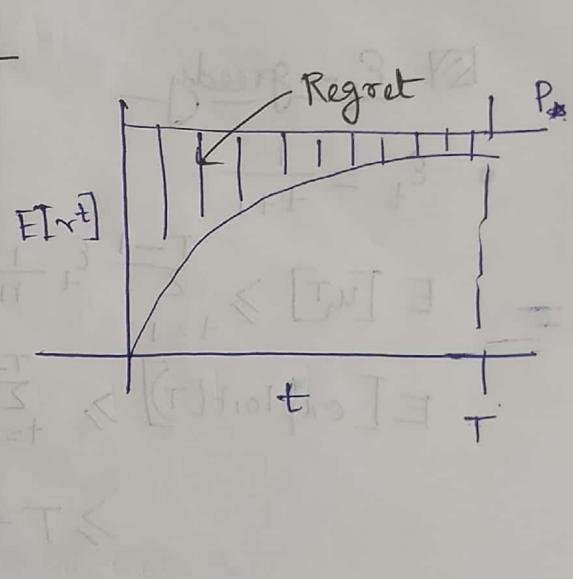


Will the ϵ -greedy line meet the max line at $t = \infty$?

$$E(r_\infty) = \epsilon \left(\frac{P_1 + P_2 + \dots + P_n}{n} \right) + (1-\epsilon) \cdot \max(P_1, P_2, \dots, P_n)$$

Expected cumulative vs. cumulative regret:

$$R^T = P_* T - \sum_{t=0}^{T-1} E[r_t]$$



IV) Is it possible to achieve a sub-linear regret: $\circ(T)$

Can we design an algorithm for which $R^T = \Theta(T)$ or equivalently.

$$\lim_{T \rightarrow \infty} \frac{R^T}{T} = 0$$

Let u_a^T denote the number of pulls of arm a in T pulls.

Let $\text{exploit}(T)$ denote the number of times in T pulls that an arm with the highest empirical average get pulled.

For sub-linear regret we need:

$$\forall a \in A, \lim_{T \rightarrow \infty} E[u_a^T] = 0 \quad [\text{Doing Infinite Exploration}]$$

$$\lim_{T \rightarrow \infty} \frac{E[\text{exploit}(T)]}{T} = 1. \quad [\text{Greedy in the Limit}]$$

◻ $\varepsilon_t = \text{greedy}$

$$\bullet \quad \varepsilon_t = \frac{1}{t+1} \quad E[u_a^T] \geq \sum_{t=1}^{T-1} \varepsilon_t + \frac{1}{n} = \frac{1}{n} \sum_{t=0}^{T-1} \frac{1}{t+1} \geq c_1 \log T.$$

$$E[\text{exploit}(T)] \geq \sum_{t=0}^{T-1} (1 - \varepsilon_t) + \frac{\varepsilon_t}{n} \geq T = c_2 \log T$$

$$\bullet \quad \varepsilon_t = e^{-t}$$

$$E[u_a^T] \geq \sum_{t=0}^{T-1} \frac{e^{-t}}{n}$$

$$= \frac{1}{n} (1 + e^{-1} + e^{-2} + e^{-3} + \dots)$$

$$= \frac{1}{n} \frac{1}{1-e^{-1}} \quad [\text{Finite}]$$

$$= \frac{1}{n} \frac{1}{1-e^{-1}} \quad [\text{Finite}]$$

4. Lower Bound (Lai and Robbins, 1985)

Let L be an algorithm such that for every bandit instance I and for every $\alpha > 0$, as $T \rightarrow \infty$

$$R^T(L, I) = o(T^\alpha)$$

Then, for every bandit instance I as $T \rightarrow \infty$

$$R^T(L, I) \geq \left[\sum_{a \in A: P_a(I) \neq P_\star(I)} \frac{P_\star(I) - P_a(I)}{KL(P_a(I), P_\star(I))} \right] \ln T$$

where

$$KL(x, y) = x \ln \frac{x}{y} + (1-x) \ln \frac{1-x}{1-y}$$

■ UCB Algorithm

- Sample every arm at least once.

- For $t = n, n+1, \dots$

Let \hat{p}_a^t denote the empirical mean of arm a

Let u_a^t be a 's number of pulls

$$\text{Define } ucb_a^t = \hat{p}_a^t + \sqrt{\frac{2}{u_a^t} \ln t}$$

Pull $\operatorname{argmax}_{a \in A} ucb_a^t$

$$R^T = O\left(\left(\sum_{\substack{a \\ p_a \neq p_*}} \frac{1}{p_* - p_a}\right) \log T\right) \text{ for UCB}$$

■ Hoeffding's Inequality (1963)

Let x be a random variable that takes values in $[a, b]$, with $E[x] = p$. Fix $u > 1$ and $\varepsilon > 0$. Let r^1, r^2, \dots, r^n be samples drawn i.i.d. from x and let $\hat{p} = \frac{\sum_{i=1}^n r^i}{n}$.

Then:

$$P\{\hat{p} > p + \varepsilon\} \leq e^{-\frac{2u\varepsilon^2}{(b-a)^2}}$$
 and $P\{\hat{p} \leq p - \varepsilon\} \leq e^{-\frac{2u\varepsilon^2}{(b-a)^2}}$

③ KL-UCB (2011)

Run UCB as above, but use:

$\text{act}_a^t \stackrel{\text{def}}{=} \text{the largest element } q \in [\hat{p}_a^t, 1] \text{ that satisfies}$

$$\text{KL}(\hat{p}_a^t, q) \leq \frac{\text{Int} + 3\ln \text{Int}}{\hat{p}_a^t}$$

$$R^T = O\left(\sum_{a: \hat{p}_a^t \neq p_a} \frac{p_a - \hat{p}_a}{\text{KL}(p_a, \hat{p}_a)} \text{Int}\right)$$

$$q \leq \hat{p} + \sqrt{\frac{2}{n} \text{Int}}$$

this also holds

◻ $\text{KL}(p_a, p_\star) \geq 2(p_\star - p_a)^2$

Pinsker's Inequality

■ Thompson Sampling (1933), 2011-12

For $t = 0, 1, 2, \dots$

For $a \in A$

Let s_a^t denote the number of "successes". (1-reward)

Let f_a^t denote the number of "failures". (0-reward)

Draw $x_a^t \sim \text{Beta}(s_a^t + 1, f_a^t + 1)$

Pull argmax x_a^t

Mean $\frac{\alpha}{\alpha + \beta}$, Variance

$$\frac{\alpha \beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}$$

$$\Delta_a \stackrel{\text{def}}{=} P_{\star} - P_a$$

z_a^t be the event that arm a is pulled at time t .

Let z_a^t be a random variable that takes value 1 if arm a is pulled a time t and 0 otherwise.

$$E[z_a^t] = P\{z_a^t\}(1) + (1 - P\{z_a^t\})(0) = P\{z_a^t\}$$

$$u_a^t \stackrel{\text{def}}{=} \sum_{i=0}^{t-1} z_a^t$$

$$\bar{u}_a^t \stackrel{\text{def}}{=} \left[\frac{8}{(\Delta_a)^2} \ln(T) \right]$$

$$R^T = T P_{\star} - \sum_{t=0}^{T-1} E[r^t]$$

$$= T P_{\star} - \sum_{t=0}^{T-1} \sum_{a \in A} P(z_a^t) E[r^t | z_a^t]$$

$$= T P_{\star} - \sum_{t=0}^{T-1} \sum_{a \in A} \cancel{E[z_a^t]} P_a E[z_a^t] P_a$$

$$= \left(\sum_{a \in A} E[u_a^T] \right) P_{\star} - \sum_{a \in A} E[u_a^T] P_a$$

$$= \sum_{a \in A} E[u_a^T] (P_{\star} - P_a)$$

$$P_a \neq P_{\star}$$

$$= \sum_{a \in A} E[u_a^T] \Delta_a$$

$$P_a \neq P_{\star}$$

$$E[u_a^T] = \sum_{t=0}^{T-1} E[z_a^t]$$

$$= \sum_{t=0}^{T-1} P\{z_a^t \text{ and } u_a^t < \bar{u}_a^T\} + \sum_{t=0}^{T-1} P\{z_a^t \text{ and } (u_a^t > \bar{u}_a^T)\}$$

$$= A + B$$

$$A = \sum_{t=0}^{T-1} P\{z_a^t \text{ and } (u_a^t < \bar{u}_a^T)\}$$

$$= \sum_{t=0}^{T-1} \sum_{m=0}^{\bar{u}_a^T-1} P\{z_a^t \text{ and } (u_a^t = m)\}$$

$$= \sum_{m=0}^{\bar{u}_a^T-1} \sum_{t=0}^{T-1} P\{z_a^t \text{ and } (u_a^t = m)\}$$

$$= \sum_{m=0}^{\bar{u}_a^T-1} P\{(z_a^0 \text{ and } (u_a^0 = m)) \text{ or } (z_a^1 \text{ and } (u_a^1 = m)) \dots\}$$

$$< \sum_{m=0}^{\bar{u}_a^T-1} 1$$

$$< \bar{u}_a^T$$

$$B = \sum_{t=0}^{T-1} P\{z_a^t \text{ and } (u_a^t > \bar{u}_a^T)\}$$

$$\leq \sum_{t=0}^{T-1} P\left\{ (\hat{P}_a^t + \sqrt{\frac{2}{u_a^t} \ln(t)}) > \hat{P}_*^t + \sqrt{\frac{2}{u_*^t} \ln(t)} \right\}$$

$$\text{and } (u_a^t \geq \bar{u}_a^T)$$

$$\leq \sum_{t=0}^{T-1} \sum_{x=\bar{u}_a^T}^t \sum_{y=1}^t P\left\{ \hat{P}_a^t(x) + \sqrt{\frac{2}{x} \ln t} > \hat{P}_*^t(y) + \sqrt{\frac{2}{y} \ln t} \right\}$$

$$\hat{P}_a(x) + \sqrt{\frac{2}{x} \ln t} \geq \hat{P}_*(y) + \sqrt{\frac{2}{y} \ln t}$$

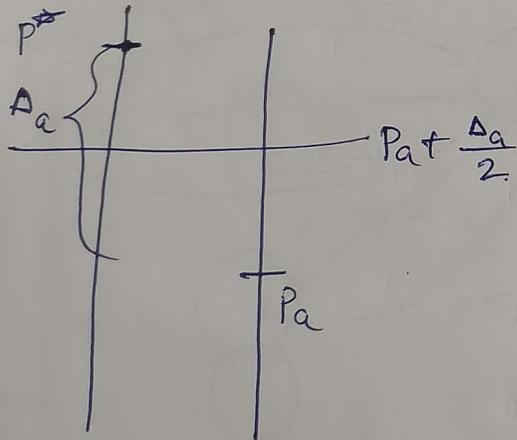
$$\Rightarrow (\hat{P}_a(x) + \sqrt{\frac{2}{x} \ln t}) \geq P_* \text{ or } (\hat{P}_*(y) + \sqrt{\frac{2}{y} \ln t}) < P_*$$

$x > \bar{u}_a T$, we have $\sqrt{\frac{2}{x} \ln t} \leq \sqrt{\frac{2}{\bar{u}_a T} \ln t} \leq \frac{\Delta_a}{2}$ and so,

$$\leq \sqrt{\frac{2 \Delta_a^2}{8 \zeta} \ln t} \ln t$$

$$\leq \frac{\Delta_a}{2}$$

$$\Rightarrow \hat{P}_a(x) + \sqrt{\frac{2}{x} \ln t} \geq P_* \Rightarrow \hat{P}_a(x) \geq P_a + \frac{\Delta_a}{2}$$



In summary:

$$\hat{P}_a(x) + \sqrt{\frac{2}{x} \ln t} \geq \hat{P}_*(y) + \sqrt{\frac{2}{y} \ln t} \geq \left(\hat{P}_a(x) \geq P_a + \frac{\Delta_a}{2} \right)$$

$$\text{or } \left(\hat{P}_*(y) < P_* - \sqrt{\frac{2}{y} \ln t} \right)$$

$$B \leq \sum_{t=0}^{T-1} \sum_{x=u_a T}^t \sum_{y=1}^t P\left\{\hat{P}_a(x) \geq p_a + \frac{\Delta a}{2}\right\} + P\left\{\hat{P}_a(y) \leq p_a - \sqrt{\frac{2}{\gamma} \ln t}\right\}$$

$$\leq \sum_{t=0}^{T-1} \sum_{x=u_a T}^t \sum_{y=1}^t \left(e^{-2x} \left(\frac{\Delta a}{2}\right)^2 + e^{-2y} \left(\sqrt{\frac{2}{\gamma} \ln t}\right)^2 \right)$$

$$\leq \sum_{t=0}^{T-1} \sum_{x=u_a T}^t \sum_{y=1}^t \left(e^{-4 \ln t} + e^{-4 \ln t} \right)$$

$$\leq \sum_{t=0}^{T-1} t^2 \left(\frac{1}{t^4}\right)$$

$$\leq \sum_{t=0}^{T-1} \frac{1}{t^2} = \frac{\pi^2}{3}$$

$$1 + \frac{1}{2^2} + \frac{1}{3^2} + \dots = \frac{\pi^2}{3}$$

Markov Decision Problems

S set of states

A set of Actions

$$A = \{\bar{a}_1, \bar{a}_2\}$$

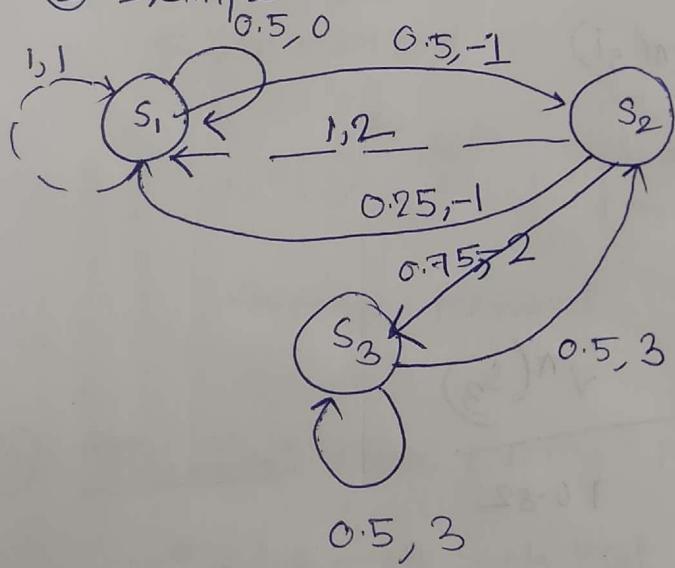
T transition function

T: $S \times A \rightarrow$ set of probability distribution over S.

R Reward function $R: S \times A \times S \rightarrow [-R_{\max}, R_{\max}]$

γ Discount factor $\gamma \in [0, 1]$

② Example



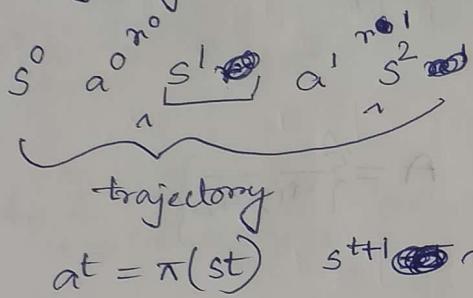
③ Policy $\pi: S \rightarrow A \rightarrow$ Not a part of the MDP.

Deterministic

Markovian - Action depends on only state

~~Deterministic~~ Stationary - Not changing with time

④ Trajectory over time



$$a_t = \pi(s_t) \quad s^{t+1} \sim T(s_t, a_t) \quad r^t = R(s_t, a_t, s^{t+1})$$

⑤ Value function of π

$$\begin{aligned} v^\pi: s &\rightarrow \mathbb{R} \\ v^\pi(s) &\stackrel{\text{def}}{=} \mathbb{E}[r^0 + \gamma r^1 + \gamma^2 r^2 + \dots] \end{aligned}$$

$$s^0 = s, \forall i = 0, 1, 2, \dots \quad : a_i = \pi(s_i)$$

⑥ Evaluation

| π | $v^\pi(s_1)$ | $v^\pi(s_2)$ | $v^\pi(s_3)$ |
|-------|--------------|--------------|--------------|
| 111 | 4.45 | 6.55 | 10.82 |
| 112 | -5.61 | -5.75 | -4.05 |
| 121 | 2.76 | 4.48 | 9.12 |
| 122 | 2.76 | 4.48 | 3.48 |
| 211 | 10 | 9.34 | 13.1 |
| 212 | 10 | 7.25 | 10 |
| 221 | 10 | 11 | 14.45 |
| 222 | 10 | 11 | 10 |

Continuing tasks and episodic tasks

i) $v^\pi(s) \stackrel{\text{def}}{=} E_\pi [r^0 + \gamma r^1 + \gamma^2 r^2 + \dots | s^0 = s]$
 Infinite discounted reward.

ii) $v^\pi(s) \stackrel{\text{def}}{=} E_\pi [r^0 + r^1 + \dots | s^0 = s]$
 Total reward.

Episodic tasks → Have a final state

iii) $v^\pi(s) \stackrel{\text{def}}{=} E_\pi [r^0 + r^1 + \dots | s^0 = s]$
 Finite horizon

Horizon: specified along with the MDP

iv) $v^\pi(s) \stackrel{\text{def}}{=} E \left[\lim_{m \rightarrow \infty} \frac{r^0 + r^1 + \dots + r^{m-1}}{m} \right]$

Average Reward

(SATR)

② Theorem: For every MDP, there exists a policy π^*

$\pi^*: s \rightarrow A$ such that $\forall s \in S, \forall \pi \in \Pi$.

$$v^{\pi^*}(s) > v^\pi(s)$$

π^* is an optimal policy. There can be many optimal policies, but they all have the same value function.

Finite horizon reward :-

$$\pi^*: S \times \{1, 2, \dots, \#\} \rightarrow A$$

Policy Evaluation

$$v^\pi(s) = E_\pi [r^0 + \gamma r^1 + \dots + r^n \mid s^0 = s]$$

$$= E_\pi [r^0 + s^0 = s] + \gamma E[r^1 + \gamma r^2 + \dots]$$

$$= \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) \cdot R(s, \pi(s), \bar{s})$$

$$+ \gamma \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) E_\pi [r^1 + \gamma r^2 + \gamma^2 r^3 + \dots]$$

$$= \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) \cdot R(s, \pi(s), \bar{s})$$

$$+ \gamma \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) E_\pi [r^1 + \gamma r^2 + \gamma^2 r^3 + \dots]$$

$$= \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) R(s, \pi(s), \bar{s}) + \gamma \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s})$$

$$E_\pi [r^0 + \gamma r^1 + \gamma^2 r^2 + \dots]$$

$\forall s \in S :$

$$v^\pi(s) = \sum_{\bar{s} \in S} T(s, \pi(s), \bar{s}) [R(s, \pi(s), \bar{s}) + \gamma v^\pi(\bar{s})]$$

$|S|$ equations in $|S|$ unknown

Bellman's Equations

$$Ax = B$$

■ Bellmann's equations are guaranteed to have solutions.

IV Action value function

$$q^\pi: S \times A \rightarrow \mathbb{R}$$

$$q^\pi(s, a) = \sum_{s' \in S} T(s, a, s') \left\{ R(s, a, s') + \gamma v^\pi(s') \right\}$$

① Bellman's Equations

$$\forall s \in S, \forall \pi \in \Pi$$

$$v^\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') \left\{ R(s, \pi(s), s') + \gamma v^\pi(s') \right\}$$

$$\forall s, \forall a \in A, \forall \pi \in \Pi$$

$$q^\pi(s, a) = \sum_{s' \in S} T(s, a, s') \left\{ R(s, a, s') + \gamma v^\pi(s') \right\}$$

Policy Evaluation

② Bellman's Optimality Equations

$$v^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s') \left\{ R(s, a, s') + \gamma v^*(s') \right\}$$

$v^*(s)$

$|S| = n$ equations in n unknowns

$$g^*(s, a)$$

$$v^* \rightarrow g^* \rightarrow \pi^*$$

$$\pi^* \rightarrow v^* \rightarrow g^*$$

◻ Bellman Operator

For $\pi \in \Pi$, $B^\pi: (S \rightarrow \mathbb{R}) \rightarrow (S \rightarrow \mathbb{R})$

is defined as follows. Let $x: S \rightarrow \mathbb{R}$

Then

$$(B^\pi(x))(s) \stackrel{\text{def}}{=} \sum_{s' \in S} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma x(s') \}$$

v^π is the fixed point of the Bellmann Operator.

◻ Bellman Optimality Operator

$B^*: (S \rightarrow \mathbb{R}) \rightarrow (S \rightarrow \mathbb{R})$ is defined as follows.

Let $x: S \rightarrow \mathbb{R}$

$$(B^*(x))(s) \stackrel{\text{def}}{=} \max_{a \in A} \sum_{s' \in S} T(s, a, s') \{ R(s, a, s') + \gamma x(s') \}$$

◻

Value Iteration

$$\lim_{t \rightarrow \infty} (B^*)^t(x) = v^*$$

Let $(V, \|\cdot\|)$ be a normed vector space. A mapping $T: V \rightarrow V$ is called a contraction mapping if there exists $L < 1$ such that $\forall u, v \in V$

$$\|Tu - Tv\| \leq L\|u - v\|.$$

Banach's Fixed Point Theorem:

Let $(V, \|\cdot\|)$ be a "complete" normed vector space and $T: V \rightarrow V$ a contraction mapping with constant L . Then T has a unique fixed point v . For $v_0 \in V$ and

$$v_{i+1} = T v_i, i = 0, 1, \dots$$

$$\|v_i - v\| \leq L^i \|v_0 - v\|$$

$$\text{Hence, } \lim_{i \rightarrow \infty} v_i = v$$

B^π and B^σ are contraction mappings in $\|\cdot\|_\infty$ with constant γ .

Take $X: S \rightarrow \mathbb{R}$ and $Y: S \rightarrow \mathbb{R}$

$$\|B^\pi(x) - B^\pi(y)\|_\infty \leq \gamma \|x - y\|_\infty$$

$$\|B^\pi(x) - B^\pi(y)\|_\infty = \max_{s \in S} |B^\pi(x)(s) - B^\pi(y)(s)|$$

$$= \max_{s \in S} \left| \sum_{s' \in S} T(s, \pi(s), s') (x(s') - y(s')) \right|$$

$$\leq \gamma \max_{s \in S} \sum_{s' \in S} T(s, \pi(s), s') |x(s') - y(s')|$$

$$\leq \gamma \max_{s \in S} \sum_{s' \in S} T(s, \pi(s), s') \|x - y\|_\infty$$

$$= \gamma \|x - y\|$$

■ $\|v_t - v_{t-1}\| \leq \epsilon.$

If v_t is close to v_{t-1} , is it close to v^* ?

Ans $\rightarrow \|v_t - v_{t-1}\| \leq \epsilon$

$$\|v_{t+1} - v_t\| \leq \gamma \epsilon$$

$$\|v_{t+2} - v_{t+1}\| \leq \gamma^2 \epsilon$$

$$\|v_T - v_{T-1}\| \leq \gamma^{T-t} \epsilon$$

$$\|v_T - v_{T-1}\| + \|v_{T-1} - v_{T-2}\| + \dots + \|v_{t+1} - v_t\|$$

$$\leq \epsilon (\gamma + \gamma^2 + \dots + \gamma^{T-t})$$

$$\Rightarrow \|v_T - v_t\| \leq \epsilon \gamma (1 + \gamma + \gamma^{T-t-1})$$

$$= \epsilon \gamma \cdot \frac{\gamma^{T-t-1}}{\gamma - 1}$$

$$\Rightarrow \lim_{t \rightarrow \infty} \|v_T - v_t\| \Rightarrow \|v^* - v_t\| = \frac{\epsilon \gamma}{\gamma - 1}$$

Linear Programming

Minimize $\sum_{s \in S} v(s)$

subject to $v(s) \geq \sum_{s' \in S} T(s, a, s') (R(s, a, s') + \gamma v'(s)) \quad \forall s, \forall a$

We can also minimize $5v_1 + 2v_2 + 3v_3$

$$5v(1) + 2v(2) + 3v(3)$$

~~Any~~ Any positive weight will do the work.

~~Answers~~

For $x : S \rightarrow \mathbb{R}$ and $y : S \rightarrow \mathbb{R}$, we define $x \geq y$ if $\forall s \in S : x(s) \geq y(s)$ and we define $x > y$ if $x \geq y$ and $\exists s \in S : x(s) > y(s)$.

If $x \geq y$, then $B^\pi(x) \geq B^\pi(y)$

$$B^\pi(x)(s) - B^\pi(y)(s)$$

$$= \sum_{s' \in S} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma x(s') \}$$

$$- \sum_{s' \in S} T(s, \pi(s), s') \{ R(s, \pi(s), s') + \gamma y(s') \}$$

$$= \gamma \sum_{s' \in S} T(s, \pi(s), s') \{ x(s') - y(s') \} \geq 0$$



Proof of Policy Improvement theorem:

Observe that for $\pi, \pi' \in \Pi$, $\forall s \in S$, $B^{\pi'}(v^{\pi})(s) = Q^{\pi}(s, \pi'(s))$



Howard's Policy Iteration

- Switch all with arbitrary action switch

Best time complexity $\rightarrow O\left(\frac{K^n}{n}\right)$

Other variants \rightarrow

Deterministic $O(K^{0.72n})$

Randomised $(O(\log K))^n$