

Naser Kaze mi 99/02/05

(Q1)

A) The purpose of the first expression  $E_\mu[\alpha] \cdot E_{\hat{\pi}_\mu}[\alpha]$  is to find the  $\mu$  such that maximizes the  $E_\mu[\alpha]$ , and the minimizing it, preventing having  $(s, a)$  pairs which cause a pick in  $\alpha$ -value. i.e pushing down the  $\alpha$ -values. The  $-E_{\hat{\pi}_\mu}[\alpha]$  is for the opposite, maximizing the  $\alpha$ -value for observed data, i.e pushing up on  $(s, a)$  samples. so we just minimize the  $\alpha$ -value for unobserved data.

B) We have

$$C(\alpha) = \min_{\alpha} \max_{\mu} \alpha(E_{s \sim D}, a \sim \mu(a|s)) [\alpha(s, a)] - E_{s \sim D, a \sim \hat{\pi}_\mu(a|s)} [\alpha(s, a)] + \frac{1}{2} E_{s, a, s' \sim D} [(Q(s, a) - \hat{Q}^{\pi_\mu}(s, a))^2] + H(\mu)$$

so at first we have to find  $\mu$  such that

$$\max_{\mu} E_{s \sim D, a \sim \mu(a|s)} [\alpha(s, a) - \log \mu(a|s)]$$

$$\text{s.t. } \sum_a \mu(a|s) = 1 \quad \text{for each } s.$$

We can solve the above optimization problem for each  $s \in D$  separately, since each  $\mu(\cdot|s)$  is only conditioned on  $s$ . Then for a fixed  $s$  we have the dual problem is:

$$E_{a \sim \mu(a|s)} [Q(s,a) - \log \mu(a|s)] - \lambda \sum_a \mu(a|s) - \lambda$$

so by the KKT conditions we have:

$$\nabla_\lambda \sum_a \mu^*(a|s) (Q(s,a) - \log \mu^*(a|s)) - \nabla_\mu \lambda \sum_a \mu(a|s) = 0$$

$$\Rightarrow Q(s,a) - \log \mu^*(a|s) - 1 - \lambda = 0$$

$$\Rightarrow \log \mu^*(a|s) = Q(s,a) - 1 - \lambda$$

$$\Rightarrow \mu(a|s)^* = \frac{\exp(Q,s,a)}{e^{1+\lambda}}$$

$$\text{and } \sum_a \mu^*(a|s) = 1 \Rightarrow e^{1+\lambda} = \sum_a \exp(Q,s,a)$$

$$\Rightarrow \mu^*(a|s) = \frac{\exp(Q(s,a))}{\sum_{a'} \exp(Q(s,a'))}$$

then we will have:

$$CQL = \min_Q E_{s \sim D} \left[ E_{a \sim \mu^*} \left[ Q(s,a) + \log \sum_{a'} \exp(Q(s,a')) - Q(s,a) \right] - E_{a \sim \hat{\pi}_\beta} [Q(s,a)] \right] + \dots$$

$$= \min_Q E_{s \sim D} \left[ \log \sum_a \exp(Q(s,a)) - E_{a \sim \hat{\pi}_\beta} [Q(s,a)] \right]$$

$$+ \frac{1}{2} E_{s,a,s' \sim D} [(Q - \hat{\beta}^{\pi_k} \hat{Q}^k)^2]$$

Adding this regularizing term, other than simplifying the maximization part, introduces the exploration to our learned policy, avoiding the concentration on large  $Q$ -values, leading to a conservative policy.

c) We have

$$CQL = \min_Q \alpha E_{s \sim D, a \sim \mu^*} [Q(s, a)] + \frac{1}{2} E_{s, a, s' \sim D} [(Q - \hat{B}^\pi \hat{Q}^k)^2] + R(\mu^*)$$

for a fixed  $s, a$  we have

$$\alpha \mu^*(a|s) Q(s, a) + \frac{1}{2} \hat{\sigma}_\beta(a|s) (Q(s, a) - \hat{B}^\pi \hat{Q}^k(s, a))^2$$

$$\nabla_Q \rightarrow \alpha \mu^*(a|s) + \hat{\sigma}_\beta(a|s) (Q(s, a) - \hat{B}^\pi \hat{Q}^k(s, a)) = 0$$

$$\Rightarrow \hat{Q}^{k+1}(s, a) \sim Q(s, a) = \hat{B}^\pi \hat{Q}^k(s, a) - \alpha \frac{\mu^*(a|s)}{\hat{\sigma}_\beta(a|s)}$$

As we see, in each iteration of the Bellman operator, we have

$\hat{Q}^{k+1}(s, a) \leq \hat{B}^\pi \hat{Q}^k(s, a)$ , and we really are underestimating the next  $Q$ -value

D)  $\forall a, s, a \in D$ ,  $|\hat{B}^\pi Q(s, a) - B^\pi Q(s, a)| \leq C_\delta(s, a)$ . By the above update rule,

$$(B^\pi Q = r + \gamma P^\pi Q)$$

$$\hat{Q}^{k+1}(s, a) \leq B^\pi \hat{Q}^k(s, a) - \alpha \frac{\mu^*(a|s)}{\hat{\sigma}_\beta(a|s)} + C_\delta(s, a)$$

$$\Rightarrow \hat{Q}^{k+1} \leq (I - \gamma P^\pi)^{-1} \left[ r - \alpha \frac{\mu^*}{\hat{\sigma}_\beta} + C_\delta \right]$$

Reward,  $Q$ -value function

$$\Rightarrow \hat{Q}^{k+1}(s, a) \leq \hat{Q}^k(s, a) - \alpha \left[ (I - \gamma P^\pi)^{-1} \left[ \frac{\mu^*}{\hat{\sigma}_\beta} \right] \right](s, a) + \left[ (I - \gamma P^\pi)^{-1} C_\delta \right](s, a)$$

Q2

A) BC only tries to mimic the experts behavior. But GAIL learns the reward function and a policy at least as good as the expert in a robust manner. GAIL also mitigates the compounding error problem. GAIL has a better generalization to unseen states, since it focuses on matching the overall distribution of trajectories. GAIL has a better exploration ability due to the entropy regularization term.

B) In IRL, multiple reward functions can describe the same expert policy. This is called the inherent ambiguity in IRL. The minimum entropy principle, states the among the reward functions describing the same optimal policy, the one with the highest trajectory entropy.

The method:

The trajectory distribution is considered to be of the form:

$$P(T|R) \propto \exp\left(\sum_{t=0}^T R(s_t, a_t)\right)$$

and the objective is to maximize:

$$\max_R \sum_{T \in D} \log P(T|R) - \lambda \sum_p p(s, a) R(s, a)$$

This principle ensures that the policy is as random as possible, while remaining consistent with the observed behaviour, avoiding unnecessary assumptions about the expert.

C) We have:  $J(\pi_E, f) \equiv E_{\pi_E} [\log(1 - D(s, a))]$  where the  $D$ , is the discriminator function and plays the role of "f" in  $J$ .

as for the  $J(\pi, f) = \lambda \mathcal{H}(\pi) - E_{\pi} [\log(D(s, a))]$ .

From another point of view we have :

$D$ : as the discriminator, tries to minimize the discrimination for the policy with respect the expert, while maximizing it for the learned policy.

This is adversarial part. On the other hand, the  $\pi$  tries to minimize this discrimination, while maximizing  $\mathcal{H}(\pi)$ . Here,  $\mathcal{H}$  addresses the ambiguity problem, which by ManEnt, we prefer the policy with max  $\mathcal{H}(\pi)$ .

D) These methods, suffer from the limit in choice of features, and the complexity of solving a more complicated optimization problem.

In contrast, GAIL is more flexible in terms of features. e.g the discriminator  $D$  can be a neural network. Also it deals with the discrimination directly.

Adversarial training also improves both policy and discriminator.

Q3

A) Pure model methods suffer from model inaccuracies, with generalization errors for unseen data. This leads to suboptimal policies. Pure environmental methods require lots of computations with extensive interaction with the environment. It can be time-costly and sample inefficient.

Hybrid methods, leverage the strength in both. e.g using the model to generate and simulate environment interaction, increasing the sample efficiency and reducing computation costs, while holding to the ground truth environment, mitigating the model inaccuracies over long-term.

B) MBPO integrates model simulation rollouts with real-world interactions.

It selects random points on the real world trajectories and starts to rollout using the model, from the selected points, and generate new trajectories.

MBPO usually uses rollout with 1-3 steps. Longer rollouts lead to accumulated errors by the model, causing unreal hallucination. But shorter rollouts may result in data with less diversity.

C) Since MBPO relies on real-time interactions with the environment and the model simulation, it needs to operate online.

D) MOrEL measures the model uncertainty by computing the variance in prediction using an ensemble of learned models. When a model is uncertain about a state-action pair, treats it with caution. As mentioned, a method for measuring the uncertainty, is to use the variance of an ensemble of models on the state-action pairs predictions.

E) COMBO leverages the Model-based Dyna-style algorithm to generate data, and adds it to the total data buffer, which is a mixture of real-world data and the model rollouts. Then CQL is applied to penalize the overestimation of Q-values by enforcing a conservative objective. This is for ensuring the robustness of the model when dealing with the data generated by the Dyna component.

F)

Conservative methods:

- pros: Robustness to model inaccuracies. Reduced overestimation. Avoiding risky actions.
- cons: computational complexities. Slow convergence. Hyperparameter sensitivity.

Uncertainty based methods:

- pros: Exploration encouragement. Handling uncertainty.
- cons: Overexploration. Complex implementation.

Conservative methods tend to be more exploitative, and robust to inaccuracy, but sensitive to hyperparameter. Uncertainty methods on the other hand are explorative, and less hyperparameter sensitive.

Q4

A) The Hoeffding's inequality states that given a sequence of random variables  $X_1, \dots, X_n$ , with  $S_n = \sum_{i=1}^n X_i$ , we have

$$P(|S_n - E[S_n]| \geq t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$$

where  $a_i \leq X_i \leq b_i$ .

If  $X_i$ 's are i.i.d., we can rewrite the inequality as

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{2n\varepsilon^2}{(b-a)^2}\right)$$

In our case, since rewards are in  $[-1, 1]$ , thus  $(b-a)^2 = 4$ , and we can have

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq \varepsilon\right) \leq 2 \exp\left(-\frac{n\varepsilon^2}{2}\right)$$

Letting  $\varepsilon = \sqrt{\frac{2 \log(1/\delta)}{n}}$ , we obtain:

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - E[X]\right| \geq \sqrt{\frac{2 \log(1/\delta)}{n}}\right) \leq 2 \exp\left(-\log\frac{1}{\delta}\right) = 2\delta$$

So

$$P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \leq \sqrt{\frac{2 \log(1/\delta)}{n}}\right) \geq 1 - 2\delta$$

here  $\frac{1}{n} \sum_{i=1}^n X_i \equiv \hat{\mu}_{i(t-1)}$ ,  $n \equiv T_i(t-1)$  and  $\mu$  is the true mean reward for action  $i$ . Hence:

$$P(|\hat{\mu}_i(t-1) - \mu_i| \leq \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}}) \geq 1 - 2\delta$$

So w.h.p  $(1-2\delta)$ ,  $\hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}}$  is an upper bound for  $\mu_i$ .

B) Let arm  $i$  have highest reward. Let's define event  $G_i$  as

$$G_i = \left\{ \mu_i \leq \min_{1 \leq t \leq n} UCB_i(t, \delta) \right\} \cap \left\{ \hat{\mu}_{i,u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} < \mu_i \right\}$$

For some constant  $u_i$ .

$G_i$  is the event of not underestimating  $\mu_i$ , while the UCB for the  $i$ th arm is below the mean reward of the optimal arm  $\mu_1$ .

Now we can write  $E[T_i(n)]$  as:

$$E[T_i(n)] = E[1(G_i)T_i(n)] + E[1(G_i^c)T_i(n)]$$

Let's first compute an upper bound for  $E[1(G_i)T_i(n)]$ . We first show that when  $G_i$  holds, then  $T_i(n) \leq u_i$ . Let  $G_i$  hold. Suppose  $T_i(n) > u_i$ .

This means arm  $i$  has been played more than  $u_i$  times. So there is some round

$t \in \{1, \dots, n\}$ , which  $T_i(t-1) = u_i$  and  $A_t = i$ . Therefore

$$\begin{aligned} UCB_i(t-1, \delta) &= \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}} \\ &= \hat{\mu}_{i,u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \\ &< u_i \end{aligned}$$

$$< UCB_1(t-1, \delta)$$

So  $A_t = \arg\max_i UCB_j(t-1, s) \neq i$ , which is a contradiction. On the other hand, we have:

$$G_j^c = \underbrace{\{u_i > \min_{1 \leq t \leq n} UCB_j(t, s)\}}_A \cup \underbrace{\{\hat{\mu}_i u_i + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \geq \mu_i\}}_B$$

For A we have:  $(T_j(t) < t)$

$$\begin{aligned} A &= \{u_i > \min_{1 \leq t \leq n} UCB_j(t, s)\} \subset \{\mu_i > \min_{1 \leq t \leq n} \hat{\mu}_i + \sqrt{\frac{2 \log(1/\delta)}{t}}\} \\ &= \bigcup_{1 \leq t \leq n} \{\mu_i > \hat{\mu}_i + \sqrt{\frac{2 \log(1/\delta)}{t}}\} \end{aligned}$$

By the union bound we have

$$\begin{aligned} P(\mu_i > \min_t UCB_j(t, t)) &\leq P\left(\bigcup_{1 \leq t \leq n} \{\mu_i > \hat{\mu}_i + \sqrt{\frac{2 \log(1/\delta)}{t}}\}\right) \\ &\leq \sum_t P(\{\mu_i > \hat{\mu}_i + \sqrt{\frac{2 \log(1/\delta)}{t}}\}) \end{aligned}$$

Lemma 1: Let  $X \in [-2, 2]$ , then we have:

$$P(X \geq \varepsilon) = P(\exp(tX) \geq \exp(t\varepsilon))$$

$$\leq E[\exp(tX)] \exp(-t\varepsilon) \quad (\text{Markov's inequality})$$

$$\leq E[\exp(2t - t\varepsilon)] \quad \text{Letting } t = \frac{\varepsilon^2}{2(\varepsilon - 2)}$$

$$= \exp(-\varepsilon^2/2)$$

Thus

$$P(\{\mu_i > \hat{\mu}_i + \sqrt{\frac{2 \log(1/\delta)}{t}}\}) \leq \exp\left(-\frac{\sqrt{2 \log(1/\delta)}}{2}\right)^2 = \delta$$

And

$$P(\mu_i \geq \min_t VCB_i(t, t)) \leq n\delta$$

As for  $B$ , let  $u_i$  be chosen large enough so that:

$$\Delta_i - \sqrt{\frac{2 \log(1/\delta)}{u_i}} \geq c\Delta_i \quad c \in (0, 1) \quad (1)$$

then since  $\hat{\mu}_i = \mu_i + \Delta_i$ , then again using lemma 1,

$$\begin{aligned} P(\hat{\mu}_{i|u_i} + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \geq \mu_i) &= P(\hat{\mu}_{i|u_i} - \mu_i + \sqrt{\frac{2 \log(1/\delta)}{u_i}} \geq u_i - \mu_i) \\ &= P(\hat{\mu}_{i|u_i} - \mu_i \geq \Delta_i - \sqrt{\frac{2 \log(1/\delta)}{u_i}}) \\ &\leq P(\hat{\mu}_{i|u_i} - \mu_i \geq c\Delta_i) \\ &\leq \exp\left(-\frac{u_i c^2 \Delta_i^2}{2}\right) \quad (\text{by just setting } t = \frac{u_i \varepsilon^2}{2(\varepsilon-2)} \text{ in lemma 1}) \end{aligned}$$

So overally

$$P(G_i) \leq n\delta + \exp\left(-\frac{u_i c^2 \Delta_i^2}{2}\right)$$

And

$$E[T_i(n)] \leq u_i + n(n\delta + \exp\left(-\frac{u_i c^2 \Delta_i^2}{2}\right))$$

choosing the smallest  $u_i$  so that (1) holds, we get  $u_i = \lceil \frac{2 \log(\frac{1}{\delta})}{(1-c^2)\Delta_i^2} \rceil$

and  $\delta = \frac{1}{n^2}$ , we obtain

$$E[T_i(n)] \leq u_i + 1 + n c^{-2} = \lceil \frac{2 \log(\frac{1}{\delta})}{(1-c^2)\Delta_i^2} \rceil + 1 + n c^{-2}$$

and setting  $c = \frac{1}{n^2}$ , leads to

$$\lceil \frac{4 \log n}{n^2 \Delta_i^2} \rceil \leq \frac{16 \log n}{\Delta_i^2}$$

and

$$\begin{aligned} n e^{-\frac{u_i c^2 \Delta_i^2}{2}} &\leq e^{\log n - \frac{u_i c^2 \Delta_i^2}{2}} \\ &\leq n^{1 - \frac{c^2 \Delta_i^2 + 4 \log n}{2 \log n (1-c^2) \Delta_i^2}} \\ &= n^{1 - \frac{2c^2}{(1-c^2)}} \end{aligned}$$

by just setting  $c \geq 2: \leq 1$

so

$$E[T_i(n)] \leq 3 + \frac{16 \log n}{\Delta_i^2}$$

(c)

Lemma 2: for any policy  $\pi$ , countable action set  $A$  and horizon  $n$ ,  
the regret  $R_n$  for this policy can be decomposed as

$$R_n = \sum_{a \in A} \Delta_a E[T_a(n)]$$

Proof: For any fixed  $t$  we have  $\sum_{a \in A} \mathbb{1}\{A_t = a\} = 1$ . Therefore

$$S_n = \sum_t X_t = \sum_t \sum_a X_t \mathbb{1}\{A_t = a\} \text{. Thus}$$

$$h_n = n \mu^* - E[S_n] = \sum_{t=1}^n \sum_{a \in A} E[(\mu^* - X_t) \mathbb{1}\{A_t = a\}]$$

$$E[(\mu^* - X_t) \mathbb{1}\{A_t = a\}] = E[E[(\mu^* - X_t) \mathbb{1}\{A_t = a\} | A_t]]$$

$$= E[\mathbb{1}\{A_t = a\}]E[(\mu^* - X_t) | A_t]$$

$$= E[\mathbb{1}\{A_t = a\}](\mu^* - \mu_{A_t=a})$$

$$= E[\mathbb{1}\{A_t = a\}]\Delta_a$$

Therefore

$$h_n = \sum_{a \in A} \sum_{t=1}^n E[\mathbb{1}\{A_t = a\}]\Delta_a = \sum_{a \in A} E[T_a(n)]\Delta_a$$

So by Lemma 2 we have  $h_n = \sum_{i=1}^K E[T_i(n)]\Delta_i$ . In part B

we had  $\delta = \frac{1}{n^2}$ . So :

$$\begin{aligned} h_n &= \sum_{i=1}^K E[T_i(n)]\Delta_i \leq \sum_{i=1}^K \left( 3 + \frac{16 \log n}{\Delta_i^2} \right) \Delta_i \\ &= \sum_{i=1}^K 3\Delta_i + \sum_{i:\Delta_i > 0} \frac{16 \log n}{\Delta_i} \end{aligned}$$

D) By part B, C and lemma 2 we have :

$$h_n = \sum_{i=1}^K \Delta_i E[T_i(n)] = \sum_{i, \Delta_i < \Delta} \Delta_i E[T_i(n)] + \sum_{i, \Delta_i \geq \Delta} \Delta_i E[T_i(n)]$$

(For some  $\Delta_0$  to be determined later)

$$\leq n\Delta + \sum_{i, \Delta_i \geq \Delta} \left( 3\Delta_i + \frac{16 \log n}{\Delta_i} \right)$$

$$\leq n\Delta + 3 \sum_i \Delta_i + \frac{16K \log n}{\Delta}$$

$$\leq 8\sqrt{nK \log n} + 3 \sum_i \Delta_i$$

By choosing

$$\Delta = \sqrt{\frac{16K \log n}{n}}$$