

# 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  $\pi_0$  to be a random policy.
2. For  $1 \leq i \leq H$ , Follow  $\pi_{i-1}$  with an  $\epsilon_i$ -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  $\epsilon$ -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  $\epsilon$  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  $\epsilon$ -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  $(0, 1)$  (proof needed).

Since it is concave, we can set derivative to 0 to find the  $\epsilon$  that maximizes the probability.

We get that

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

We then plug back into  $p$  and take the reciprocal to get the expected number of episodes for  $\mathcal{M}_i$ , which is

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

.

Note that  $x \leq \frac{x^2}{x-1} \leq 4x$

so this also leads to a bound of  $O(kH^2)$  if we sum across all  $1 \leq i \leq H$ .

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  $\pi_0$  to be a random policy.
2. For  $1 \leq i \leq H$ , Follow  $\pi_{i-1}$  with an  $\epsilon_i$ -greedy policy until an optimal policy is found.

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

*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  $P_i = P$ .

Then for each  $\mathcal{M}_i$ , the probability that the  $\epsilon$ -greedy algorithm will achieve the optimal policy is

$$\begin{aligned}
p(\epsilon) &= \binom{i-1}{0} (1-\epsilon)^{i-1} \left(\frac{\epsilon}{k}\right) + \binom{i-1}{1} (1-\epsilon)^{i-2} \left(\frac{\epsilon}{k}\right)^2 + \dots + \binom{i-1}{i-1} \left(\frac{\epsilon}{k}\right)^i \\
&= \left(\frac{\epsilon}{k}\right) \sum_{j=0}^{i-1} \binom{i-1}{j} (1-\epsilon)^{i-1-j} \left(\frac{\epsilon}{k}\right)^j
\end{aligned}$$

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

For  $i \geq 2$ ,  $p$  is concave (need a proof of this), so we can find the maximal probability by setting the derivative to 0.

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

So for  $\mathcal{M}_i$ , the probability of reaching the optimal policy in each episode is

$$\begin{aligned} p_i(\epsilon^*) &= \left(\frac{\frac{k}{i(k-1)}}{k}\right)\left(1 - \frac{k}{i(k-1)} + \frac{\frac{k}{i(k-1)}}{k}\right)^{i-1} \\ &= \frac{i-1}{i^2(k-1)} \end{aligned}$$

■

## 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  $\pi_1^*, \dots, \pi_n^*$ . We want to show that by using  $\pi_{i-1}^*$  with  $\epsilon$ -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 nonstationary MDP  $M = (\mathcal{S}, \mathcal{A}, H, \mathbb{P}, R, \rho)$ . Let  $L$  be a collection of states and actions  $L \subset \mathcal{S} \times \mathcal{A} \times [H]$  corresponding to the set of states marked unknown in the RMAX algorithm.

**Definition** ( $M_L$ ). Let  $M$  and  $L$  be defined as above. We define  $M_L$  as a MDP with the same states and actions as  $M$ , and with the same rewards and transitions for  $(s, a, h) \notin L$ . For  $(s, a, h) \in L$ , the the reward will be equal to  $R_{max}$  and the transition will lead to an absorbing state that always gives reward  $R_{max}$ .

The RMAX algorithm works by starting with  $L_0 = \mathcal{S} \times \mathcal{A}$ , and using the optimal policy for  $M_{L_0}$  on  $M$ . Once enough samples have been collected from a state-action pair  $(s, a, h)$ , it is removed from  $L_0$  to get  $L_1 = L_0 \setminus \{(s, a, h)\}$ , and then using the optimal policy for  $M_{L_1}$ . This is done until there are no unknown states left (i.e.,  $L$  is empty).

Note that

$$\mathcal{S} \times \mathcal{A} \times [H] = L_0 \supset L_1 \supset \dots \supset L_{SAH} = \emptyset$$

Denote  $M_{L_j}$  as  $M_j$ . The following describes the main theorem.

**Theorem.** Fix  $0 < \delta < 1, \epsilon > 0$ . Then we can find an  $\epsilon$ -optimal policy with probability higher than  $1 - \delta$  by performing the following steps on the sequence of MDP's  $M_1, \dots, M_{SAH}$  with corresponding optimal policies  $\pi_1^*, \pi_{SAH}^*$ :

- Initialize  $\pi_0$  to be a random policy.
- For  $t = 1, \dots, n$ , follow  $\pi_{t-1}$  with  $\epsilon$ -greedy exploration to collect  $K_2$  trajectories. Compute the optimal policy  $\pi_t$  from the model learned by the trajectories.
- Output  $\pi_{SAH}$ .

**Lemma.** Let  $\alpha > 0$  and  $1 \leq j \leq SA$ . Denote  $(s, a)$  to be the entry in  $L_{j-1}$  but not in  $L_j$ . Consider playing the policy  $\pi_{j-1}^*$  with  $\epsilon$ -greedy on  $M_j$ . Let  $V_M(\pi)$  denote the value function on MDP  $M$  of policy  $\pi$ , and let  $V_M^*$  be the value of any optimal policy on  $M$ . Then either

1.  $V_{M_j}(\pi_{j-1}) > V_{M_j}^* - \alpha$
2.  $(s, a, h)$  will be played with probability of at least  $\frac{\alpha}{H^* R_{max}}$ .

*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  $M_j$  and let  $T'$  be the set of all trajectories that pass through  $(s, a, h)$ . Let  $P_M^\pi(t)$  denotes the probability that the trajectory  $t$  will be reached when playing policy  $\pi$  on the MDP  $M$ . Since (2) does not hold we have that

$$\sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) < \frac{\alpha}{H R_{max}}.$$

We want to show that

$$V_{M_j}^* - V_{M_j}(\pi_{j-1}) < \alpha.$$

Note that the only difference between  $M_j$  and  $M_{j-1}$  is that in  $M_{j-1}$ , the state  $(s, h)$  is an absorbing state and any action from this state gives a reward of  $R_{max}$ , but in  $M_j$  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

$$V_{M_j}^* \leq V_{M_{j-1}}^*.$$

$$\Rightarrow V_{M_j}^* - V_{M_j}(\pi_{j-1}) \leq V_{M_{j-1}}^* - V_{M_j}(\pi_{j-1})$$

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

$$|V_{M_j}(\pi_{j-1}) - V_{M_{j-1}}^*| = \left| \sum_{t \in T} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right|$$

$$\begin{aligned}
&= \left| \sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) + \sum_{t \in T \setminus T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) - \sum_{t \in T \setminus T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right| \\
&\leq \left| \sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right| + \left| \sum_{t \in T \setminus T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T \setminus T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right|
\end{aligned}$$

Since the rewards and transition probabilities for  $t \in T \setminus T'$  are the same for MDPs  $M_j$  and  $M_{j-1}$ , the second expression above is equal to 0. Thus we are left with the first expression

$$= \left| \sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right|.$$

Finally note that since the probabilities for all states before step  $h$  have the same transition probability for  $M_j$  and  $M_{j-1}$ , we have that  $\sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) = \sum_{t \in T'} P_{M_{j-1}}^{\pi_{j-1}}(t)$ . We also have that  $\forall t, 0 \leq V_{M_j}(t) \leq HR_{max}$  and  $0 \leq V_{M_{j-1}}(t) \leq HR_{max}$ . Thus we have

$$\left| \sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) V_{M_j}(t) - \sum_{t \in T'} P_{M_{j-1}}^{\pi_{j-1}}(t) V_{M_{j-1}}(t) \right| \leq \sum_{t \in T'} P_{M_j}^{\pi_{j-1}} HR_{max}.$$

*Note.* Not sure if the above argument is valid for this inequality, but a similar inequality was proved in the paper and the details were omitted.

However from our assumption we have that

$$\sum_{t \in T'} P_{M_j}^{\pi_{j-1}}(t) < \frac{\alpha}{HR_{max}}.$$

Substituting this into the previous expression we have that

$$V_{M_j}^* - V_{M_j}(\pi_{j-1}) < \alpha$$

Which is the same statement as (1). ■

**Lemma.** Fix  $0 < \delta < 1$ . Then if we obtain  $K_1 = \max((\frac{4SHR_{max}}{\epsilon})^3, -6 \ln^3(\frac{\delta}{6SA^2})) + 1$  samples for  $(s, a)$  then the estimate for the transition probability for  $(s, a)$  is within  $\frac{\epsilon}{2*SHR_{max}}$  of the actual probabilities.

*Proof.* Use Chernoff bound and pigeonhole principle. ■

**Lemma.** Fix  $0 < \delta < 1$  and  $\epsilon > 0$ . Then we can find an  $\epsilon$ -optimal policy with probability higher than  $1 - \delta$  for MDP  $M_j$  after collecting trajectories by running policy  $\pi_{j-1}$  with  $\epsilon$ -greedy exploration for  $K_2$  episodes where  $K_2^{\frac{1}{3}} + K_2 \frac{\alpha}{H*HR_{max}} > A^2SK_1$  and computing the optimal policy for that model.

*Note.* We may require less than this number since we only have a single unknown state.

*Proof.* By the previous two lemmas, this number of steps ensures that  $K_1$  samples are obtained for  $(s, a)$  with high (greater than  $1 - \delta$ ) probability. Thus, we can accurately estimate the MDP  $M_j$  with a TV distance of less than  $\frac{\epsilon}{2*SHR_{max}}$ . By lemma 4 in the R-max paper, this means that the optimal value function of the estimated MDP and  $M_j$  are  $\epsilon$  close to each other. So calculate the optimal policy for the estimated MDP to get the resulting policy. ■

*Note.* In this case we don't really use  $\epsilon$ -greedy (we set  $\epsilon$  to zero).

*Proof (of Theorem).* Follows from previous lemma. ■

## 2.2 UCRL