

Another broad class of policies for MDPs are **index policies**

- Suppose state and action spaces decompose as $S = S_1 \times S_2 \times \dots \times S_k$ and $A = \{a_1, a_2, \dots, a_k\}$. Then an index policy comprises a set of fns $\phi_1, \phi_2, \dots, \phi_m$ with $\phi_i: S_i \rightarrow \mathbb{R}$ (indices) s.t. $A(s) = \underset{k \in [m]}{\text{argmax}} \{ \phi_i(s_i) \}$ for $s = (s_1, s_2, \dots, s_m)$

- In other words, for each 'part' of the state, we compute a fn, and then 'act on the part' with the highest value

Eg - (single machine scheduling with discounting)

- There are n jobs, with each job i having known processing time t_i and reward r_i upon completion. Want to schedule them on a single machine to maximize discounted sum of rewards (discount factor β)

- If $S_k \subseteq [n] \equiv$ set of remaining jobs, then
$$V(S_k) = \max_{i \in S_k} [r_i \beta^{t_i} + \beta^{t_i} V(S_k \setminus i)]$$

- If $n = 2$, $V(\{1, 2\}) = \max [r_1 \beta^{t_1} + r_2 \beta^{t_1+t_2}, r_2 \beta^{t_2} + r_1 \beta^{t_1+t_2}]$
 \Rightarrow we first serve $\underset{i \in \{1, 2\}}{\text{argmax}} \{ r_i \beta^{t_i} / (1 - \beta^{t_i}) \}$

- This problem has an easy soln via an **interchange argument**. Suppose order is $i_1 \rightarrow i_2 \rightarrow \dots \rightarrow j \rightarrow k \rightarrow \dots \rightarrow i_m$
 \Rightarrow Reward is of form $R_1 + \beta^{T+t_j} r_j + \beta^{T+t_j+t_k} r_k + R_2$
by previous argument, comparing $i_1 \rightarrow i_2 \rightarrow \dots \rightarrow j \rightarrow k \rightarrow \dots \rightarrow i_m$ and $i_1 \rightarrow i_2 \rightarrow \dots \rightarrow k \rightarrow j \rightarrow \dots \rightarrow i_m$ is same as comparing $r_j \beta^{t_j} / (1 - \beta^{t_j})$ and $r_k \beta^{t_k} / (1 - \beta^{t_k})$

\Rightarrow **OPT policy** \equiv serve jobs in decreasing order of $r_j \beta^{t_j} / (1 - \beta^{t_j})$

Bayesian Multi-Armed Bandits

- Now we consider a vast generalization of the above
 - There are m 'arms', where each arm i is a Markov chain $X_i[t]$ on state S_i .
 - In each round, we can play a single arm, i.e. $A = [m]$
 - If $A[t] = i$, then $R(t) = R_i(X_i[t])$ and $X_i[t] \rightarrow X_i[t+1]$
 - (**Non-restlessness**) The state of arm i changes only when it is 'played' (i.e. $X_i[t+1] = X_i[t] \forall i$ s.t. $A[t] \neq i$)
 - (**Discounted infinite-horizon objective**) $\max \sum_{t=0}^{\infty} \beta^t (\sum_{i=1}^m \mathbb{1}_{\{A[t]=i\}} R_i(X_i[t]))$
- Eg (MAB with Beta-Bernoulli priors)
 - m actions, where action i gives reward $R_i(\theta_i)$
 - θ_i unknown, but assume $\theta_i \sim \text{Beta}(N_i, S_i)$ (prior)
 - If we play action k , then post-action $\theta_i = \begin{cases} \text{Beta}(N_i+1, S_i+1) & \text{if } R_i=1 \\ \text{Beta}(N_i, S_i) & \text{if } R_i=0 \end{cases}$
 - Aim: $\max \sum_{t=0}^{\infty} \beta^t (\sum_{i=1}^m \mathbb{1}_{\{A[t]=i\}} R_i(t))$

The 1.5 arm problem - Suppose we just have 2 arms

- 1) MC $X[t]$ with reward $R_i(X_i[t])$
- 2) Constant reward arm with reward γ

Now opt policy is a stopping problem: play arm 1 till some time T (stopping time), then play 2 forever (since no new info)

$$\Rightarrow R = \sup_{T \geq 0} \mathbb{E} \left[\sum_{t=0}^{T-1} R_i(X_i[t]) \beta^t + \beta^T \gamma / (1-\beta) \right]$$

The **Gittins Index** of arm 1 in state $X_i[0]=2$ is the smallest constant reward γ s.t. you are indifferent between playing arm 1 in state $X_i[0]$ and arm 2

Formally for any $x \in S$,

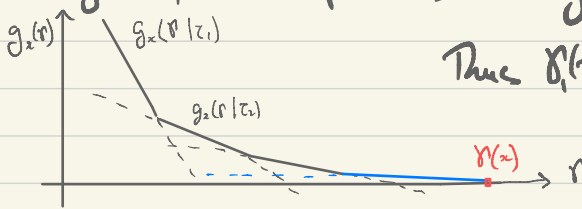
$$V_1(x) = \sup \left\{ v : \frac{v}{1-\beta} \leq \sup_{\tau > 0} \mathbb{E} \left[\sum_{t=0}^{\tau-1} \beta^t R_1(x_t) + \frac{\beta^\tau v}{1-\beta} \mid x_0 = x \right] \right\}$$

This can also be viewed as the **minimum per-bound charge** for pulling arm 1 s.t. you are indifferent between pulling once or not when $x_0 = x$

$$V_1(x) = \sup \left\{ v : 0 \leq \sup_{\tau > 0} \mathbb{E} \left[\sum_{t=0}^{\tau-1} \beta^t [R_1(x_t) - v] \mid x_0 = x \right] \right\} \quad g_x(v)$$

Note that $g_x(v)$ is decreasing and convex in v

- To see this, note that for a fixed sample path x_1, x_2, \dots and fixed τ , $\sum_{t=0}^{\tau-1} \beta^t [R_1(x_t) - v]$ is linear decreasing
- Taking expectation preserves linearity, and sup over τ makes it convex



Thus $V_1(x)$ has a unique soln

Also for the optimal τ , we have $\mathbb{E} \left[\sum_{t=0}^{\tau-1} \beta^t R_1(x_t) - v(x) \sum_{t=0}^{\tau-1} \beta^t \mid x_0 = x \right] = 0$

Thus

$$V_1(x) = \sup_{\tau > 0, \text{ stopping time}} \frac{\mathbb{E} \left[\sum_{t=0}^{\tau-1} \beta^t R_1(x_t) \mid x_0 = x \right]}{\mathbb{E} \left[\sum_{t=0}^{\tau-1} \beta^t \mid x_0 = x \right]}$$

Gittins Index for arm 1
in state x

Eg - Suppose $X_i[0] = X_i[1] = X_i[2] = \dots = \begin{cases} M & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases}$ 'Collapsing' arm
 (play once to learn)
 $R_i(x) = x$

Then $V('unknown') = \sup[V] \quad V/1-\beta \leq pM/1-\beta + (1-p)V/\beta$
 $= \frac{pM}{1-\beta(1-p)}$

Eg - Single job with reward r_i , processing time t_i
 (amount of job processed) $R_i(0) = \sup_{z > 0} \frac{r_i \beta^{t_i} \mathbb{1}\{\tau > t_i\}}{\sum_{t=1}^{\infty} \beta^t} = \frac{r_i \beta^{t_i}}{1-\beta^{t_i}} (1-\beta)$
 the index we get via interchange

Thm (Gittins '79) - For finite arms $[m]$, and bounded rewards $R_i(x) \in [-C, C] \forall i \in [m], x \in S_i$:
 A policy is optimal if and only if it always selects arm i at time t with highest Gittins index $V_i(x_i[t])$.

There are actually more general conditions for when an index policy is optimal. When is it not, though?
 - Independence of irrelevant alternatives (IIA): A policy π satisfies IIA if for any set of arms $[m]$ and $i \in [m]$, if $\pi([m]) = i$, then $\pi([n]) = i$ for any $[n] \supseteq [m]$.

An index policy is optimal if and only if opt policy is IIA.