

---

# 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:**  $A = A_{\text{pack}} \cup A_{\text{trade}}$  where:

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

**Reward Function:**

$$R(S_t; a_t; S_{t+1}) = R_{\text{collection}} + R_{\text{budget}}$$

where:

- $R_{\text{collection}} = 10 \cdot \text{completion}$
- $R_{\text{budget}} = \begin{cases} 0 & \text{if } b_t > 0 \\ 0.1 \cdot b_t & \text{if } b_t < 0 \end{cases}$
- completion is the change in number of needed cards collected
- $b_t$  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  $\text{Dir}(\cdot_k)$  is a conjugate prior of  $\text{Mult}(K; p)$  and has a similar  
 65 Bayesian update rule to that of the Beta distribution, which is just to increment  $\frac{j}{k}$  by one each time  
 66 pack  $k$  outputs card  $j$ .

67 Finally, we define the reward of a pack pull  $n_t$  as the total value of all the cards drawn, i.e.  $r_t := n_t \cdot v$ ,  
 68 where  $v \in \mathbb{R}^J$  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  $T \gg K$

---

```

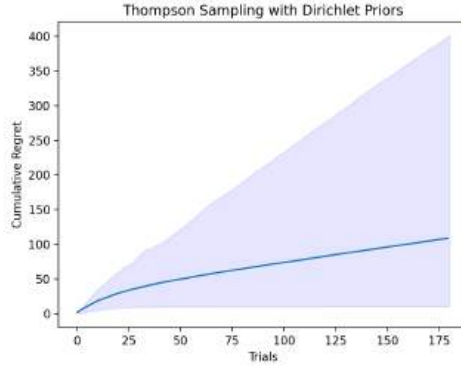
1: Initialize  $\hat{\theta}_k^{(j)} = 1 \forall k \in 1:K, j \in 1:J$ 
2: for  $t = 0; 1; \dots; T - 1$  do
3:   for  $k = 1; 2; \dots; K$  do
4:     Sample  $\hat{\theta}_k \sim \text{Dir}(\hat{\theta}_k)$ 
5:     Compute reward map  $\hat{r}_k = r(\hat{\theta}_k)$ , where  $r$  is the function that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

9: end for
```

---

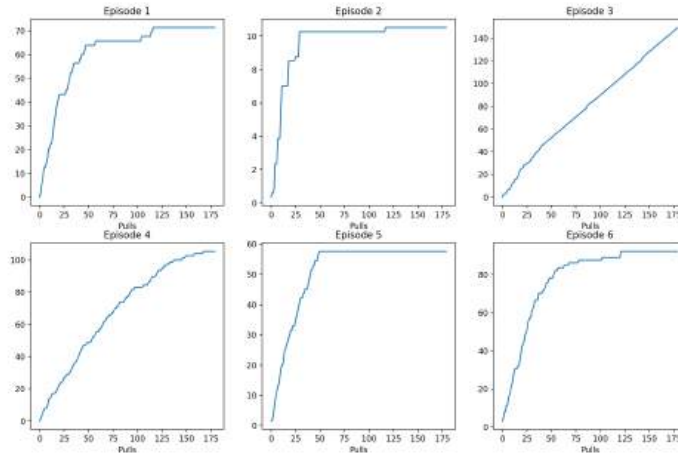
## 72 3.3 Data generation

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

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

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

## 4.2 Algorithm

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

In addition to the reward function described in Section 1.2, we also experimented with an exploration bonus that would encourage exploration. Rather than base the exploration bonus on a Hoeffding bound (as was done in the UCB-VI algorithm covered in class), we formulated a new bonus based on our generalization of Thompson sampling from Phase 2 but still based on the principle of "optimization in the face of uncertainty" that underlies the UCB algorithm. In particular, we wanted the policy to give extra weight to packs that still may have a good chance to yield the card that is still most needed at time  $t$ ,  $c_t^j = \text{argmax}(\text{target})$ , within the degree of confidence granted by the prior pack openings. To do so, we maintained a set of Bayesian priors on the pack distributions that we updated after each pack opening as in Phase 1. Then, we calculated the upper credible interval of  $p(c_t^j)$  for each pack, weighted by a hyperparameter  $C$  that determines how greedy the policy should be. A higher value of  $C$  encourages exploration; typically a value of between 10 and 20 was needed to make the reward bonus comparable to the baseline rewards of the environment, as the large number of cards in environment led to individual card draw probabilities that were quite low.

Mathematically, the marginal distribution over the Dirichlet priors for  $c_t^j$  winds up being  $p_{c_t^j}^{(a)}$  Beta( $\theta_j; \theta$ ) where in terms of the vector parameterizing the Dirichlet distribution

$$\theta = \sum_{c \in c_t} X_c^j$$

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

$$P(p_{c_t^j}^{(a)} \leq U_a) := 1 \quad (1)$$

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

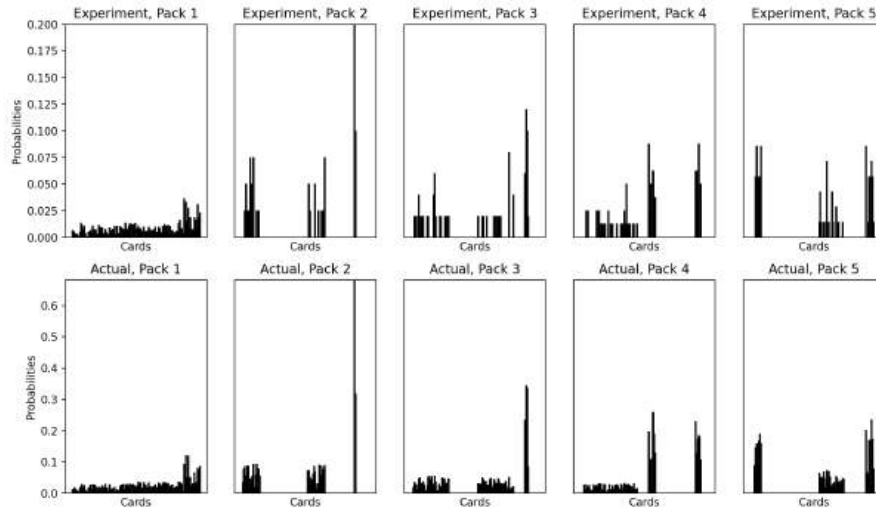
We also used the Bayesian idea to parameterize our policy. As an aside, one of the themes of the course was the generality of MDPs, given that past history can be incorporated into the current state to make the environment Markovian again. Incorporating all information about the past transitions into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

$$\pi(s) := \text{argmax}_a \sum_{j \in c_t} E_{p(j|s;a) = \text{Dir}(\theta_a)} [r(s; a)] + C \cdot b_t(s; a) \quad (2)$$

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 to the true pack parameters  $p$ , as we expect that  $\| \hat{p} - p \|$  asymptotically as more  
 151 samples are collected. Mathematically, our distance measure is written

$$d := \sum_i \|\hat{p}_i - p_i\| \quad (3)$$

152 where  $\hat{p}_i := \frac{1}{n} \sum_{j=1}^n \mathbb{I}(c_j = i)$ . 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  $C$ :









