Understanding Curriculum Learning in Policy Optimization for Online Combinatorial Optimization

Runlong Zhou

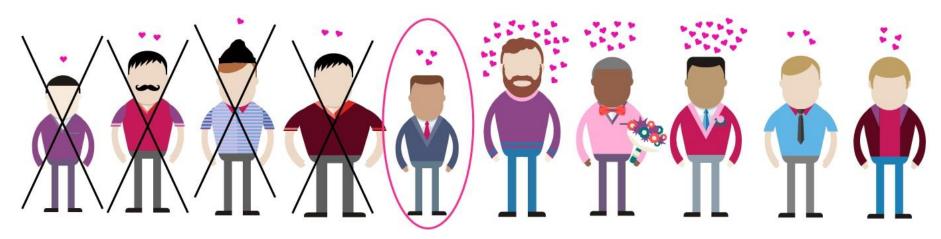
Joint work with Yuandong Tian, Yi Wu and Simon Du

Motivation

- Machine Learning is good at Combinatorial Optimization problems
 - Is (any part of) this success explainable?
- Online CO matches the nature of Reinforcement Learning
 - Sequential decision-making
- Theoretical understanding of RL techniques on online CO
 - Curriculum Learning

Example 1

- Secretary Problem (a.k.a. Marriage Problem, Sultan's Dowry Problem)
 - Hire one secretary among n candidates, each with (different) ability score
 - Arrive sequentially, but the order is unknown
 - Once reject someone, cannot come back and re-hire
 - Once hire someone, ends
 - Maximize the probability of hiring the best candidate



Example 1 (cont'd)

- Secretary Problem
 - Isomorphic up to only the relative rank
 - n! Permutations (instances), each with equal probability
 - Find a policy working averagely well on the instance distribution

Example 2

- Online Knapsack
 - *n* items arrive sequentially
 - Each with value v_i and size s_i revealed upon arrival
 - Total budget B
 - Once reject something, cannot come back and re-pick
 - To pick something, requires total size $\leq B$
 - Maximize total value of picked items
 - Distribution of instances defined via each $(v_i, s_i) \sim F_v \times F_s$ i.i.d.

Formulation

- Latent MDPs [Kwon et al.,2021]
 - A distribution of MDPs $\mathcal{M} = \{\mathcal{M}_1, ..., \mathcal{M}_M\}$
 - Each with weight (probability) $w_1, ..., w_M$
 - All MDPs share **states** *S*, **actions** *A*, **horizon** *H*
 - But have their own initial state distribution ν_m reward function r_m (bounded by [0,1]), transition P_m
- Why use this setting?
 - Characterize instances and working averagely well

Formulation (cont'd)

- Policy class
 - Stationary $\pi: S \to \Delta(A)$
 - Log-linear parameterization
 - Assume we have a fixed feature mapping $\phi: S \times A \to \mathbb{R}^d$
 - Learn parameter θ , and policy

$$\pi(a|s) = \frac{\exp(\theta^{\mathsf{T}}\phi(s,a))}{\sum_{a'\in A} \exp(\theta^{\mathsf{T}}\phi(s,a'))}$$

Formulation (cont'd)

Value functions and entropy regularization

•
$$V_{i,h}^{\pi,\lambda}(s) = \mathbb{E}\left[\sum_{t=0}^{h-1} r_i(s_t, a_t) + \lambda \ln \frac{1}{\pi(a_t|s_t)} \mid \mathcal{M}_i, \pi, s_0 = s\right]$$

• $Q_{i,h}^{\pi,\lambda}(s) = \mathbb{E}\left[\sum_{t=0}^{h-1} r_i(s_t, a_t) + \lambda \ln \frac{1}{\pi(a_t|s_t)} \mid \mathcal{M}_i, \pi, (s_0, a_0) = (s, a)\right]$

•
$$A_{i,h}^{\pi,\lambda}(s) = Q_{i,h}^{\pi,\lambda}(s) - V_{i,h}^{\pi,\lambda}(s)$$

- Goal
 - Find $\pi_{\lambda}^{\star} = \operatorname{argmax}_{\pi} \sum_{i} w_{i} \sum_{s_{0}} \nu_{i}(s_{0}) V_{i,h}^{\pi,\lambda}(s_{0})$

Prerequisites

- State-action visitation distribution
 - $d_{i,h}^{\pi}(s) = \mathbb{P}(s_h = s \mid \mathcal{M}_i, \pi)$
 - $d_{i,h}^{\pi}(s,a) = d_{i,h}^{s_0,\pi}(s)\pi(a|s)$
- Function approximation loss

$$L(g; \theta, v) = \sum_{i=1}^{M} w_i \sum_{h=1}^{H} \underset{s, a \sim v_{i, H-h}}{\mathbb{E}} \left[\left(A_{i, h}^{\pi_{\theta}, \lambda}(s, a) - g^{\top} \nabla \ln \pi_{\theta}(a|s) \right)^2 \right]$$

Fisher information matrix

$$\Sigma_{v}^{\theta} = \sum_{i=1}^{M} w_{i} \sum_{h=1}^{H} \underset{s, a \sim v_{i, H-h}}{\mathbb{E}} \left[\nabla_{\theta} \ln \pi_{\theta}(a|s) \left(\nabla_{\theta} \ln \pi_{\theta}(a|s) \right)^{\top} \right]$$

Natural Policy Gradient

- Initial param: θ_0
- Update:

$$\theta_{t+1} = \theta_t + \eta g_t$$

• $g_t \approx \operatorname{argmin}_g L(g; \theta_t, d^t)$

```
Algorithm 3 NPG: Sample-based NPG (full version).
```

```
1: Input: Environment E; learning rate \eta; episode number T; batch size N; initialization \theta_0; sampler \pi_s;
      regularization coefficient \lambda; entropy clip bound U; optimization domain \mathcal{G}.
 2: for t \leftarrow 0, 1, ..., T - 1 do
          Initialize \widehat{F}_t \leftarrow 0^{d \times d}, \widehat{\nabla}_t \leftarrow 0^d.
         for n ← 0, 1, . . . , N − 1 do
             for h ← 0, 1, . . . , H - 1 do
                  if \pi_s is not None then
                      s_h, a_h, \widehat{A}_{H-h}(s_h, a_h) \leftarrow \text{Sample}(E, \pi_s, \text{True}, \pi_t, h, \lambda, U) \text{ (see Alg. 4)}.
                     //s, a \sim \widetilde{d}_{m,h}^{\pi_s}, estimate A_{m,H-h}^{t,\lambda}(s,a).
                  else
                      s_h, a_h, \widehat{A}_{H-h}(s_h, a_h) \leftarrow \text{Sample}(E, \pi_t, \text{False}, \pi_t, h, \lambda, U).
                     //s, a \sim d_{m,h}^{\theta_t}, estimate A_{m,H-h}^{t,\lambda}(s,a).
                  end if
10:
              end for
11:
              Update:
12:
                                                       \widehat{F}_t \leftarrow \widehat{F}_t + \sum_{h=0}^{n-1} \nabla_{\theta} \ln \pi_{\theta_t}(a_h|s_h) \left( \nabla_{\theta} \ln \pi_{\theta_t}(a_h|s_h) \right)^{\top},
                                                       \widehat{\nabla}_t \leftarrow \widehat{\nabla}_t + \sum_{h=0}^{H-1} \widehat{A}_{H-h}(s_h, a_h) \nabla_{\theta} \ln \pi_{\theta_t}(a_h | s_h).
          end for
          Call any solver to get \widehat{g}_t \leftarrow \arg\min_{g \in \mathcal{G}} g^{\top} \widehat{F}_t g - 2g^{\top} \widehat{\nabla}_t.
          Update \theta_{t+1} \leftarrow \theta_t + \eta \widehat{g}_t.
16: end for
17: Return: \theta_T.
```

Analysis

Assumption

- g_t^{\star} is the true minimizer of L at time t
- d^* is the visitation distribution of the optimal policy
- $||\phi(s,a)||_2 \le B$

Definition 4. Define for $0 \le t \le T$:

- (Excess risk) $\epsilon_{\text{stat}} := \max_t \mathbb{E}[L(g_t; \theta_t, d^t) L(g_t^{\star}; \theta_t, d^t)];$
- (Transfer error) $\epsilon_{\text{bias}} := \max_t \mathbb{E}[L(g_t^{\star}; \theta_t, d^{\star})];$
- (Relative condition number) $\kappa := \max_t \mathbb{E}\left[\sup_{x \in \mathbb{R}^d} \frac{x^\top \Sigma_{d^*}^{\theta_t} x}{x^\top \Sigma_t x}\right]$. Note that term inside the expectation is a random quantity as θ_t is random.

The expectation is with respect to the randomness in the sequence of weights g_0, g_1, \ldots, g_T .

Analysis (cont'd)

Main result

Theorem 6. With Def. 4, 5 and 8, our algorithm enjoys the following performance bound:

$$\mathbb{E}\left[\min_{0\leq t\leq T}V^{\star,\lambda}-V^{t,\lambda}\right]\leq \frac{\lambda(1-\eta\lambda)^{T+1}\Phi(\pi_0)}{1-(1-\eta\lambda)^{T+1}}+\eta\frac{B^2G^2}{2}+\sqrt{H\epsilon_{\text{bias}}}+\sqrt{H\kappa\epsilon_{\text{stat}}},$$

where $\Phi(\pi_0)$ is the Lyapunov potential function which is only relevant to the initialization.

Interpretation

- Linear convergence ($\lambda = 0$ gets $1/\sqrt{T}$ rate)
- $\epsilon_{\rm bias}$ depends on the **design of feature mapping**, and is hardly removable
- $\epsilon_{\rm stat}$ depends on the closeness between g_t and g_t^{\star} , can be reduced with a bigger batch-size
- κ shows the closeness between d^{\star} and d^{t} , and is of our special interest

Curriculum Learning [Bengio et al., 2009]

- First learn small scale problems, then use this model to initialize large scale learning
- Mentioned as "Bootstrapping" in [Kong et al.,2019]
 - They suggested a series of curricula: 10, 20, ..., 200
 - No need to do this, because **reducing** κ **is sufficient**
- Even works for changing distribution
 - Uniform over 10! instances
 - Arbitrary complex for 200! Instances
 - As long as the optimal policies are "Similar", it transfers well

Curriculum Learning for SP

- Distribution modeling
 - Suppose for each candidate i, it has probability P_i to be the so-far best
 - For classical SP, $P_i = 1/i$
- Optimal policy
 - Always a **p-threshold policy**: accept if and only if i/n > p and is so-far best
 - For classical SP, p = 1/e

Curriculum Learning for SP (cont'd)

• Comparing κ under with/without curriculum learning

Theorem 7. Assume that each candidate is independent of others and the i-th candidate has a probability P_i of being the best so far (Sec. 4.1). Assume the optimal policy is a p-threshold policy and the sampling policy is a q-threshold policy. There exists a policy parameterization such that:

$$\kappa_{\text{curl}} = \Theta\left(\begin{cases} \prod_{j=\lfloor nq\rfloor+1}^{\lfloor np\rfloor} \frac{1}{1-P_j}, & q \leq p, \\ 1, & q > p, \end{cases}\right),$$

$$\kappa_{\text{na\"{i}ve}} = \Theta\left(2^{\lfloor np\rfloor} \max\left\{1, \max_{i \geq \lfloor np\rfloor+2} \prod_{j=\lfloor np\rfloor+1}^{i-1} 2(1-P_j)\right\}\right),$$
(1)

where κ_{curl} and $\kappa_{\text{na\"ive}}$ are κ of the sampling policy and the na\"ive random policy, respectively.

Curriculum Learning for SP (cont'd)

- Classical case
 - The target problem is the classical SP

$$\kappa_{ ext{curl}} = \left\{ egin{array}{l} rac{\lfloor n/\mathrm{e}
floor}{\lfloor nq
floor}, & q \leq rac{1}{\mathrm{e}}, \ 1, & q > rac{1}{\mathrm{e}}, \end{array}
ight. \quad \kappa_{ ext{na\"{i}ve}} = 2^{n-1} rac{\lfloor n/\mathrm{e}
floor}{n-1}.$$

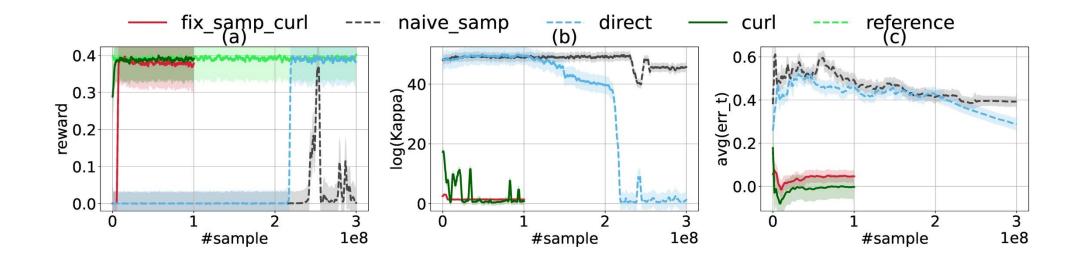
- General case
 - The target problem satisfies $P_i \leq 1/2$

$$\kappa_{ ext{curl}} \leq \left\{ egin{array}{ll} 2^{\lfloor np \rfloor - \lfloor nq
floor}, & q \leq p, \ 1, & q > p, \end{array}
ight. \quad \kappa_{ ext{na\"{i}ve}} \geq 2^{\lfloor np
floor}.$$

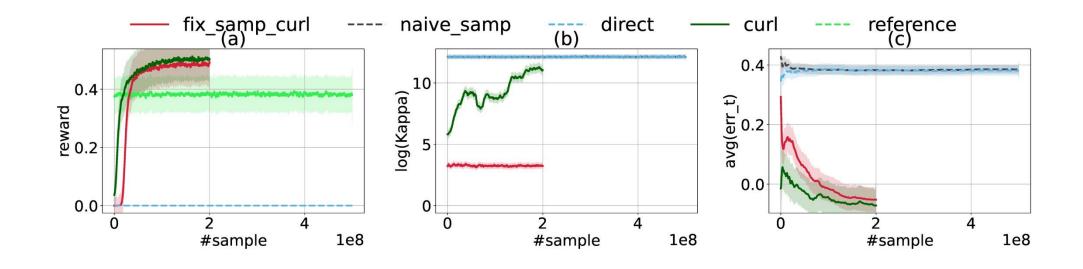
- Failure case
 - Best candidate always come as the last one $\kappa_{\text{curl}} = \infty$, $\kappa_{\text{larger than }\kappa_{\text{na\"{i}ve}}} = 2^{n-1}$.

$$\kappa_{\rm curl} = \infty$$
, larger than $\kappa_{\rm na\"{i}ve} = 2^{n-1}$.

Experiments – Secretary Problem



Experiments – Online Knapsack



References

- Semih Cayci, Niao He, and R. Srikant. Linear convergence of entropy-regularized natural policygradient with linear function approximation, 2021.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. InProceedings of the 26th Annual International Conference on Machine Learning, ICML '09, page41–48, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605585161.doi: 10.1145/1553374.1553380. URLhttps://doi.org/10.1145/1553374.1553380.
- Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policygradient methods: Optimality, approximation, and distribution shift, 2020.
- Weiwei Kong, Christopher Liaw, Aranyak Mehta, and D. Sivakumar. A new dog learns old tricks: Rlfinds classic optimization algorithms. InICLR, 2019.
- Jeongyeol Kwon, Yonathan Efroni, Constantine Caramanis, and Shie Mannor. RI for latent mdps:Regret guarantees and a lower bound, 2021.