
CS/STAT 184(0) Final Report

Christopher Lapop Salazar
Eric Peterson
Michelle Nhung Le
clapopsalazar@college.harvard.edu
peterson@g.harvard.edu
michellele@college.harvard.edu

Abstract

1 We developed a reinforcement learning model for hobby card collecting, focusing
2 on optimizing card pulls and completing specific collections within a budget. Using
3 modified Thompson Sampling and Proximal Policy Optimization, our approaches
4 achieve sublinear regret and efficient policy learning. The results show robust
5 performance for both diverse and rare card collections.

6 1 Project Phases

7 1.1 Phase 1: Unknown pack pull rates, minimize regret with different pack choices

8 In Phase 1, we adapt a Bernoulli Bandit problem to hobby card-collecting, where each pack pull
9 generates cards of varying market values. Our goal is to maximize the value of collected cards by
10 learning optimal pull strategies using a reinforcement learning agent.

11 We extend Thompson Sampling for multinomial outcomes, using Dirichlet priors to model pull
12 probabilities. The algorithm iteratively updates these priors based on observed rewards.

13 1.2 Phase 2: Maximize probability of completing a specific deck of cards within set budget

14 Phase 2 introduces selling individual cards at market value, allowing the agent to extend the time
15 horizon. This marks a change from an MAB problem towards an MDP problem, as the current budget
16 and cards drawn so far are used to make decisions. The state space encodes the cards that have been
17 pulled so far and the remaining budget for buying cards. The action space consists of the arms to pull
18 from Phases 1-2, and additionally the option to sell cards to a merchant to recoup money for more
19 card packs.

20 Actions that result in acquiring needed cards could generate rewards, while actions that don't result
21 in acquiring needed cards would generate either no reward or smaller rewards if they give high-value
22 cards that can be sold to buy more card packs. The distributions of cards granted by the different
23 packs are encoded in transition probabilities of the MDP. The state space consists of all possible
24 combinations of cards, which is quite high-dimensional, so we will look into different exploration
25 and policy gradient methods for learning the MDP and finding an optimal policy. An additional
26 ability to purchase cards at market value price is being considered, which would allow guaranteed
27 (but possibly overly costly) steps towards completing the specific deck of cards.

We can optimize card collection utility by applying Proximal Policy Optimization. Environment settings are: **State Space:** $s_t = [\mathbf{x}_t, b_t]$, where:

$\mathbf{x}_t \in \mathbb{N}^N$ is current collection (count of each card) $b_t \in \mathbb{R}$ is current budget (N is total number of unique cards)

Action Space: $\mathcal{A} = \mathcal{A}_{\text{pack}} \cup \mathcal{A}_{\text{trade}}$ where:

- Buy pack: $\mathcal{A}_{\text{pack}} = a_{\text{pack}}(k) \mid k \in [K]$ for K different pack types
- Trade cards: $\mathcal{A}_{\text{trade}} = a_{\text{buy}}(i), a_{\text{sell}}(i) \mid i \in [N]$ for buying/selling each card

Reward Function:

$$R(s_t, a_t, s_{t+1}) = R_{\text{collection}} + R_{\text{budget}}$$

where:

- $R_{\text{collection}} = 10.0 \cdot \Delta\text{completion}$
- $R_{\text{budget}} = \begin{cases} 0 & \text{if } \Delta b_t > 0 \\ -0.1 \cdot \Delta b_t & \text{if } \Delta b_t < 0 \end{cases}$
- $\Delta\text{completion}$ is the change in number of needed cards collected
- Δb_t is the change in budget

28 This reward function effectively measures the progress made towards completing the collection by
 29 considering the reduction in the remaining cards needed after updating the collection with the drawn
 30 cards and cares less about cost because cost is scaled by -0.1 .

31 2 Comparison to other works

32 The second phase of our project addresses the challenge of collecting a series of unique objects, each
 33 with an independent and unchanging probability of being obtained. This setup serves as a variation of
 34 the classical coupon collector problem. Our "completionist" test case functions in the same manner,
 35 with the objective of pulling a minimum of one of each unique trading card. When comparing our
 36 implementation with the report by Emma Brunskill and Lihong Li, the report develops an algorithm
 37 FORCEDEXP, designed for Lifelong Multitask Reinforcement Learning. FORCEDEXP tackles the
 38 problem by emphasizing the design of an exploration strategy that can identify and adapt to new
 39 tasks. In contrast, our work with a Bayesian policy model and PPO applies single-task reinforcement
 40 learning methods, prioritizing a specific task model. By refining and optimizing an existing policy, it
 41 focuses on the completion of a single predefined task - our specified collection.

42 3 Phase 1

43 3.1 Environment

44 We wrote an environment class called PACKPOOL from scratch, which will also serve as a core for
 45 the environment in the later phases of our project. The class contains num_packs "packs" objects,
 46 which correspond to the card packs with randomly-sampled cards. In order to make the environment
 47 correspond to some realistic card-collecting scenarios, there are 3 sub-categories of cards ("common,"
 48 "uncommon," and "foil") which each have different ranges of values. The class stores the lists of cards
 49 that each pack can possibly contain, the probabilities of drawing them, and each of their individual
 50 values. The PACKPOOL class also includes methods to open packs (effectively, sample according to
 51 the card probabilities within each respective pack) and calculate the values of sets of cards.

52 3.2 Algorithm

53 The problem of trying to maximize the value of collected cards by opening packs with a priori
 54 unknown card distributions is conceptually similar to the multi-armed bandit problem that we
 55 studied in class. The most effective approach we learned for this kind of problem was Thompson
 56 sampling, but in order to apply this algorithm to our problem we need to extend it to deal with an
 57 important distinction from Bernoulli bandits: the arms now output a vector drawn from a multinomial
 58 distribution, with the reward being a weighted sum of the categories drawn (here, the categories are
 59 the cards and the weights are their individual values).

60 Since the essence of the Thompson sampling algorithm is to maintain Bayesian priors on the arm
 61 distributions and sample them according to which ones seem most promising given all existing
 62 information, this suggests that we can straightforwardly adapt it if we have a conjugate prior for
 63 the multinomial distribution and some way of ranking arms to compute the $\arg \max_k$ over packs. It
 64 turns out that the Dirichlet distribution $\text{Dir}(\alpha_k)$ is a conjugate prior of $\text{Mult}(K, \mathbf{p})$ and has a similar
 65 Bayesian update rule to that of the Beta distribution, which is just to increment α_k^j by one each time
 66 pack k outputs card j .

67 Finally, we define the reward of a pack pull \mathbf{n}_t as the total value of all the cards drawn, i.e. $r_t := \mathbf{n}_t \cdot \mathbf{v}$,
 68 where $\mathbf{v} \in \mathbb{R}^J$ is a vector containing all of the individual card values. Then, the regret is just the
 69 cumulative sum of differences between the expected value of the pulled packs and the pack with the
 70 highest expected value.

71 In summary, our new algorithm can be written as:

Algorithm 1 Modified Thompson Sampling for Packs

Require: Number of packs K , number of cards J , horizon $T \geq K$

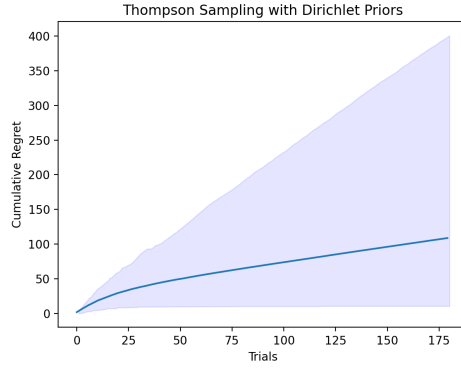
```

1: Initialize  $\alpha_k^{(j)} = 1 \in \mathbb{R}^{J,K}$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   for  $k = 1, 2, \dots, K$  do
4:     Sample  $\hat{\theta}_k \sim \text{Dir}(\alpha_k)$ 
5:     Compute reward map  $\hat{r}_k = r(\hat{\theta}_k)$ , where  $r$  is the function that maps the hand to its value,
        $\mathbb{R}^J \rightarrow \mathbb{R}$ 
6:   end for
7:   Open pack  $a_t = \arg \max_{k \in \{1, \dots, K\}} \hat{r}_k$ , where ties are broken uniformly at random
8:   Update the belief for pack  $a_t$  using the drawn cards  $\mathbf{n}_t$ :
       
$$(\alpha_{a_t}^{(1)}, \alpha_{a_t}^{(2)}, \dots, \alpha_{a_t}^{(J)}) \leftarrow (\alpha_{a_t}^{(1)} + n_t^{(1)}, \alpha_{a_t}^{(2)} + n_t^{(2)}, \dots, \alpha_{a_t}^{(J)} + n_t^{(J)})$$

9: end for
```

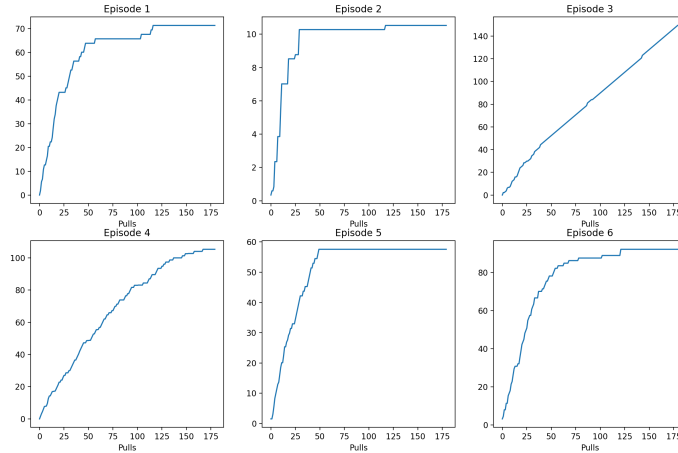
72 3.3 Data generation

73 We ran our algorithm on the PACKPOOL class for 100 episodes and a horizon of $T = 180$ pack pulls.
 74 The mean regret with 95 percent confidence intervals is plotted below:



75

76 The algorithm is evidently able to achieve sublinear regret, although with a fairly large variance.
 77 Plotting the first 6 episodes suggests that algorithm's behavior is fairly "streaky:" it is often able to
 78 quickly identify the optimal pack and almost completely stops exploring (Episodes 1, 2, 5 and 6) but
 79 occasionally stops exploring with a high-confidence prior on the incorrect optimal pack (Episode 3):



80 4 Phase 2

81 4.1 Environment

82 In approaching Phase 2, the environment makes use of the `PACKPOOL` class with modifications to
 83 the `open_pack_list` function to provide the specific cards produced from pack openings rather
 84 than only producing the card values. Most notably, Phase 2 involves the creation of a class
 85 `DefinedCollection` whose reward function is built to provide a full reward in cases where a
 86 pulled card is found and then removed in a living target list and partial rewards based on card
 87 values above a certain threshold. The number of pack openings will also be modified, based on the
 88 subtraction of a consistent pack price from a budget, allowing cards to contribute toward opening an
 89 additional pack.

90 We set the variable `cost_per_pack` to be sufficiently high that it would ensure that our MDP is
 91 finite (if `cost_per_pack` were lower than the expected value of any of the packs, then it would
 92 be theoretically possible for the budget to grow indefinitely without any of the target cards being
 93 collected, resulting in an infinite-horizon MDP). This was done so as not to preclude some of the
 94 dynamic programming-based methods that we learned in class (such as UCB-VI), though in the

95 end we didn't take advantage of dynamic programming (we discuss this further in the Conclusions
96 section).

97 While it is easy to write down the states and actions of this MDP, it is very high-dimensional, with the
98 state space consisting of all possible combinations of budget and drawn cards within the constraints
99 of the initial budget and target collection. The transition probabilities $P(s'|s, a)$, related to the pack
100 distributions, are also unknown at the beginning of the algorithm execution, so it is necessary to
101 explore the MDP to converge towards an optimal policy. In the next section, we describe how we
102 adapt the Bayesian method of Phase 1 to overcome these challenges.

103 4.2 Algorithm

104 The problem of maximizing probability to complete the specified collection initially led to us
105 considering Proximal Policy Evaluation in the Project Proposal. However, UCB-VI is being highly
106 considered due to how sparse the rewards provided from the packs may be, allowing UCB-VI's
107 exploration to shine in making sure card packs are not only chosen with confidence but are selected
108 without neglecting the potential of less explored packs.

109 In addition to the reward function described in Section 1.2, we also experimented with an exploration
110 bonus that would encourage exploration. Rather than base the exploration bonus on a Hoeffding bound
111 (as was done in the UCB-VI algorithm covered in class), we formulated a new bonus based on our
112 generalization of Thompson sampling from Phase 2 but still based on the principle of "optimization
113 in the face of uncertainty" that underlies the UCB algorithm. In particular, we wanted the policy
114 to give extra weight to packs that still may have a good chance to yield the card that is still most
115 needed at time t , $c'_t = \text{argmax}(\text{target})$, within the degree of confidence granted by the prior pack
116 openings. To do so, we maintained a set of Bayesian priors on the pack distributions that we updated
117 after each pack opening as in Phase 1. Then, we calculated the upper credible interval of $p(c'_t)$ for
118 each pack, weighted by a hyperparameter C that determines how greedy the policy should be. A
119 higher value of C encourages exploration; typically a value of between 10 and 20 was needed to
120 make the reward bonus comparable to the baseline rewards of the environment, as the large number
121 of cards in environment led to individual card draw probabilities that were quite low.

122 Mathematically, the marginal distribution over the Dirichlet priors for c'_t winds up being $p_{c'_t}^{(a)} \sim$
123 $\text{Beta}(\alpha', \beta')$ where in terms of the vector α parameterizing the Dirichlet distribution

$$\begin{aligned}\alpha' &= \alpha_{c'_t} \\ \beta' &= \sum_{c \neq c'_t} \alpha_c\end{aligned}$$

124 Then, the exploration bonus assigned to each pack a is given by $b_t(s, a) := C \cdot U_a$, where U_a is
125 calculated from

$$P(p_{c'_t}^{(a)} \leq U_a) := 1 - \delta \quad (1)$$

126 We set $\delta = 0.95$ and didn't use it as a separate hyperparameter.

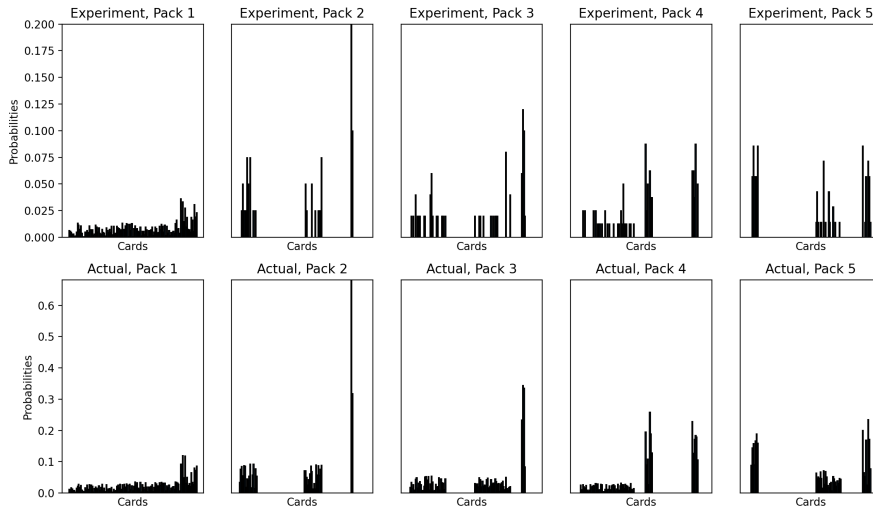
127 We also used the Bayesian idea to parameterize our policy. As an aside, one of the themes of the
128 course was the generality of MDPs, given that past history can be incorporated into the current state
129 to make the environment Markovian again. Incorporating all information about the past transitions
130 into a Bayesian prior is a particularly memory-efficient way of doing so. It also provides one
131 way of overcoming the large dimensionality of our state space S in designing policies $\pi(s)$, by
132 parameterizing the policies in terms of the priors α . We used a policy based on the expected reward
133 $r(s, a)$ of pulling pack a in state s , which represents all information about of the cards collected so
134 far as well as the current budget. Including the reward bonus discussed previously, our policy is

$$\pi_{\alpha, C}(s) := \arg \max_a \left[\mathbb{E}_{\hat{P}(\cdot|s, a) = \text{Dir}(\alpha_a)} [r(s, a)] + C \cdot b_t(s, a) \right] \quad (2)$$

135 The two terms in this policy complement each other – the first prioritizes greedy actions, while the
 136 second encourages exploration.

137 4.3 Results

138 We designed the card packs to guarantee that all cards existed in at least one pack. In particular, any
 139 card could be drawn from Pack 1 with low probability, while other packs had probabilities more
 140 concentrated around specific subsets of cards. In each case, we tracked all cards collected throughout
 141 the episode. Both of these points are illustrated in the following figure, where the empirical pack
 142 distributions after running the algorithm once are plotted above their actual distributions:



143 We measured the performance of our algorithm in several different ways. First, we measured the
 144 cumulative reward throughout the trajectory. We plotted this alongside the cumulative budget and
 145 fraction of `target` cards collected. Additionally, since we are interested in measuring how well the
 146 algorithm learns the MDP, we plotted a measure of the statistical distance between the Bayesian
 147 prior distributions based on the opened pack data and the "true" distributions of the packs. We
 148 initially considered using the KL-divergence (discussed in class in the context of PG methods) as
 149 a statistical distance measure, but eventually opted for a simpler metric that directly compares the
 150 prior parameters α to the true pack parameters p , as we expect that $\alpha \rightarrow p$ asymptotically as more
 151 samples are collected. Mathematically, our distance measure is written

$$d := |\alpha - p| \quad (3)$$

152 where $\bar{\alpha}_i := \alpha_i / |\alpha|$. We computed how d evolved for each pack separately throughout the trajectory.
 153 We tested our algorithm on two different `target` collections. In the first case, `target` consisted of
 154 one of each possible card (`completionist`). This test was somewhat easy, and the policy sampled
 155 more-or-less evenly from each pack even for a small value of C :

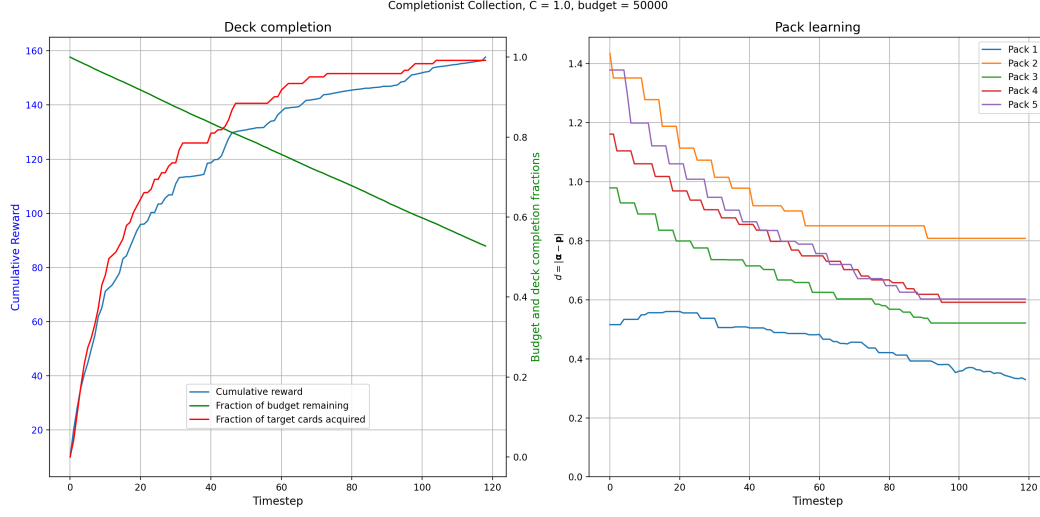


Figure 1: Performance of Bayesian policy on completionist goal with low exploration bonus, $C = 1.0$

156 Tuning C to encourage exploration over exploitation resulted in worse performance for this target
 157 collection:

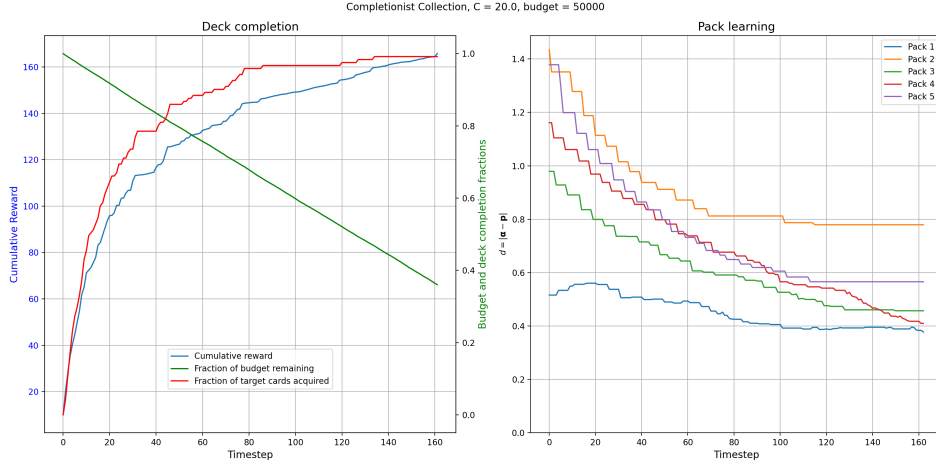


Figure 2: Performance of Bayesian policy on completionist goal with high exploration bonus, $C = 20.0$. Note the longer time required for deck completion and lower remaining budget upon completion

158 Next, we tested our policy on another target collection that consisted of many copies of a single
 159 rare card that can only be found in one deck (we called this test rare). Intuitively, one would expect
 160 forced exploration to be necessary for good performance, as it is not known in advance which pack
 161 the target card can be found in and without any exploration bonus b_t or priors on where to find the
 162 rare card, $\pi_{\alpha, C}(s)$ is determined only by the values of the cards in the packs, which are unrelated to
 163 the odds of finding the rare card. A test confirms this intuition:

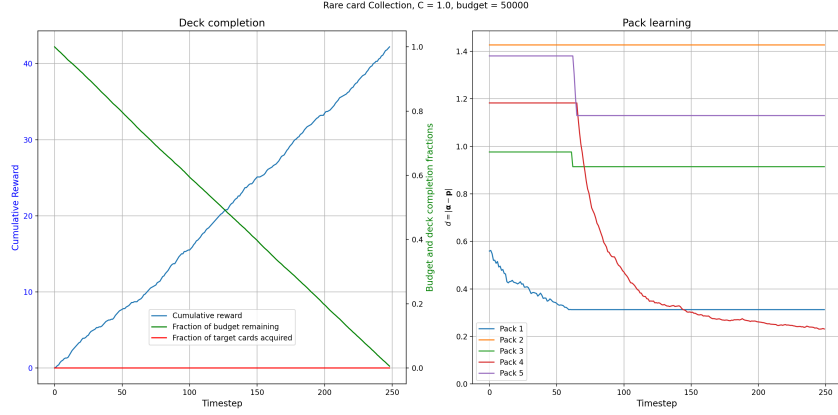


Figure 3: Performance of Bayesian policy on rare goal with low exploration bonus, $C = 1.0$. Note that the performance of this policy is highly-sensitive to whether it gets "lucky" and finds the rare card early by chance. This instance was run with `random.seed(1884)`.

- 164 In this instance, the policy never learns that the rare card is in Pack 2, and it never tries to pull that
 165 pack as its cards are of lower value on average and there is no reason to believe that it contains the
 166 rare card.
- 167 The situation is improved substantially when C is set to 20.0:

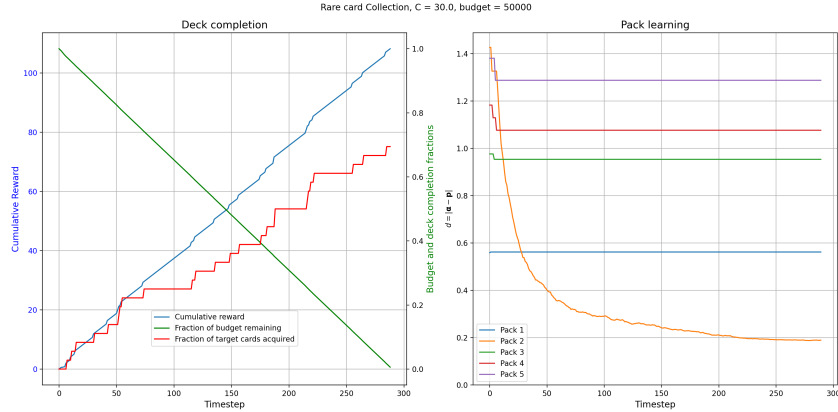


Figure 4: Performance of Bayesian policy on rare goal with high exploration bonus, $C = 1.0$

- 168 In this case, the policy pulls the rare card after about 10 timesteps and correctly focuses on Pack 2
 169 thereafter.

170 4.4 MDP Agent - Proximal Policy Optimization

- 171 Above, we borrow the idea of Exploration bonus from UCB-VI and ideas from Thompson sampling
 172 to solve the problem of completing the specific target collection while maximizing the budget, and
 173 in the process learn information about the card distributions. The reward defined above is based on
 174 a predefined target collection and includes a partial reward for high-value cards that can be sold to
 175 expand the budget. In this section, we'll take a different approach to the card collector, training it like
 176 a MDP agent, where it learns a general policy for optimizing pack openings and card buying/selling
 177 decisions, by maximizing a reward function that focus more heavily on completing a predefined
 178 collection. subsection 1.2

The PPOCollector class implements the Proximal Policy Optimization (PPO) algorithm as in the lecture notes. The key components are:

Policy Network $\pi_\theta(a|s)$: This neural network maps states to action probabilities. It is used to sample actions during trajectory collection and is updated via the PPO objective to learn the optimal policy.

- Inputs: state s
- Outputs: action probabilities $\pi_\theta(a|s)$ for each possible action a
- Update: Gradient ascent on PPO objective (Eq 6.50 in notes)

$$\mathbb{E}_{s, a \sim \pi_k} \left[\frac{\pi_\theta(a|s)}{\pi_k(a|s)} A^{\pi_k}(s, a) - \lambda \text{KL}(\pi_k || \pi_\theta) \right]$$

Value Network $V_\pi(s)$: This neural network estimates the expected return (value) from each state. It is used as a baseline in calculating advantages. Based on Lecture note's Baseline & Advantage.

- Inputs: state s
- Outputs: estimated value $V_\pi(s)$
- Update: Gradient descent on MSE loss between estimated and actual returns

4.5 Performance of PPO

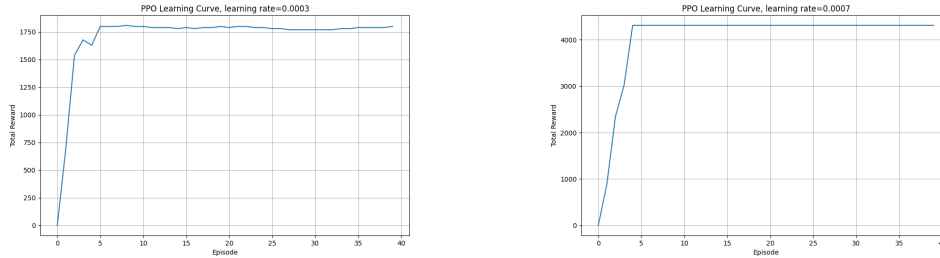


Figure 5: PPO with 2 different learning rates in its 2 Neural Network

The PPO agent's learning curve demonstrates a clear and effective learning progression across training episodes. Starting from near-zero rewards in the first 10 episodes during its initial exploration phase, the agent quickly experiences a dramatic improvement between episodes 10-20, where it rapidly discovers effective strategies for card collection, for both learning rates. This learning period shows some fluctuation especially learning rate 0.0003, particularly around episode 15 where there's a notable spike and dip, suggesting the agent found a promising strategy but needed time to stabilize its policy. After episode 20, the agent's performance converges to a steady reward level, around 1750 for learning rate 0.0003 and 4100 for learning rate 0.0007, maintaining this consistent performance through the remaining episodes up to 100. Learning rate 0.0003, after convergence the curve not as stable as a straight line like in learning rate 0.0007, which very likely indicates that it's stuck at a local minimum.

5 Conclusions

We implemented two conceptually different reinforcement learning approaches for a hobby card collecting problem. The first approach borrowed ideas from the Thompson sampling and UCB algorithms to learn the pack distributions and try to optimize the policy within a single epoch by aggregating new data about the MDP throughout the trajectory. The second approach used a policy gradient method (Proximal Policy Optimization) to find the optimal policy for assembling a target deck with less focus on cost and more focus on deck completion. Two methods don't use the same reward functions so there is no method of comparisons; however, each performed as well as expected.

210 **References**

- 211 [1] Brunskill, E., and Li, L. The Online Coupon-Collector Problem and Its Application to Lifelong Reinforcement
212 Learning. *arXiv preprint arXiv:1506.03379*, 2015. Available at <https://arxiv.org/abs/1506.03379>.