

Notes for 2022 exam

abandon

Q1)

states

C0

C1 - not good

C1 - suitable

C1 - perfect

C2 - not good

C2 - suitable

C2 - perfect

- { C0 is starting state
- abandon and recruit are terminal states

recruit

actions

- A (abandon): available in all states but the terminal ones, lead to abandon state with prob. 100% and reward 0
- R (recruit): available in C1, C2 - suitable and -perfect states, lead to recruit state with prob. 100% and reward 200 for C1, C2 - perfect and 50 for C1, C2 - suitable
- I (interview): available in C0 and C1 states, leads to C1 - not good wp. 25%, C1 - suitable wp. 50% or C1 - perfect wp. 25% with reward -30 from C0 and to C2 - not good wp. 25%, C2 - suitable wp. 50% or C2 - perfect wp. 25% with reward -30 from any of the C1 states.

Q2) Bellman equations for v_{π}

$$\begin{cases} v_{\pi}(s_0) = (-1 \times 0.3) + 0.9 [0.3 v_{\pi}(s_0) + 0.7 v_{\pi}(s_1)] \\ v_{\pi}(s_1) = (-1 \times 0.85) + 0.9 [0.85 v_{\pi}(s_1) + 0.15 v_{\pi}(s_2)] \\ v_{\pi}(s_2) = 10 \end{cases}$$

Q3) Policy improvement: in s_1 , the greedy action w.r.t. v_{π} is

$$\underset{a}{\operatorname{argmax}} \left[r(s_1, a) + \gamma \sum_{s'} p(s'|s_1, a) v_{\pi}(s') \right]$$

is obviously achieved when $a = \text{east}$ for which

$$\begin{aligned} v_{\tilde{\pi}}(s_1) &= (-1 \times 0.3) + 0.9 [0.3 v_{\tilde{\pi}}(s_1) + 0.7 \times 10] \\ &= 6.575 > 2.128 \end{aligned}$$

Hence $\tilde{\pi}$ which differs from π by choosing action east in s_1 has higher value function

Q4) This time it holds both in s_0 and s_1 that

$$v_{\tilde{\pi}}(s) = \underset{a}{\operatorname{argmax}} \left[r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\tilde{\pi}}(s') \right]$$

showing that $v_{\tilde{\pi}} = v^*$ and hence that $\tilde{\pi}$ is optimal

Q5) Usually this would be shown by backward induction (Bellman equations for optimal policy in finite-horizon case) but here there are only 8 possible sequence of actions (see Q7 below) and it is easily checked that the best ones are (2, 1, *) both yielding a reward of 2.

Q 6) $\begin{cases} \pi(1|s) = 1 & \text{stays in state 1} \\ \pi(2|s) = 1 & \text{stays in states } \{1, 2\} \end{cases}$
 and neither of them reaches the third state which provides rewards

Q 7) See part of course on REINFORCE algorithm

Q 8)	Sequence	$\sum_{i=0}^{\infty} R_{i+1}$	$\sum_{i=0}^{\infty} \frac{\log \pi_\theta(A_i s_i)}{\theta}$	prob. sequence
	(2, 1, 1)	2	$1/\theta - 2/(1-\theta)$	$\theta(1-\theta)^2$
	(2, 1, 2)	2	$2/\theta - 1/(1-\theta)$	$\theta^2(1-\theta)$
	(1, 2, 1)	1	$1/\theta - 2/(1-\theta)$	$\theta(1-\theta)^2$

Q 9) The REINFORCE formula yields

$$\begin{aligned}
 & 2 \left[(1-\theta)^2 - 2\theta(1-\theta) \right] \\
 & + 2 \left[2\theta(1-\theta) - \theta^2 \right] \\
 & + 1 \left[(1-\theta)^2 - 2\theta(1-\theta) \right] \\
 & = 3\theta^2 - 8\theta + 3 \quad \text{for the gradient of} \\
 & \quad \text{the policy}
 \end{aligned}$$

The optimal policy corresponds to the root that is in $(0, 1)$:

$$\theta^* = \frac{4 - \sqrt{7}}{6} \approx 0,451$$

Q 10) There are $k(k-1)$ actions consisting of the (ordered) pairs (j, k) where $j \in \{1, \dots, k\}$, $k \in \{1, \dots, k\}$ with $j \neq k$

The best action is (j^*, k^*) where

$$\begin{cases} \theta_{j^*} > \theta_j \text{ for } j \neq j^* \text{ (highest } \theta_i) \\ \theta_{k^*} > \theta_k \text{ for } k \neq k^*, j^* \text{ (second highest } \theta_i) \end{cases}$$

Q11) See course for the definition of regret and its fundamental rewriting. Here

$$\mathbb{E}[R_T] = \sum_{(j,k) \neq (j^*, k^*)} [(\theta_{j^*} - \theta_j) + \alpha (\theta_{k^*} - \theta_k)] \mathbb{E}[N_{j,k}(T)]$$

$$\text{where } N_{j,k}(T) = \sum_{t=1}^T \mathbf{1}_{\{A_t = (j, k)\}}$$

Q12) $\bar{X}_{j,k}(t) = \frac{1}{N_{j,k}(t)} \sum_{s=1}^t X_s \mathbf{1}_{\{A_s = (j, k)\}}$

UCB plays

$$\arg \max_{(j,k)} \bar{X}_{j,k}(t) + \sqrt{\frac{\gamma \log t}{2 N_{j,k}(t)}}$$

Q13) $\Delta_{j,k} = (\theta_{j^*} - \theta_j) + \alpha (\theta_{k^*} - \theta_k)$ and the result from the course is that for UCB

$$\mathbb{E}[R_T] \leq \sum_{(j,k) \neq (j^*, k^*)} C \frac{\log T}{\Delta_{j,k}} + O(1)$$

which grows as $\log T$, as expected, but

also as K^2 (depending on the gaps) which is surprising (and in fact not optimal) as there are only K unknown values of θ_i (and not really $K(K-1)$ unrelated arms).

Q14) Action (f, b) is associated with vector

$$A_b = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \cdot \begin{array}{l} f^{\text{th}} \text{ position} \\ b^{\text{th}} \text{ position} \end{array}$$

such that

$$\mathbb{E}[X_t | A_t] = A_b^T \theta = \theta_f + \alpha \theta_b$$

where

$$\theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_K \end{pmatrix}$$

Q15) $\bar{X}_1(\tau)$ and $\bar{X}_2(\tau)$ are independent and follow, respectively, $N(\mu_1, \frac{\sigma^2}{\tau/2})$ and $N(\mu_2, \frac{\sigma^2}{\tau/2})$ distributions, hence

$$\bar{X}_1(\tau) - \bar{X}_2 \sim N(\Delta, \frac{4\sigma^2}{\tau})$$

Q16) The error (ie. selecting arm 2) occurs when $\bar{X}_2(\tau) > \bar{X}_1(\tau)$, that is,

when $\bar{X}_1(\tau) - \bar{X}_2(\tau) < 0$

Hence

$$P(\text{error}) = P\left(\underbrace{\bar{X}_1(\tau) - \bar{X}_2(\tau) - \Delta}_{\sim N(0, \frac{46^2}{\tau})} < -\Delta\right) \leq e^{-\frac{\Delta^2 \tau}{86^2}}$$

which is less than S when

$$\frac{\Delta^2 \tau}{86^2} \geq \log \frac{1}{S}$$

that is

$$\tau \geq \frac{86^2}{\Delta^2} \log \frac{1}{S}$$

(Q17) The error occurs when

$$\bar{X}_1(\tau) - \bar{X}_2(\tau) < \frac{46^2 \log(1/S)}{\Delta \tau}$$

Hence

$$P(\text{error}) \leq e^{-\underbrace{\left(\Delta + \frac{46^2 \log(1/S)}{\Delta \tau}\right)^2 \frac{\tau}{86^2}}_{\geq \frac{86^2 \log(1/S)}{\tau}}}$$

thus

$$P(\text{error}) \leq e^{-\log(\frac{1}{S})} = S$$

(Q18) The exact probability of correct (i.e. choose arm 1) is given by

$$P(\bar{X}_1(\tau) - \bar{X}_2(\tau) > \frac{46^2 \log(1/S)}{\Delta \tau})$$

where $\bar{X}_1(T) - \bar{X}_2(T) \sim N(\Delta, \frac{\delta^2}{T})$

By symmetry of the Gaussian distribution this is larger than $1/2$ iff.

$$\frac{462 \log(1/\delta)}{\Delta T} \leq \Delta$$

that is

$$T \geq \frac{462}{\Delta^2} \log\left(\frac{1}{\delta}\right)$$

Q19)

$$KL(p||q) = \mathbb{E}_{x \sim p} \left[\log \frac{p(x)}{q(x)} \right], \text{ here}$$

$$\begin{aligned} &= \mathbb{E}_{x \sim N(\mu, \sigma^2)} \left[\underbrace{\frac{1}{2\sigma^2} \{(x-\mu')^2 - (x-\mu)^2\}}_{= 2(x - \frac{\mu' + \mu}{2})(\mu - \mu')} \right] \\ &= \frac{(\mu - \mu')^2}{2\sigma^2} \end{aligned}$$

Q20) (Hard!) Using the general inequality with

$$\underline{\text{model } V}: \quad \mu_1, \mu_2 < \mu_1$$

$$\underline{\text{model } V'}: \quad \mu'_1 = \mu_1, \mu'_2 > \mu_1$$

yields

$$KL(\mu_1, \mu'_2) \mathbb{E}_V [N_2(T)] \geq d(\delta, 1-\delta)$$

by selecting E as the event that arm 2 get selected, which is an error under model V which occurs w.p. at most δ whereas it is the correct outcome under model V' (as $\mu'_2 > \mu'_1$) and hence occurs w.p. at least $1-\delta$

(Notice that when $\delta < 1/2$ $d(p, q) > d(\delta, 1-\delta)$ when $p < \delta$ and $q > 1-\delta$)

Hence

$$\frac{(\mu_2 - \mu'_2)^2}{\Delta^2} \mathbb{E}_r[N_2(T)] \geq d(\delta, 1-\delta)$$

for all $\mu'_2 > \mu_1$ so that it also holds when $\mu'_2 = \mu_1$ by continuity and thus

$$\mathbb{E}_r[N_2(T)] \geq \frac{\varepsilon \delta^2}{\Delta^2} d(\delta, 1-\delta)$$

Proceed similarly with the second change of distribution to obtain the inequality for $\mathbb{E}_r[N_1(T)]$ (when $\mu_1 < \mu_2$).