

Curriculum Learning for Reinforcement Learning

Runxuan Jiang

1 Deterministic MDP

We first consider the case of using curriculum learning on an MDP with deterministic reward and transitions.

Proposition (Existence and bound for deterministic finite MDP.). Let $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, P, R, \rho)$ be a deterministic MDP with sparse reward, and that $R(s, a) = \mathcal{X}(s = g)$ where for some $g \in A$.

Then there exists a finite sequence of H MDP's $\mathcal{M}_i = (\mathcal{S}, \mathcal{A}, H, P_i, R_i, \rho_i)$ such that the following algorithm

1. Initialize π_0 to be a random policy.
2. For $1 \leq i \leq H$, Follow π_{i-1} with an ϵ_i -greedy policy until an optimal policy is found.

will find an optimal policy for \mathcal{M} in $O(H^2k)$ episodes in expectation, where $k = |\mathcal{A}|$.

Proof. Let $s_0, s_1, \dots, s_H = g$ be a trajectory that the optimal policy for \mathcal{M} takes. Construct \mathcal{M}_i such that $R_i(s, a) = \mathcal{X}(s = s_i)$ and $P_i(s_i | s_i, a) = 1 \ \forall a \in A$, otherwise $P_i = P$ for $1 \leq i \leq H$.

Let $p_i(\epsilon)$ denote the probability that the ϵ -greedy algorithm will achieve the optimal policy for MDP \mathcal{M}_i . Then we have

$$p_i(\epsilon) \geq f(\epsilon) := (1 - \epsilon)^{i-1} \left(\frac{\epsilon}{k} \right).$$

This comes from the fact that the optimal trajectory can be achieved by playing the greedy option $i - 1$ times, and then choosing the exploratory option and choosing the action that reaches the correct state, which has probability $\frac{\epsilon}{k}$.

For $i \geq 2$, f is concave, so we can find the maximal probability by setting the derivative to 0.

Note. Will need a prove on why this is concave.

$$p'_i(\epsilon) = 0$$

$$\begin{aligned}
&\Rightarrow \frac{(1-\epsilon)^{i-1}}{k} - (i-1)(1-\epsilon)^{i-2}(\frac{\epsilon}{k}) = 0 \\
&\Rightarrow \frac{(1-\epsilon)^{i-1}}{k} = (i-1)(1-\epsilon)^{i-2}(\frac{\epsilon}{k}) \\
&\frac{1-\epsilon}{k} = (i-1)(\frac{\epsilon}{k}) \Rightarrow \epsilon = \frac{1}{i}.
\end{aligned}$$

By plugging in this value for ϵ in f for \mathcal{M}_i , p is bounded below by

$$(1 - \frac{1}{i})^{i-1} \frac{1}{ki}.$$

So the expected number of episodes to reach the optimal policy for \mathcal{M}_i is

$$\begin{aligned}
&\frac{1}{(1 - \frac{1}{i})^{i-1} (\frac{1}{ki})} \\
&= \frac{ki}{(1 - \frac{1}{i})^{i-1}} = (i-1)k(\frac{i}{i-1})^i
\end{aligned}$$

Notice that since $i \geq 2$ and $x \mapsto (\frac{x}{x-1})^x$ is decreasing for $x > 0$, the above expression is bounded above by

$$4k(i-1).$$

So the number of episodes for the full curriculum is in expectation

$$\sum_{j=2}^H 4k(j-1) \leq \frac{H}{2}(4kH)$$

which is $O(kH^2)$. ■

Remark (Tighter bound). Note that the above proposition only considers the case where for each \mathcal{M}_i , the ϵ -greedy algorithm is greedy for $i-1$ steps, and then gets the correct state for the single final exploratory step. However, the curriculum learning also helps in the situation where some of the steps are exploratory and "lucky" enough to get the same action as the optimal policy, while the remaining steps are greedy.

However, even using this the bound is still $O(kH^2)$, so this is the tightest bound we can achieve under these conditions (probably).

Proof. Using a similar setup as the above proof, we now consider using the probability for achieving an optimal policy:

$$p(\epsilon) \geq f(\epsilon) := \binom{i-1}{0} (1-\epsilon)^{i-1} (\frac{\epsilon}{k}) + \binom{i-1}{1} (1-\epsilon)^{i-2} (\frac{\epsilon}{k})^2 + \dots + \binom{i-1}{i-1} (\frac{\epsilon}{k})^i$$

$$\begin{aligned}
&= \left(\frac{\epsilon}{k}\right) \sum_{j=0}^{i-1} \binom{i-1}{j} (1-\epsilon)^{i-1-j} \left(\frac{\epsilon}{k}\right)^j \\
&= \left(\frac{\epsilon}{k}\right) (1-\epsilon + \frac{\epsilon}{k})^{i-1}.
\end{aligned}$$

Note that this is concave on $(0, 1)$.

Note. **Proof of concavity needed.**

Since it is concave, we can set the derivative to 0 to find the ϵ that maximizes f . We have

$$\begin{aligned}
f'(\epsilon) &= \frac{(1-\epsilon + \frac{\epsilon}{k})^{i-1}}{k} + \left(\frac{\epsilon}{k}\right)(i-1)(1-\epsilon + \frac{\epsilon}{k})^{i-2} \left(\frac{1}{k} - 1\right) = 0 \\
&\Rightarrow \frac{(1-\epsilon + \frac{\epsilon}{k})^{i-1}}{k} = (1-i)\left(\frac{\epsilon}{k}\right)\left(\frac{1}{k} - 1\right)(1-\epsilon + \frac{\epsilon}{k})^{i-2} \\
&\Rightarrow (1-\epsilon + \frac{\epsilon}{k})^{i-1} = \epsilon(1-i)\left(\frac{1}{k} - 1\right)(1-\epsilon + \frac{\epsilon}{k})^{i-2} \\
&\Rightarrow 1-\epsilon + \frac{\epsilon}{k} = \epsilon(1-i)\left(\frac{1}{k} - 1\right) \\
&\Rightarrow 1-\epsilon + \frac{\epsilon}{k} = \frac{\epsilon}{k} - \epsilon + \frac{\epsilon i}{k} + i\epsilon \\
&\Rightarrow 1 = \epsilon\left(i - \frac{i}{k}\right) \\
&\Rightarrow \epsilon = \frac{1}{i - \frac{i}{k}} = \frac{k}{i(k-1)}
\end{aligned}$$

We then plug this value of ϵ back into f to get the maximal point of f :

$$\left(\frac{\frac{k}{i(k-1)}}{k}\right) \left(1 - \frac{k}{i(k-1)} + \frac{\frac{k}{i(k-1)}}{k}\right)^{i-1} = \frac{1}{i(k-1)} \left(1 - \frac{k-1}{i(k-1)}\right) = \frac{1 - \frac{1}{i}}{i(k-1)}$$

and take the reciporocal to get the expected number of episodes for \mathcal{M}_i , which is

$$\frac{i(k-1)}{1 - \frac{1}{i}} = \frac{i^2(k-1)}{i-1}$$

.

Note that $x \leq \frac{x^2}{x-1} \leq 4x$, so since $2 \leq i \leq H$ we get that

$$\sum_{j=2}^H \frac{j^2(k-1)}{j-1} = (k-1) \sum_{j=2}^H \frac{j^2}{j-1}$$

which is $O(kH^2)$.

To prove this is a tight bound, we can show that the probability given by $p(\epsilon)$ is a tight lower bound to the actual probability. ■

Proposition (Existence and bound for deterministic finite MDP). Let $\mathcal{M} = (S, A, H, P, R, \rho)$ be a deterministic MDP with sparse reward, and that $R(s, a) = \mathcal{X}(s = g)$ for some $g \in A$.

Then there exists a finite sequence of H MDP's $\mathcal{M}_i = (S, A, H, P_i, R_i, \rho_i)$ such that the following algorithm

1. Initialize π_0 to be a random policy.
2. For $1 \leq i \leq H$, Follow π_{i-1} with an ϵ_i -greedy policy until an optimal policy is found.

will find an optimal policy with probability $1 - \delta$ in at least x episodes.

Note. **TODO: Prove this theorem (same as above but finding the number of episodes to get an optimal policy with high probability).**

2 Curriculum Learning with Model-Based algorithms

Given a model-based reinforcement learning algorithm (such as RMAX or UCB), we want to show that we can solve MDP's using the curriculum learning paradigm on the MDP's used by the algorithm to "model" the MDP to be solved. Formally, such algorithms generate a sequence of MDP's M_1, \dots, M_n with corresponding optimal policies π_1^*, \dots, π_n^* . We want to show that by using π_{i-1}^* with ϵ -greedy on M_i we can achieve optimal regret bounds for M_i . Through induction, we therefore have a "curriculum" for learning the MDP.

2.1 RMAX

We start with an episodic nonstationery MDP $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, R, \rho)$. The R-max algorithm works by encouraging exploration on "unknown" states by setting the reward for those states to a maximal value R_{max} .

2.1.1 Original R-Max Algorithm

We first describe the RMAX algorithm for MDP's:

Algorithm 1 R-max for MDP's

Require: MDP $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, R, \rho), \epsilon > 0, 0 \leq \delta \leq 1, R_{max} > 0$.

```
1: procedure R-MAX( $M, \epsilon, \delta, R_{max}, K_1, K_2$ )
2:    $Rewards[S][A][H]$  be a new array.
3:    $Transitions[S][A][H]$  be a new array.
4:    $Known[S][A][H]$  be a new array.
5:    $\hat{M} \leftarrow M$  with all states replaced by absorbing states with reward  $R_{max}$  for any action.
6:    $Done \leftarrow False$ 
7:   while not  $Done$  do
8:      $Done \leftarrow True$ 
9:      $\pi \leftarrow$  optimal policy for  $\hat{M}$ 
10:    while  $\pi$  has not been run for more than  $K_1$  steps do
11:      obtain trajectories by running  $\pi$  on  $M$ .
12:      if a state  $(s, a, h)$  has been visited more than  $K_2$  times and  $Known[s][a][h]$  is
false then
13:         $Known[s][a][h] \leftarrow True$ 
14:        set the estimated transition probabilities  $Transitions[s][a][h]$  and rewards
 $Rewards[s][a][h]$  based on the average of the observed transitions and rewards for  $(s, a, h)$ 
15:        replace the transition probabilities and reward in  $(s, a, h)$  in  $\hat{M}$  by the
estimated values  $Transitions[s][a][h]$  and  $Rewards[s][a][h]$ 
16:         $Done \leftarrow False$ 
17:        break
18:      end if
19:    end while
20:  end while
21:  return The optimal policy for  $\hat{M}$ 
22: end procedure
```

Note. In the original R-max algorithm, the optimal policy for each M_j is not used since M_j is not known fully. Instead, the optimal policy for the estimated MDP for M_j , which is based on the samples collected previously, is used.

Remark (Is this curriculum learning?). Notice that this algorithm appears similar to curriculum learning. However there are some differences. In curriculum learning, the algorithm is given a sequence of tasks and trains an agent on the current task by using a (near)-optimal policy for the previous task. So given a sequence of MDP's M_1, \dots, M_n you explore M_j with the aid of the policy π_{j-1} . However, in the R-max algorithm we are exploring M directly using π_j .

2.1.2 R-max Algorithm with Curriculum Learning

We first introduce some definitions. Define $\mathcal{L} = \mathcal{P}(\mathcal{S} \times [H])$ to be the set of all possible sets of (nonstationary) states.

Definition. Let M be a nonstationary finite horizon MDP as defined above and let $L \in \mathcal{L}$.

Define the MDP \underline{M}_L to be the MDP with the same states and actions as M , and with the same rewards for all actions for all states $(s, h) \notin L$. For $(s, h) \in L$, the state is an absorbing state and the reward is set to R_{max} for all actions.

We now introduce the curriculum learning algorithm utilizing the exploration principles of the R-max algorithm. For simplicity, we assume that the starting state for M , s_0 is fixed: $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, R, s_0)$. We also assume that we have access to the MDP M_L for all $L \in \mathcal{L}$.

Algorithm 2 Curriculum Learning with R-max

Require: MDP $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, R, s_0)$, $\epsilon > 0$, $\delta \leq 1$, $R_{max} > 0$ and $\mathcal{M} = \{M_L : L \in \mathcal{L}\}$.

- 1: Initialize $L_0 = (\mathcal{S} \times [H]) \setminus \{(s_0, 1)\}$; $M_0 = M_{L_0}$; $\bar{\pi}_0$ being uniformly random policy; $(\bar{s}_0, \bar{h}_0) = (s_0, 1)$.
- 2: $\mathcal{H} \leftarrow \{\}$
- 3: **for** $t = 1, \dots, S \times H$ **do**
- 4: **for** $i = 1, \dots, K_1$ **do**
- 5: Run $\bar{\pi}_{t-1}$ with ϵ -greedy on M_{t-1} and add the trajectory to \mathcal{H} for states in $L \setminus L_{t-1}$.
- 6: **end for**
- 7: Call offline policy optimization on \mathcal{H} that returns $\bar{\pi}_t$, a near-optimal policy for M_{t-1} .
- 8: $\mathcal{H}_t \leftarrow \{\}$
- 9: **for** $i = 1, \dots, K_2$ **do**
- 10: Run π_t on M and add trajectory to \mathcal{H}_t
- 11: **end for**
- 12: **if** $\exists (s, h) \in L : (s, h)$ is reached by at least 1 trajectory in \mathcal{H}_t **then**
- 13: $(\bar{s}_t, \bar{h}_t) \leftarrow$ most visited state in L_{t-1}
- 14: **else**
- 15: **break**
- 16: **end if**
- 17: $L_t = L_{t-1} / (\bar{s}_t, \bar{h}_t)$; $M_t = M_{L_t}$.
- 18: **end for**
- 19: Output $\bar{\pi}_t$.

We now introduce the main theoretical guarantee for the above algorithm.

Theorem. Given a nonstationary finite horizon MDP M with fixed initial state, an $\epsilon > 0$, a $0 \leq \delta \leq 1$, the maximum reward for M , R_{max} and the \mathcal{M} corresponding to M , then running

algorithm 2 with $K_1 = \frac{-\ln(\frac{\delta}{2S})}{2(\frac{\epsilon}{SH^2R_{max}})^2} A^{\frac{1}{\frac{-\ln(\frac{\delta}{2S})}{2(\frac{\epsilon}{SH^2R_{max}})^2} A}}$, $K_2 = \frac{-2S^2H^2\ln(\delta)}{(\frac{\epsilon}{HR_{max}})^2}$, the algorithm will output an ϵ -optimal policy with probability $1 - \delta$. Furthermore, the sample complexity of the algorithm is $SH(K_1 + K_2)$.

We start off with some preliminary results about estimating the MDP, which will lead to a result justifying K_1 in line 5.

Theorem. Let M be a MDP and let M' be another MDP such that $\|P_M - P_{M'}\|_\infty < \frac{\epsilon}{SH^2R_{max}}$.

Then for all policies π ,

$$|V_M(\pi) - V_{M'}(\pi)| \leq \epsilon$$

Proof. Lemma 4 of the R-max paper. ■

Lemma. Let M be an MDP. To approximate the transition probability for a state s, a, h in M with error less than ϵ with probability higher than $1 - \delta$, we need $n \geq \frac{-\ln(\delta/2S)}{2\epsilon^2}$ samples.

Proof. Let n be the number of samples collected where the state (s, a, h) leads to the state s' . Let X_i be an indicator variable that is 1 if (s, a, h) leads to $(s', h + 1)$ and 0 otherwise in i th sample. Let p be the actual transition probability of transitioning from (s, a, h) to $(s', h + 1)$. Let $S = \frac{X_1 + \dots + X_n}{n}$. Note that $0 \leq \frac{X_i}{n} \leq \frac{1}{n}$. Since there are S possibilities for (s', h) , we want the probability that $|S - p| \geq \epsilon$ to be less than $\frac{\delta}{S}$. By Hoeffding's inequality we have

$$P(|S - p| \geq \epsilon) \leq 2 \exp\left(\frac{-2\epsilon^2}{n(\frac{1}{n})^2}\right) = 2 \exp(-2n\epsilon^2).$$

We want

$$\begin{aligned} 2 \exp(-2n\epsilon^2) &\leq \delta/S \\ \Rightarrow -2n\epsilon^2 &\leq \ln\left(\frac{\delta}{2S}\right) \\ \Rightarrow n &\geq \frac{-\ln(\frac{\delta}{2S})}{2\epsilon^2}. \end{aligned}$$
■

Lemma (Justification for K_1 (loose bound)). Fix $0 \leq \delta \leq 1$. In order to collect n samples for each action played at state (\bar{s}_t, \bar{h}_t) using $\bar{\pi}_{t-1}$ in the loop in line 4 with probability $1 - \delta$ or more, we need to set $K_1 = nA \frac{\ln(1-\delta \frac{1}{nA})}{\ln(\frac{p\epsilon}{SAH})}$.

Proof. The probability of $\bar{\pi}_{t-1}$ of reaching the state (\bar{s}_t, \bar{h}_t) is at least $\frac{p}{SH}$ with probability $1 - \delta$. The probability of playing any specific action from (\bar{s}_t, \bar{h}_t) given that (\bar{s}_t, \bar{h}_t) is reached is $\frac{\epsilon}{A}$. Thus the probability of reaching (\bar{s}_t, \bar{h}_t) and playing a specific action a is

$$\frac{\epsilon p}{SAH}.$$

We want to play each action at least n times. In this case, the probability for this is equivalent to playing a single action nA times. Let k be the proposed number of trajectories to collect. Let $X_i, 1 \leq i \leq k$ be a Bernoulli random variable that corresponds to a trajectory visiting (s, h) and performing action a on it. Then we have

$$P\left(\sum_{i=1}^k X_i \geq nA\right) \geq \prod_{j=1}^{nA} P\left(\sum_{i=1}^{\frac{k}{nA}} X_{j+i} \geq 1\right)$$

$$\begin{aligned}
&= \prod_{j=1}^{nA} (1 - P(\sum_{i=1}^{\frac{k}{nA}} X_{j*i} = 0)) \\
&= (1 - (\frac{p\epsilon}{SAH})^{\frac{k}{nA}})^{nA} \geq \delta \\
&\Rightarrow 1 - (\frac{p\epsilon}{SAH})^{\frac{k}{nA}} \geq \delta^{\frac{1}{nA}} \\
&\Rightarrow 1 - \delta^{\frac{1}{nA}} \geq (\frac{p\epsilon}{SAH})^{\frac{k}{nA}} \\
&\Rightarrow \ln(1 - \delta^{\frac{1}{nA}}) \geq \frac{k}{nA} \ln(\frac{p\epsilon}{SAH}) \\
&\Rightarrow k \geq nA \frac{\ln(1 - \delta^{\frac{1}{nA}})}{\ln(\frac{p\epsilon}{SAH})}
\end{aligned}$$

■

Next we prove a lemma showing that if the policy in line 9 is not already an optimal policy for the MDP M , then it must have a guaranteed probability of exploring L .

Lemma. Let $\alpha > 0$ let $L \subset \mathcal{S} \times [H]$. Let π_L be an optimal policy for M_L . Let $V_M(\pi)$ denote the value function on MDP M of policy π , and let V_M^* be the value of any optimal policy on M . Then either

1. $V_M(\pi_L) > V_M^* - \alpha$
2. A state in L will be played with probability of at least $\frac{\alpha}{H * R_{max}}$ when playing π_L on M (or M_L).

Proof. This proof follows the proof of lemma 6 in the original R-max paper.

Suffices to show that if (2) does not hold then (1) holds. Suppose that (2) does not hold. Let T be the set of all possible trajectories on M and let T' be the set of all trajectories that pass through L . Let $P_M^\pi(t)$ denotes the probability that the trajectory t will be reached when playing policy π on the MDP M . Since (2) does not hold we have that

$$\sum_{t \in T'} P_M^{\pi_L}(t) < \frac{\alpha}{H R_{max}}.$$

We want to show that

$$V_M^* - V_M(\pi_L) < \alpha.$$

Note that the only differences between M and M_L is that in M_L , the states in L are absorbing states and any action from this state gives a reward of R_{max} , but in M the reward is replaced by the actual reward and the transition probabilities are also replaced by the actual transition probabilities from M . Thus we have

$$\begin{aligned}
&V_M^* \leq V_{M_L}^* \\
&\Rightarrow V_M^* - V_M(\pi_L) \leq V_{M_L}^* - V_M(\pi_L)
\end{aligned}$$

We now further decompose each of these value functions over all possible trajectories. Let $V_M(t)$ denote the total reward of going through trajectory t on MDP M . Then we have

$$\begin{aligned}
|V_M(\pi_L) - V_{M_L}^*| &= \left| \sum_{t \in T} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right| \\
&= \left| \sum_{t \in T'} P_M^{\pi_L}(t) V_M(t) + \sum_{t \in T \setminus T'} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) - \sum_{t \in T \setminus T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right| \\
&\leq \left| \sum_{t \in T'} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right| + \left| \sum_{t \in T \setminus T'} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T \setminus T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right|
\end{aligned}$$

Since the rewards and transition probabilities for $t \in T \setminus T'$ are the same for MDPs M and M_L , the second expression above is equal to 0. Thus we are left with the first expression

$$= \left| \sum_{t \in T'} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right|.$$

Finally note that since the probabilities for all states before reaching something in L have the same transition probability for M and M_L , we have that $\sum_{t \in T'} P_M^{\pi_L}(t) = \sum_{t \in T'} P_{M_L}^{\pi_L}(t)$. We also have that $\forall t, 0 \leq V_M(t) \leq HR_{max}$ and $0 \leq V_{M_L}(t) \leq HR_{max}$. Thus we have

$$\left| \sum_{t \in T'} P_M^{\pi_L}(t) V_M(t) - \sum_{t \in T'} P_{M_L}^{\pi_L}(t) V_{M_L}(t) \right| \leq \sum_{t \in T'} P_M^{\pi_L}(t) HR_{max}.$$

However from our assumption we have that

$$\sum_{t \in T'} P_M^{\pi_L}(t) < \frac{\alpha}{HR_{max}}.$$

Substituting this into the previous expression we have that

$$V_M^* - V_M(\pi_L) \leq V_{M_L}^* - V_M(\pi_L) < \alpha$$

Which is the same statement as (1). ■

Lemma ((Justification for K_2)). In the loop in line 9, to guarantee that the resulting $\bar{\pi}_t$ would have a probability of reaching \bar{s}_t, \bar{h}_t of more than $\frac{p}{2SH}$, we need $K_2 \geq \frac{-2S^2H^2 \ln(\delta)}{p^2}$.

Proof. Let k be the number of times that the for loop in line 9 is run. Let $(s, h) \in L_{i-1}$. Define the indicator random variable X_i as follows:

$$X_i = \begin{cases} 1 & \text{trajectory passes through } (s, h) \\ 0 & \text{otherwise} \end{cases}.$$

Now define $S_k = \sum_{i=1}^k \frac{X_i}{k}$. Then we have $E[S_k] = P(\text{trajectory passes through } (s, h))$. By Hoeffding's inequality we have that

$$\begin{aligned}
P(E[S_k] - S_k \geq \frac{p}{2SH}) &\leq \exp\left(\frac{-2(\frac{p}{2SH})^2}{k(\frac{1}{k})^2}\right) \leq \delta \\
&\Rightarrow \frac{-kp^2}{2S^2H^2} \leq \ln(\delta) \\
&\Rightarrow \frac{-kp^2}{\ln(\delta)} \leq 2S^2H^2 \\
&\Rightarrow -k \leq \frac{2S^2H^2 \ln(\delta)}{p^2} \\
&\Rightarrow k \geq \frac{-2S^2H^2 \ln(\delta)}{p^2}
\end{aligned}$$

■

Proof (of main theorem). By the above lemma, for any $L \subset \mathcal{S} \times [H]$ we either have that the optimal policy for M_L is ϵ -optimal for M or a state in L will be played with probability of at least $\frac{\epsilon}{HR_{max}}$ using the current optimal policy. So we set $p = \frac{\epsilon}{HR_{max}}$.

Next, let M'_L be our estimated MDP for M_L . We need $\|P_{M_L} - P_{M'_L}\|_\infty < \frac{\epsilon}{SH^2R_{max}}$ which means that for each state we need (here we use σ -greedy)

$$n \geq \frac{-\ln(\frac{\delta}{2S})}{2(\frac{\epsilon}{SH^2R_{max}})^2}$$

samples. In order to satisfy this amount for every state in M_L that is not in L , and for every action played on each state, we would need

$$K_1 = nA \frac{\ln(1 - \delta^{\frac{1}{A}})}{\ln(\frac{p\sigma}{SAH})} = \frac{-\ln(\frac{\delta}{2S})}{2(\frac{\epsilon}{SH^2R_{max}})^2} A \frac{\ln(1 - \delta^{\frac{1}{2(\frac{\epsilon}{SH^2R_{max}})^2 A}})}{\ln(\frac{\epsilon}{SAH} \sigma)}$$

total samples. Since we need this number of samples for each iteration, the total sample complexity of the loop in line 4 is

$$\frac{-\ln(\frac{\delta}{2S})}{2(\frac{\epsilon}{SH^2R_{max}})^2} SAH \frac{\ln(1 - \delta^{\frac{1}{2(\frac{\epsilon}{SH^2R_{max}})^2 A}})}{\ln(\frac{\epsilon}{SAH} \sigma)}$$

Next we have that

$$K_2 \geq \frac{-2S^2H^2 \ln(\delta)}{p^2} = \frac{-2S^2H^2 \ln(\delta)}{(\frac{\epsilon}{HR_{max}})^2}.$$

Since there are SH iterations at most, the total sample complexity of the loop in line 9 is

$$\frac{-2S^3H^3 \ln(\delta)}{\left(\frac{\epsilon}{HR_{max}}\right)^2}$$

The total sample complexity is therefore

$$\frac{-2S^3H^3 \ln(\delta)}{\left(\frac{\epsilon}{HR_{max}}\right)^2} + \frac{-\ln\left(\frac{\delta}{2S}\right)}{2\left(\frac{\epsilon}{SH^2R_{max}}\right)^2} SAH \frac{\ln\left(1 - \delta^{\frac{1}{2\left(\frac{\epsilon}{SH^2R_{max}}\right)^2 A}}\right)}{\ln\left(\frac{\epsilon}{SH^2R_{max}}\right)}$$

The algorithm is optimal because the resulting M_L would have a value function ϵ -close to M 's value function, and the estimated M'_L would have a value function ϵ -close to M_L 's value function, all with probability $1 - \delta$. ■