

Adaptive Submodular Policy Optimization

Anonymous authors
Paper under double-blind review

Keywords: submodularity, adaptive submodularity, policy gradients

Summary

We propose KL-regularized policy optimization for adaptive submodular maximization. Adaptive submodularity is a framework for decision-making under uncertainty with submodular rewards. The benefit of policy optimization is that we can learn controllers for large action spaces that can utilize state-of-the-art large language model (LLM) priors. The benefit of submodularity are more efficient policy gradient updates because the gradient associated with an action only affects its immediate gain. When the reward model is correctly specified, we prove that our policies monotonically improve as the regularization diminishes and converge to the optimal greedy policy. Our experiments show major gains in statistical efficiency, in both synthetic problems and LLMs.

Contribution(s)

1. We propose KL-regularized policy optimization for adaptive submodular maximization.
Context: There are prior works on gradient-based optimization of submodular (not adaptive) functions. See Paragraph 2 in Section 6. There are prior works on policy gradients in more general settings. See Paragraphs 1 and 3 in Section 6.
2. We derive more efficient policy gradient estimators than in more general settings, with $O(n)$ terms as opposing to $O(n^2)$, where n is the horizon.
Context: None
3. We prove that our policy converges to the optimal greedy policy for adaptive submodular maximization as the regularization diminishes (Theorem 1).
Context: None
4. We prove that our policies monotonically improve over reference policies used for their regularization as the regularization diminishes (Theorem 4).
Context: None
5. We demonstrate the efficiency of new policy gradient estimators empirically, on both synthetic problems and LLMs (Section 5).
Context: None

Adaptive Submodular Policy Optimization

Anonymous authors

Paper under double-blind review

Abstract

1 We propose KL-regularized policy optimization for adaptive submodular maximization.
 2 Adaptive submodularity is a framework for decision-making under uncertainty with
 3 submodular rewards. The benefit of policy optimization is that we can learn controllers
 4 for large action spaces that can utilize state-of-the-art large language model (LLM) priors.
 5 The benefit of submodularity are more efficient policy gradient updates because the
 6 gradient associated with an action only affects its immediate gain. When the reward
 7 model is correctly specified, we prove that our policies monotonically improve as the
 8 regularization diminishes and converge to the optimal greedy policy. Our experiments
 9 show major gains in statistical efficiency, in both synthetic problems and LLMs.

10 1 Introduction

11 Many real-world problems have *diminishing returns*. The number of influenced people in a social
 12 network increases sublinearly with the number of influencers (Kempe et al., 2003). The information
 13 gain due to adding a sensor decreases if other sensors have already been placed at similar locations
 14 (Krause et al., 2008). Recommending similar content to already recommended content does not
 15 increase engagement (Yue & Guestrin, 2011; Hiranandani et al., 2019). The property of diminishing
 16 returns, known as *submodularity*, allows for efficient optimization. Specifically, a greedy algorithm
 17 for maximizing submodular functions in n steps is $(1 - 1/e)$ -optimal (Nemhauser et al., 1978).

18 We study adaptive decision making with submodular functions. *Adaptive submodularity* (Golovin &
 19 Krause, 2011) is a generalization of submodularity where the expected gain in reward after taking an
 20 action, in expectation over its observation, is a submodular function. One application of adaptive
 21 submodularity is preference elicitation (Gabilon et al., 2013), which is a special case of question-
 22 answering games (Dasgupta, 2005; Karbasi et al., 2012). These problems are submodular because
 23 the information gain due to asking a question diminishes with more previously asked questions. A
 24 greedy algorithm for adaptive submodular maximization in n steps, which takes the action with the
 25 highest expected gain conditioned on the history, is $(1 - 1/e)$ -optimal (Golovin & Krause, 2011).

26 The goal of this work is to bring together submodular and policy optimization, to their mutual benefit.
 27 In particular, *policy gradients* (Williams, 1992) arose as a versatile tool for reinforcement learning
 28 (Sutton & Barto, 1998) and are behind the recent advances in learning *large language models (LLMs)*
 29 (Schulman et al., 2015; 2017; Ouyang et al., 2022). The benefit of casting submodular maximization
 30 as policy learning is that we can learn controllers for large action spaces, of all responses of the LLM.
 31 The benefit of submodularity in optimization are more efficient policy gradient updates, because the
 32 gradient associated with an action only affects its immediate gain. This is in contrast to more general
 33 recent frameworks, such as submodular reinforcement learning (Prajapat et al., 2024).

34 We make the following contributions:

- 35 1. We propose KL-regularized policy optimization for adaptive submodularity (Section 3). The bene-
 36 fit of formulating adaptive submodular maximization in this way is that we can learn controllers
 37 for large action spaces that can leverage state-of-the-art pre-trained policies, such as LLMs. Our

- 38 main contribution to policy optimization are more efficient policy gradient updates, because the
 39 gradient associated with an action only affects its immediate gain.
- 40 2. We analyze our policies and prove two claims. First, we show that our policy converges to the
 41 optimal greedy policy for adaptive submodular maximization as the regularization diminishes.
 42 Second, we show that our policies monotonically improve over reference policies used for their
 43 regularization as the regularization diminishes. The main contribution in our analysis is bringing
 44 together techniques for analyzing KL-regularized policies and adaptive submodular maximization.
 45 This requires generalization of existing concepts of near-optimal adaptive submodular policies to
 46 stochastic policies, for instance.
- 47 3. We empirically evaluate our policies for adaptive submodular maximization. They can be learned
 48 more efficiently than using a vanilla policy gradient and are applicable to LLMs.

49 2 Background

50 We start with introducing our notation. Random variables are capitalized, except for Greek letters like
 51 θ . We denote the marginal and conditional probabilities under probability measure p by $p(X = x)$
 52 and $p(X = x \mid Y = y)$, respectively. When the random variables are clear from context, we write
 53 $p(x)$ and $p(x \mid y)$. For a positive integer n , we define $[n] = \{1, \dots, n\}$. The indicator function is
 54 $\mathbb{1}\{\cdot\}$. The i -th entry of vector v is v_i . If the vector is already indexed, such as v_j , we write $v_{j,i}$.

55 We introduce our multi-step optimization notation next. An *agent* interacts with the environment for n
 56 steps. To simplify exposition, we assume that n is fixed. The agent initially observes a *context* $x \in \mathcal{X}$,
 57 where \mathcal{X} is the space of contexts. The context is a side information that could define the problem
 58 instance, for example. In step $t \in [n]$, the agent takes an *action* a_t and *observes* y_t . The difference
 59 between actions and observations is that the agent controls the actions. The observations depend on
 60 actions but are provided by the environment. The *history* of n actions and their observations is a set
 61 $h_n = \{(a_t, y_t)\}_{t \in [n]}$. We denote by $r(x, h_n) \geq 0$ the *reward* associated with context x and history
 62 h_n . The probability that action a is taken in context x and history h_{t-1} is $\pi(a \mid x, h_{t-1}; \theta)$, and is
 63 parameterized by $\theta \in \Theta$. We call θ a *policy* and Θ the space of policy parameters. The action and
 64 observation in step t are generated as $a_t \sim \pi(\cdot \mid x, h_{t-1}; \theta)$ and $y_t \sim p(\cdot \mid x, h_{t-1}, a_t)$, respectively.
 65 Since the order of the observations in the history does not matter, our setting is less general than
 66 classic reinforcement learning (Sutton & Barto, 1998) but more general than a bandit (Lattimore &
 67 Szepesvari, 2019), because both a_t and y_t depend on the history.

68 The probability of history h_n in context x under policy θ factors as

$$\pi(h_n \mid x; \theta) = \prod_{t=1}^n p(y_t \mid x, h_{t-1}, a_t) \pi(a_t \mid x, h_{t-1}; \theta). \quad (1)$$

69 This follows from the chain rule and our modeling assumptions. The value of policy θ is

$$V(\theta) = \mathbb{E}_{x, h_n \sim \pi(\cdot \mid x; \theta)} [r(x, h_n)],$$

70 where $x \sim \mathcal{D}$ is drawn from a distribution of contexts \mathcal{D} . The optimal policy and its value are

$$\theta^* = \arg \max_{\theta \in \Theta} V(\theta), \quad V^* = \max_{\theta \in \Theta} V(\theta), \quad (2)$$

71 respectively. The question-answering game in Section 1 can be formulated in our notation as follows.
 72 The questions are actions, the answers are observations, and the reward is the fraction of objects that
 73 the user does not think about, based on the questions and their answers in the history.

74 2.1 Adaptive Submodularity

75 Adaptive submodularity (Golovin & Krause, 2011) is a framework for sequential decision making
 76 under uncertainty with diminishing returns. Under this assumption, a near-optimal policy is greedy
 77 conditioned on the history and thus can be computed efficiently.

78 *Adaptive submodularity* is formally defined as follows. Let

$$\Delta(a \mid x, h_{t-1}) = \mathbb{E}_{y \sim p(\cdot \mid x, h_{t-1}, a)} [r(x, h_{t-1} + \{(a, y)\})] - r(x, h_{t-1}) \quad (3)$$

79 be the *expected gain* in reward after taking action a in context x and history h_{t-1} . We make two
80 assumptions. First, the expected gain is *non-negative*; $\Delta(a \mid x, h_{t-1}) \geq 0$ holds for any context x ,
81 history h_{t-1} , and action a . Second, the expected gain is *submodular*,

$$\Delta(a \mid x, h_{t-1}) \geq \Delta(a \mid x, h_{t-1} + \{(a', y')\})$$

82 holds for any context x , history h_{t-1} , actions a and a' , and observation y' . These assumptions are
83 analogous to those in classic submodularity (Nemhauser et al., 1978), except that the ground set are
84 actions and the assumptions are in expectation over the observations of the actions. Similarly to the
85 classic setting, they imply efficiency. Specifically, let

$$\pi_g(a \mid x, h_{t-1}) = \mathbb{1} \left\{ a = \arg \max_{a'} \Delta(a' \mid x, h_{t-1}) \right\} \quad (4)$$

86 be the greedy policy with respect to $\Delta(a \mid x, h_{t-1})$. Then its expected value is at least $(1 - 1/e)V^*$
87 (Golovin & Krause, 2011), where V^* is defined in (2).

88 2.2 KL-Regularized Policy Optimization

89 One limitation of solving adaptive submodular problems as in (4) is that the maximization is difficult
90 when the action space is large or infinite, such as in LLMs (Brown et al., 2020; Wei et al., 2022). This
91 motivates our work on solving (4) as a controller learning problem. Learning of controllers for large
92 action spaces is at the center of *reinforcement learning from human feedback (RLHF)* (Christiano
93 et al., 2017). Specifically, once a reward model is learned, the policy is optimized to maximize the
94 expected reward under the reward model using *proximal policy optimization (PPO)* (Schulman et al.,
95 2017). Specifically, the objective is

$$\mathcal{L}_{\text{PPO}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, a \sim \pi(\cdot \mid x; \theta)} \left[r(x, a) - \beta \log \frac{\pi(a \mid x; \theta)}{\pi_0(a \mid x)} \right], \quad (5)$$

96 where x is a prompt sampled from a dataset of prompts \mathcal{D} , a is its response, and $\pi(a \mid x; \theta)$ is the
97 probability of generating response a to prompt x by policy θ . The first term is the expected reward
98 for response a to prompt x . The second term penalizes for deviations of the optimized policy from a
99 *reference policy* π_0 , usually obtained by supervised fine-tuning (Mangrulkar et al., 2022; Hu et al.,
100 2022). The parameter $\beta \geq 0$ trades off the two terms. In adaptive submodularity (Section 2.1), the
101 prompt x and its response a are the history and action, respectively.

102 PPO is a popular policy-learning framework with two benefits. First, it is suitable for large action
103 spaces. Specifically, once the policy is learned, the best action is just sampled from it. Second, the
104 prior information can be integrated through the reference policy. While PPO has been popularized by
105 RLHF, we note that the idea of KL-regularized policies goes to Schulman et al. (2015), where it was
106 used to motivate trust-region policy optimization; and to Todorov (2006), where it was proposed and
107 analyzed in the context of Markov decision processes (Puterman, 1994).

108 3 Algorithm

109 We bring together adaptive submodular maximization and KL-regularized policy optimization. This
110 has two benefits. First, we extend adaptive submodular maximization to large action spaces and
111 learning from pre-trained reference policies. Second, KL-regularized policy optimization can be done
112 more efficiently by leveraging adaptive submodularity.

Algorithm 1 KL-PO

```

1: Input: Learning rate schedule  $(\alpha_i)_{i \in \mathbb{N}}$ 
2: Initialize  $\theta$  and  $i \leftarrow 1$ 
3: while not convergence do
4:   Simulate  $h_n \sim \pi(\cdot | x; \theta)$ 
5:    $\theta \leftarrow \theta + \alpha_i \sum_{t=1}^n (f_t(\theta) - \beta) \sum_{\ell=1}^t \nabla \log \pi(a_\ell | x, h_{\ell-1}; \theta)$ 
6:    $i \leftarrow i + 1$ 
7: Output: Learned policy  $\theta$ 

```

Algorithm 2 KL-SubPO

```

1: Input: Learning rate schedule  $(\alpha_i)_{i \in \mathbb{N}}$ 
2: Initialize  $\theta$  and  $i \leftarrow 1$ 
3: while not convergence do
4:   Simulate  $h_n \sim \pi(\cdot | x; \theta)$ 
5:    $\theta \leftarrow \theta + \alpha_i \sum_{t=1}^n (f_t(\theta) - \beta) \times \nabla \log \pi(a_t | x, h_{t-1}; \theta)$ 
6:    $i \leftarrow i + 1$ 
7: Output: Learned policy  $\theta$ 

```

113 **3.1 Classic Policy Optimization**

114 To understand the benefit of our method, we first introduce a classic n -step KL-regularized policy
 115 optimization. When actions in (5) are replaced with histories, we immediately obtain

$$\mathcal{L}_{\text{KL-PO}}(\theta, \beta) = \mathbb{E}_\theta \left[r(x, h_n) - \beta \sum_{t=1}^n \log \frac{\pi(h_n | x; \theta)}{\pi_0(h_n | x)} \right],$$

116 where $\mathbb{E}_\theta[\cdot] = \mathbb{E}_{x \sim \mathcal{D}, h_n \sim \pi(\cdot | x; \theta)}[\cdot]$. The problem of policy optimization is to maximize $\mathcal{L}_{\text{KL-PO}}(\theta, \beta)$
 117 with respect to θ . We call this algorithm **KL-PO** and present it in Algorithm 1. The main challenge is
 118 that the gradient of $\mathcal{L}_{\text{KL-PO}}(\theta, \beta)$ has $O(n^2)$ terms. To see this, we first note that

$$\mathbb{E}_\theta[r(x, h_n)] = \sum_{t=1}^n \mathbb{E}_{\theta, t}[\Delta(a_t | x, h_{t-1})],$$

119 where $\mathbb{E}_{\theta, t}[\cdot] = \mathbb{E}_{x \sim \mathcal{D}, h_{t-1} \sim \pi(\cdot | x; \theta), a_t \sim \pi(\cdot | x, h_{t-1}; \theta)}[\cdot]$. This follows from the factorization of $\pi(h_n |$
 120 $x; \theta)$ in (1) and the definition of $\Delta(a_t | x, h_{t-1})$ in (3). Therefore, the n -step objective is

$$\mathcal{L}_{\text{KL-PO}}(\theta, \beta) = \sum_{t=1}^n \mathbb{E}_{\theta, t}[f_t(\theta)], \quad (6)$$

121 where

$$f_t(\theta) = \Delta(a_t | x, h_{t-1}) - \beta \log \frac{\pi(a_t | x, h_{t-1}; \theta)}{\pi_0(a_t | x, h_{t-1})}.$$

122 Using basic rules of differentiation and the score identity (Aleksandrov et al., 1968), we obtain

$$\nabla \mathbb{E}_{\theta, t}[f_t(\theta)] = \mathbb{E}_{\theta, t} \left[(f_t(\theta) - \beta) \sum_{\ell=1}^t \nabla \log \pi(a_\ell | x, h_{\ell-1}; \theta) \right]. \quad (7)$$

123 Therefore, the policy gradient (Williams, 1992) of (6) involves $n(n+1)/2$ terms. This leads to an
 124 $O(n^2)$ variance in the empirical estimate in **KL-PO** (line 5). The dependence on prior actions arises
 125 because they all impact the gain in step t . This motivated many prior works on variance reduction of
 126 policy gradients (Sutton et al., 2000; Baxter et al., 2001; Baxter & Bartlett, 2001; Munos, 2006).

127 **3.2 Adaptive Submodular Policy Optimization**

128 The key idea in our algorithm is to replace the empirical gradient estimate in **KL-PO** (line 5), which
 129 involves $\sum_{\ell=1}^t \nabla \log \pi(a_\ell | x, h_{\ell-1}; \theta)$, with $\nabla \log \pi(a_t | x, h_{t-1}; \theta)$. An informal justification for
 130 this choice is that for any content x and history h_{t-1} , a near-optimal policy in (4) only maximizes the
 131 immediate gain conditioned on x and h_{t-1} .

Mathematically, this change can be viewed as follows. Suppose that (6) is replaced with

$$\mathcal{L}_{\text{KL-SUBPO}}(\theta, \beta) = \sum_{t=1}^n \mathbb{E}_{\theta, \theta_h, t} [f_t(\theta)] , \quad (8)$$

where $\mathbb{E}_{\theta, \theta_h, t} [\cdot] = \mathbb{E}_{x \sim \mathcal{D}, h_{t-1} \sim \pi(\cdot | x; \theta_h), a_t \sim \pi(\cdot | x, h_{t-1}; \theta)} [\cdot]$ and θ_h is a history-generating policy that is independent of θ . Then, using basic rules of differentiation and the score identity (Aleksandrov et al., 1968), we obtain

$$\nabla \mathbb{E}_{\theta, \theta_h, t} [f_t(\theta)] = \mathbb{E}_{\theta, \theta_h, t} [(f_t(\theta) - \beta) \nabla \log \pi(a_t | x, h_{t-1}; \theta)] . \quad (9)$$

This gradient differs from (7) because we do not differentiate with respect to θ_h . The result is a major gain in efficiency, due to replacing t terms in $\nabla \mathbb{E}_{\theta, t} [f_t(\theta)]$ by a single one.

We call the resulting algorithm **KL-SubPO** and present it in Algorithm 2. Although (9) has fewer terms than (7), the objective (8) needs to be properly justified and we do that in Section 4. Specifically, we prove that when the problem is adaptive submodular, the maximization of (8) yields near-optimal greedy policies for any history-generating policy θ_h . The learned policies monotonically improve over reference policies π_0 as $\beta \rightarrow 0$ when the reward model is correctly specified.

4 Analysis

We make the following assumptions. First, we analyze an idealized variant of **KL-SubPO**, which is formulated as a maximization of (8). Second, we assume that the optimal solution to (8) is realizable and identifiable. Finally, we assume that the reward model is known.

We start with the observation that

$$\mathbb{E}_{\theta, \theta_h, t} [f_t(\theta)] = \mathbb{E}_{x \sim \mathcal{D}, h_{t-1} \sim \pi(\cdot | x; \theta_h)} [\mathbb{E}_{a_t \sim \pi(\cdot | x, h_{t-1}; \theta)} [f_t(\theta) | x, h_{t-1}]] ,$$

The inner expectation has the same algebraic form as (5). Thus, for any context x and history h_{t-1} , the maximizer has a closed form (Todorov, 2006) of

$$\pi(a | x, h_{t-1}; \theta) = \frac{1}{Z(x, h_{t-1})} \pi_0(a | x, h_{t-1}) \exp \left[\frac{1}{\beta} \Delta(a | x, h_{t-1}) \right] , \quad (10)$$

where $Z(x, h_{t-1})$ is the normalizer. This allows us to analyze the properties of the optimal policy irrespective of θ_h . In the following, we first show that as $\beta \rightarrow 0$, the policy converges to the optimal greedy policy. Then we introduce γ -approximate policies to analyze the non-asymptotic behavior of **KL-SubPO**.

Theorem 1. *Let $\hat{\theta}(\beta) = \arg \max_{\theta} \mathcal{L}_{\text{KL-SUBPO}}(\theta, \beta)$. Let $\Delta(a | x, h_{t-1})$ be the expected gain of taking action a in context x and history h_{t-1} , as defined in (3). Let π_g be the near-optimal greedy policy in (4). Then, if the best greedy action is unique, for any x, h_{t-1} , and a ,*

$$\lim_{\beta \rightarrow 0} \pi(a | x, h_{t-1}; \hat{\theta}(\beta)) = \pi_g(a | x, h_{t-1}) .$$

Proof Sketch. When $\beta = 0$, the KL regularizer in (8) vanishes and our policy ends up maximizing $\Delta(a | x, h_{t-1})$, which is exactly the greedy policy in (4). See Appendix A for details. \square

This result confirms that as the KL regularization diminishes, our policy becomes the greedy policy that maximizes the expected marginal gain. Now we analyze the non-asymptotic behavior through the novel concept of γ -approximate greedy policies.

4.1 γ -Approximate Greedy Policies

Traditional greedy policies take actions that maximize the expected marginal gain. The solutions to (8) do that only approximately. Therefore, we extend the notion of the marginal gain from individual actions to entire policies. For a policy θ , the expected marginal gain is

$$\Delta(\theta | x, h_{t-1}) = \mathbb{E}_{a \sim \pi(\cdot | x, h_{t-1}; \theta)} [\Delta(a | x, h_{t-1})] . \quad (11)$$

166 **Definition 2** (γ -Approximate Greedy Policy). For $\gamma \geq 1$, a policy θ is γ -approximate greedy if for
 167 all contexts x and histories h_{t-1} ,

$$\Delta(\theta \mid x, h_{t-1}) \geq \frac{1}{\gamma} \max_{a'} \Delta(a' \mid x, h_{t-1}). \quad (12)$$

168 Our notion of γ -approximate greedy policies is inspired by but distinct from the approximate greedy
 169 policies in Golovin & Krause (2011). While they require every action in the policy’s support to be
 170 approximately optimal, we only require the approximate optimality of the *expected gain* with respect
 171 to a fixed policy. This relaxation is better suited for our setting because we learn stochastic policies,
 172 which have large action spaces and are regularized by pre-trained LLM priors using KL.

173 **Theorem 3** (Performance of γ -Approximate Greedy Policies). Let θ be a γ -approximate greedy policy
 174 and V^* be the expected value of the optimal n -step policy. Under the assumptions in Section 2.1,

$$V^* - V(\theta) \leq (1 - 1/e^{1/\gamma})V^*. \quad (13)$$

175 *Proof Sketch.* The proof follows a standard submodularity argument. We define the optimality gap
 176 $X_t = V^* - r(x, h_t)$ and show that $\mathbb{E}[X_t]$ decreases exponentially at rate $1/(\gamma n)$. See Appendix A
 177 for details. \square

178 This result generalizes the classic $(1 - 1/e)$ -approximation guarantee to approximate greedy policies,
 179 with the approximation factor that degrades smoothly with γ . When $\gamma = 1$, we recover the classic
 180 guarantee of the exact greedy policy.

181 4.2 Improvement Guarantees

182 Having characterized the performance of γ -approximate greedy policies generally, we now establish
 183 how KL-SubPO produces improved policies.

184 **Theorem 4** (Policy Improvement). Let the reference policy π_0 in (8) be γ -approximate greedy. Let
 185 $\hat{\theta}(\beta) = \arg \max_{\theta} \mathcal{L}_{\text{KL-SUBPO}}(\theta, \beta)$ be the optimal solution and $\hat{\pi}(\cdot \mid \cdot) = \pi(\cdot \mid \cdot; \hat{\theta}(\beta))$. Then there
 186 exists $\gamma' \in [1, \gamma]$ such that

$$V^* - r(\hat{\pi} \mid x, h_{t-1}) \leq \left(1 - \frac{1}{\gamma' n}\right) (V^* - r(x, h_{t-1})) \quad (14)$$

187 where $r(\hat{\pi} \mid x, h_{t-1}) = \mathbb{E}_{a \sim \hat{\pi}}[r(x, h_{t-1} \cup \{(a, y)\})]$, holds for all contexts x and histories h_{t-1} .
 188 Furthermore:

- 189 1. $\hat{\pi}$ is a $(1 - 1/e^{1/\gamma'})$ -optimal policy.
- 190 2. γ' decreases monotonically with the regularization parameter β .

191 This theorem establishes two important properties of our KL-SubPO policies. First, they improve a
 192 γ -approximate greedy reference policy to a policy with an approximation factor $\gamma' \leq \gamma$. Second,
 193 the regularization parameter β affects this improvement: a stronger regularization (larger β) leads to
 194 more conservative improvements, while a weaker regularization makes the policy more greedy. The
 195 core insight behind this result is the closed-form solution in (10), which indicates monotonicity. We
 196 formalize and prove it properly.

197 5 Experiments

198 We conduct three experiments. The first two experiments are synthetic and the last one is on LLMs.

199 5.1 Linear Maximization

200 In the first experiment, we study n -step maximization of a linear function with K unknown parameters.
 201 The function is represented by a vector $w \in \mathbb{R}^K$ where $w_k = (k/K)^2$. The actions are the standard

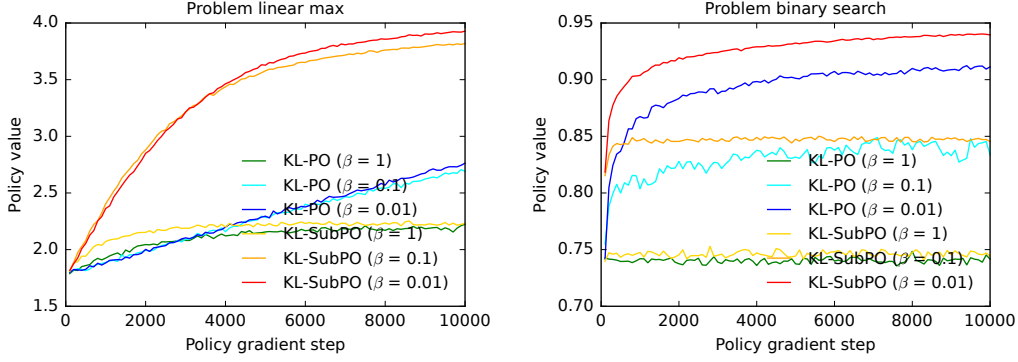


Figure 1: Experiments on the linear maximization problem in Section 5.1 and the binary search problem in Section 5.2.

202 basis in \mathbb{R}^K , $\mathcal{A} = \{e_i\}_{i=1}^K$. The non-zero entry of an action indicates the revealed entry of w . The
 203 reward is the sum of the revealed entries $r(x, h_t) = \sum_{\ell=1}^t a_\ell^\top w$. The policy is parameterized as
 204 $\pi(a \mid x, h_t; \theta) \propto \exp[\phi(h_t, a)^\top \theta]$, where $\phi(h_t, a)$ is the feature vector for history h_t and action
 205 a . The feature vector for action e_i is a zero vector if the action was taken before and e_i otherwise.
 206 Formally, for any $e_i \in \mathcal{A}$ and $k \in [K]$, $\phi_k(h_t, e_i) = e_{i,k} \prod_{\ell=1}^t (1 - a_{\ell,k})$. We set $K = 20$ and
 207 the horizon is $n = 5$. The optimal policy selects the 5 highest entries of w and its value is 4.07.
 208 We experiment with $\beta \in \{0.01, 0.1, 1.0\}$ to show a range of operating modes of **KL-SubPO**. The
 209 reference policy π_0 is uniform. All policies are optimized by Adam (Kingma & Ba, 2015).

210 Our results are reported in Figure 1a. We observe three main trends. First, **KL-SubPO** outperforms
 211 **KL-PO** for all β . This shows that **KL-SubPO** is generally more efficient than **KL-PO**. Second, **KL-SubPO**
 212 policies improve as $\beta \rightarrow 0$ when the reward model is correctly specified (Section 4). Finally, the
 213 **KL-SubPO** policy at $\beta = 0.01$ is near optimal.

214 5.2 Binary Search

215 In the second experiment, we have a binary search problem over $[K]$. A random integer k_* is chosen
 216 from $[K]$ and our goal is to identify it. The actions are all possible halving questions on $[K]$. More
 217 specifically, $\mathcal{A} = \{q_i\}_{i=1}^{K-1}$, where $q_i \in \{0, 1\}^K$ is a vector whose first i entries are ones and the rest
 218 are zeros. When the agent takes action q_i in step t , the observation is $y_t = q_{i,k_*}$. Simply put, the
 219 answer is “yes” if $k_* \leq i$ and “no” otherwise. The reward is the fraction of eliminated integers in
 220 $[K]$, that cannot be k_* based on the answers thus far,

$$r(x, h_t) = \frac{1}{K} \sum_{k=1}^K \prod_{\ell=1}^t y_\ell (1 - a_{\ell,k}) + (1 - y_t) a_{t,k}.$$

221 The policy is parameterized as in Section 5.1. The feature vector for action q_i is an outer product of
 222 the state s_t , which indicates the remaining integers, and q_i , $\phi(h_t, q_i) = \text{vec}(s_t^\top q_i)$. The state is

$$s_{t,k} = \mathbb{1} \left\{ \sum_{\ell=1}^t y_\ell a_{\ell,k} + (1 - y_t)(1 - a_{t,k}) = t \right\}.$$

223 We set $K = 32$ and the horizon is $n = 5$. The optimal policy is binary search and its value is 0.97.
 224 We experiment with the same policies as in Section 5.1.

225 Our results are reported in Figure 1b. We observe three main trends. First, **KL-SubPO** performs
 226 comparably to **KL-PO** when β is high and both policies perform poorly. Second, **KL-SubPO** policies
 227 improve as $\beta \rightarrow 0$ when the reward model is correctly specified (Section 4). Finally, the **KL-SubPO**
 228 policy at $\beta = 0.01$ is near optimal.

5.3 Twenty Questions

The last experiment is a 20Q game (Karbasi et al., 2012) with 20 animals. The agent is an LLM. It is optimized against a user represented by an LLM. The reward is the fraction of eliminated animals. The horizon is $n = 6$ questions. The experimental setup is described in detail in Appendix B. We conduct another experiment, where the animals are replaced with Amazon products, in Appendix C. We let the agent interact with the user and generate a dataset of 200 trajectories of length $n = 6$. The reward of the original LLM is 0.817 ± 0.006 . We standardize trajectory rewards to zero mean and unit variance, and learn a policy by KL-P0. Its reward is 0.815 ± 0.006 and the policy does not improve over the baseline. When the trajectory rewards are clipped at 0, the reward is 0.833 ± 0.005 (2% improvement over the baseline). We also standardize per-step gains to zero mean and unit variance, and learn a policy by KL-SubP0. Its reward is 0.829 ± 0.006 (1.5% improvement over the baseline). When the per-step gains are clipped at 0, the reward is 0.876 ± 0.004 (7% improvement over the baseline). We conclude that KL-SubP0 outperforms KL-P0 in both settings, irrespective of the rewards being clipped or not.

6 Related Work

Our work can be viewed as a special case of submodular reinforcement learning (Prajapat et al., 2024). This is because adaptive submodular functions are set functions, as opposing to functions of sequences of states and actions in Prajapat et al. (2024). These additional properties allow us to derive policy gradients that do not have a quadratic number of terms in the horizon n , unlike in Prajapat et al. (2024). The limitations of adaptive submodularity have been noted before and therefore it was extended, for instance to functions of sequences (Mitrovic et al., 2019).

Gradient-based optimization of submodular functions has also been explored before. For instance, Hassani et al. (2017) showed that stochastic projected gradient methods can provide strong approximation guarantees for maximizing continuous submodular functions with convex constraints. Bai et al. (2018) optimized deep submodular functions by gradient ascent. Our paper is the first work on gradient-based optimization of adaptive submodular functions.

Policy gradients were proposed by Williams (1992) and build on the score identity of Aleksandrov et al. (1968). It is well known that policy gradients have a high variance and therefore many variance reduction techniques have been proposed (Sutton et al., 2000; Baxter et al., 2001; Baxter & Bartlett, 2001; Munos, 2006; Kveton et al., 2020). Our contribution to these works is a policy gradient that does not have a quadratic number of terms in the horizon n .

7 Conclusions

We propose KL-regularized policy optimization for adaptive submodular maximization, a framework for decision-making under uncertainty with submodular rewards. The submodularity allows for more efficient policy gradients than in more general settings. The KL-regularization allows for learning policies for large or infinite action spaces that utilize state-of-the-art LLM priors.

Our analysis makes several simplifying assumptions, which allow us to study the problem more cleanly. First, we analyze an idealized variant of KL-SubP0, which is formulated as a maximization of (8). Second, we assume that the optimal solution to (8) is realizable and identifiable. Finally, we assume that the reward model is known. This is rarely the case in practice and the model has to be estimated. We will address these limitations in our future work.

References

V. M. Aleksandrov, V. I. Sysoyev, and V. V. Shemeneva. Stochastic optimization. *Engineering Cybernetics*, 5:11–16, 1968.

- 273 Wenruo Bai, William Stafford Noble, and Jeff Bilmes. Submodular maximization via gradient ascent:
274 The case of deep submodular functions. In *Advances in Neural Information Processing Systems*
275 *31*, 2018.
- 276 Jonathan Baxter and Peter Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial*
277 *Intelligence Research*, 15:319–350, 2001.
- 278 Jonathan Baxter, Peter Bartlett, and Lex Weaver. Experiments with infinite-horizon, policy-gradient
279 estimation. *Journal of Artificial Intelligence Research*, 15:351–381, 2001.
- 280 Tom Brown et al. Language models are few-shot learners. In *Advances in Neural Information*
281 *Processing Systems 33*, 2020.
- 282 Paul Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
283 reinforcement learning from human preferences. In *Advances in Neural Information Processing*
284 *Systems 30*, 2017.
- 285 Sanjoy Dasgupta. Analysis of a greedy active learning strategy. In *Advances in Neural Information*
286 *Processing Systems 17*, pp. 337–344, 2005.
- 287 Victor Gabillon, Branislav Kveton, Zheng Wen, Brian Eriksson, and S. Muthukrishnan. Adaptive
288 submodular maximization in bandit setting. In *Advances in Neural Information Processing Systems*
289 *26*, pp. 2697–2705, 2013.
- 290 Daniel Golovin and Andreas Krause. Adaptive submodularity: Theory and applications in active
291 learning and stochastic optimization. *Journal of Artificial Intelligence Research*, 42:427–486,
292 2011.
- 293 Hamed Hassani, Mahdi Soltanolkotabi, and Amin Karbasi. Gradient methods for submodular
294 maximization. In *Advances in Neural Information Processing Systems 30*, 2017.
- 295 Gaurush Hiranandani, Harvineet Singh, Prakhar Gupta, Iftikhar Ahamath Burhanuddin, Zheng Wen,
296 and Branislav Kveton. Cascading linear submodular bandits: Accounting for position bias and
297 diversity in online learning to rank. In *Proceedings of the 35th Conference on Uncertainty in*
298 *Artificial Intelligence*, 2019.
- 299 Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and
300 Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *Proceedings of the 10th*
301 *International Conference on Learning Representations*, 2022.
- 302 Amin Karbasi, Stratis Ioannidis, and Laurent Massoulié. Hot or not: Interactive content search using
303 comparisons. In *2012 Information Theory and Applications Workshop*, pp. 291–297, 2012.
- 304 David Kempe, Jon Kleinberg, and Eva Tardos. Maximizing the spread of influence through a social
305 network. In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge*
306 *Discovery and Data Mining*, pp. 137–146, 2003.
- 307 Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *Proceedings of*
308 *the 3rd International Conference on Learning Representations*, 2015.
- 309 Andreas Krause, Ajit Paul Singh, and Carlos Guestrin. Near-optimal sensor placements in Gaussian
310 processes: Theory, efficient algorithms and empirical studies. *Journal of Machine Learning*
311 *Research*, 9:235–284, 2008.
- 312 Branislav Kveton, Martin Mladenov, Chih-Wei Hsu, Manzil Zaheer, Csaba Szepesvari, and Craig
313 Boutilier. Differentiable meta-learning in contextual bandits. *CoRR*, abs/2006.05094, 2020. URL
314 <http://arxiv.org/abs/2006.05094>.
- 315 Tor Lattimore and Csaba Szepesvari. *Bandit Algorithms*. Cambridge University Press, 2019.

- 316 Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin
317 Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. [https://github.
318 com/huggingface/peft](https://github.com/huggingface/peft), 2022.
- 319 Marko Mitrovic, Ehsan Kazemi, Moran Feldman, Andreas Krause, and Amin Karbasi. Adaptive
320 sequence submodularity. In *Advances in Neural Information Processing Systems 32*, 2019.
- 321 Remi Munos. Geometric variance reduction in Markov chains: Application to value function and
322 gradient estimation. *Journal of Machine Learning Research*, 7:413–427, 2006.
- 323 G. L. Nemhauser, L. A. Wolsey, and M. L. Fisher. An analysis of approximations for maximizing
324 submodular set functions - I. *Mathematical Programming*, 14(1):265–294, 1978.
- 325 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
326 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
327 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and
328 Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances
329 in Neural Information Processing Systems 35*, 2022.
- 330 Manish Prajapat, Mojmir Mutny, Melanie Zeilinger, and Andreas Krause. Submodular reinforcement
331 learning. In *Proceedings of the 12th International Conference on Learning Representations*, 2024.
- 332 Martin Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
333 Wiley & Sons, New York, NY, 1994.
- 334 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
335 policy optimization. In *Proceedings of the 32nd International Conference on Machine Learning*,
336 pp. 1889–1897, 2015.
- 337 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
338 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL [https://arxiv.org/abs/
339 1707.06347](https://arxiv.org/abs/1707.06347).
- 340 Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge,
341 MA, 1998.
- 342 Richard Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods
343 for reinforcement learning with function approximation. In *Advances in Neural Information
344 Processing Systems 12*, pp. 1057–1063, 2000.
- 345 Emanuel Todorov. Linearly-solvable Markov decision problems. In *Advances in Neural Information
346 Processing Systems 19*, 2006.
- 347 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
348 Andrew Dai, and Quoc Le. Finetuned language models are zero-shot learners. In *Proceedings of
349 the 10th International Conference on Learning Representations*, 2022.
- 350 Ronald Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
351 learning. *Machine Learning*, 8(3-4):229–256, 1992.
- 352 Yisong Yue and Carlos Guestrin. Linear submodular bandits and their application to diversified
353 retrieval. In *Advances in Neural Information Processing Systems 24*, pp. 2483–2491, 2011.

354 A Proofs

355 *Proof of Theorem 1.* This is trivial. When $\beta = 0$, there is the KL-term vanishes from $\mathcal{L}_{\text{KL-SubP0}}$. So
 356 the optimal policy is the one that maximizes $\Delta(a|x, h_{t-1})$ at every x, h_{t-1} . This is exactly what
 357 greedy policy does. \square

358 **Lemma 5** (Value Upper Bound). *Let $\pi(\cdot | x, h_{t-1}, \theta)$ be a γ -approximate greedy policy and V^* be*
 359 *the expected reward of the optimal n -step policy. Then for all contexts x and histories h_{t-1} :*

$$V^* \leq r(x, h_{t-1}) + \gamma n \Delta(\theta | x, h_{t-1}),$$

360 *Proof.* The proof is based on the the usual submodular "each step can't help more than the first step"
 361 argument. Let π^* be an optimal n steps policy. Then

$$\begin{aligned} V^* - r(x, h_{t-1}) &\leq \mathbb{E}_{h_n \sim \pi^*} [r(x, h_{t-1} + h_n)] - r(x, h_{t-1}) \\ &= \sum_{k=1}^n \mathbb{E} [\Delta(a_k^* | x, h_{k-1} + h_{t-1})] \end{aligned}$$

362 where h_{k-1} is the history after $k-1$ steps under π^* and $a_k^* \sim \pi^*(\cdot | x, h_{k-1})$. By *adaptive*
 363 *submodularity*, each incremental gain satisfies

$$\begin{aligned} \Delta(a_k^* | x, h_{k-1} + h_{t-1}) &\leq \Delta(a_k^* | x, h_{t-1}) \\ &\leq \max_{a'} \Delta(a' | x, h_{t-1}) \\ &\leq \gamma \Delta(\theta | x, h_{t-1}). \end{aligned}$$

364 Summing over n steps gives

$$V^* - r(x, h_{t-1}) \leq \gamma n \Delta(\theta | x, h_{t-1}).$$

365 \square

366 **Lemma 6** (One-step Gap Reduction). *Under adaptive submodularity and for any γ -approximate*
 367 *greedy policy π , the expected reduction in the optimality gap after one step satisfies:*

$$\mathbb{E}[X_t] \leq (1 - 1/(\gamma n)) \mathbb{E}[X_{t-1}]. \quad (15)$$

368 *Proof.* For any realized history h_t , and any policy π we define the expected one-step reward as:

$$r(\pi | x, h_t) := r(x, h_t) + \mathbb{E}_{a \sim \pi(\cdot | x, h_t; \theta)} [\Delta(a | x, h_t)] \quad (16)$$

$$= \mathbb{E}_{a \sim \pi(\cdot | x, h_t; \theta)} [r(x, h_t \cup \{(a, y)\})] \quad (17)$$

369 where the second equality follows from the definition of $\Delta(a | x, h_t)$ in (3). By Lemma 5 adaptive
 370 submodularity implies:

$$V^* \leq r(x, h_{t-1}) + \gamma n \Delta(\pi | x, h_{t-1}) \quad (18)$$

371 This inequality captures the key property that the remaining value after history h_{t-1} is bounded by
 372 γn times the one-step gain.

373 Expanding using the definition of $r(\pi | x, h_{t-1})$:

$$V^* \leq r(x, h_{t-1}) + \gamma n \Delta(\pi | x, h_{t-1}) \quad (19)$$

$$= r(x, h_{t-1}) \quad (20)$$

$$+ \gamma n (r(\pi | x, h_{t-1}) - V^* + V^* - r(x, h_{t-1})) \quad (21)$$

374 Rearranging terms gives:

$$V^* - r(\pi | x, h_{t-1}) \leq (1 - \frac{1}{\gamma n})(V^* - r(x, h_{t-1})) \quad (22)$$

375 Note that this holds for every history. Therefore, the result follows by noting that $X_t = V^* - r(x, H_t)$
 376 and taking expectations. \square

377 *Proof of Theorem 3 (Performance of γ -Approximate Greedy Policies).* Let H_t denote the (random)
 378 history after t actions of policy π . Define the gap random variables $X_t = V^* - r(x, H_t)$, which
 379 measure how far we are from optimality after t steps. By Lemma 6 we have that $\mathbb{E}[X_i]$ decreases
 380 exponentially:

$$\mathbb{E}[X_t] \leq (1 - 1/(\gamma n))\mathbb{E}[X_{t-1}]. \quad (23)$$

381 Iterating this inequality from $t = 1$ to n :

$$\mathbb{E}[X_n] \leq (1 - 1/(\gamma n))^n \mathbb{E}[X_0] \quad (24)$$

382 Since $X_0 = V^* - r(\pi | x, H_0)$ where H_0 is the empty history, and $\mathbb{E}[X_0] = V^* - V(\theta)$:

$$V^* - V(\theta) \leq (1 - 1/(\gamma n))^n V^* \leq e^{-1/\gamma} V^*. \quad (25)$$

383 When $\gamma = 1$, we recover the classical $(1 - 1/e)$ -approximation of the exact greedy policy. \square

384 **Lemma 7.** Let $p(x)$ be a probability distribution, and let $g(x)$ be a real valued function. Define
 385 $\mathbb{E}_p[g(x)] = \int p(x) g(x) dx$. Now define a new distribution $p'(x)$ by reweighting $p(x)$ with the factor
 386 $e^{g(x)}$: $p'(x) = \frac{p(x) e^{g(x)}}{Z}$, where $Z = \mathbb{E}_p[e^{g(x)}] = \int p(x) e^{g(x)} dx$.

387 Then,

$$\mathbb{E}_{p'}[g(x)] \geq \mathbb{E}_p[g(x)]$$

388 *Proof.* We want to show

$$\frac{1}{Z} \mathbb{E}_p[e^{g(x)} g(x)] \geq \mathbb{E}_p[g(x)].$$

389 Equivalently,

$$\mathbb{E}_p[e^{g(x)} g(x)] \geq Z \mathbb{E}_p[g(x)] = \mathbb{E}_p[e^{g(x)}] \mathbb{E}_p[g(x)].$$

390 Thus it suffices to show

$$\mathbb{E}_p[e^{g(x)} g(x)] \geq \mathbb{E}_p[e^{g(x)}] \mathbb{E}_p[g(x)].$$

391 Let $Y = g(x)$ be a real-valued random variable under p . We claim

$$\mathbb{E}[e^Y Y] \geq \mathbb{E}[e^Y] \mathbb{E}[Y].$$

392 Rewrite this as

$$\mathbb{E}[e^Y (Y - \mathbb{E}[Y])] = \text{Cov}(e^Y, Y) \geq 0.$$

393 But $\text{Cov}(e^Y, Y) \geq 0$ holds because e^Y is a strictly increasing function of Y . By a standard result
 394 (e.g., Chebyshev's sum inequality), an increasing function of a random variable is positively correlated
 395 with that variable. \square

396 *Proof of Theorem 6 (Policy Improvement).* To establish the theorem, it suffices to show that for all
 397 contexts x and histories h_{t-1} :

$$\Delta(\hat{\pi}|x, h_{t-1}) \geq \Delta(\pi_0|x, h_{t-1}) \quad (26)$$

398 This improvement in expected marginal gain directly implies the desired approximation bounds.

Dog	Cat	Elephant	Lion	Tiger
Giraffe	Panda	Kangaroo	Horse	Penguin
Dolphin	Koala	Zebra	Wolf	Shark
Eagle	Cheetah	Bear	Monkey	Snake

Figure 2: Animals in the 20Q game.

399 From the optimality conditions of **KL-SubPO** in (10), we know that:

$$\hat{\pi}(a|x, h_{t-1}) = \frac{1}{Z(x, h_{t-1})} \pi_0(a|x, h_{t-1}) \exp\left(\frac{1}{\beta} \Delta(a|x, h_{t-1})\right), \quad (27)$$

400 where $Z(x, h_{t-1})$ is the normalization factor:

$$Z(x, h_{t-1}) = \sum_{a' \in \mathcal{A}} \pi_0(a'|x, h_{t-1}) \exp\left(\frac{1}{\beta} \Delta(a'|x, h_{t-1})\right). \quad (28)$$

401 Fix any context-history pair (x, h_{t-1}) . Let $p(a) = \pi_0(a|x, h_{t-1})$ and define $g(a) = \frac{1}{\beta} \Delta(a|x, h_{t-1})$.

402 Then $\hat{\pi}$ can be written as:

$$p'(a) = \frac{p(a) \exp(g(a))}{\sum_{a'} p(a') \exp(g(a'))} \quad (29)$$

403 By Lemma 7, we have:

$$\mathbb{E}_{a \sim p'}[g(a)] \geq \mathbb{E}_{a \sim p}[g(a)] \quad (30)$$

404 which directly implies the desired improvement property.

405 For $\beta_2 < \beta_1$, we can express $\pi(\cdot|\cdot; \hat{\theta}(\beta_2))$ as a reweighting of $\pi(\cdot|\cdot; \hat{\theta}(\beta_1))$:

$$\begin{aligned} \pi(a|x, h_{t-1}; \hat{\theta}(\beta_2)) &= \frac{1}{\hat{Z}(x, h_{t-1})} \hat{\pi}(a|x, h_{t-1}) \\ &\quad \times \exp\left(\frac{1}{\delta} \Delta(a|x, h_{t-1})\right), \end{aligned}$$

406 where $\delta = 1/\beta_2 - 1/\beta_1$. Applying our previous result twice yields:

$$\begin{aligned} \Delta(\pi(\cdot|\cdot; \hat{\theta}(\beta_2))|x, h_{t-1}) &\geq \Delta(\pi(\cdot|\cdot; \hat{\theta}(\beta_1))|x, h_{t-1}) \\ &\geq \Delta(\pi_0(\cdot|\cdot)|x, h_{t-1}). \end{aligned}$$

407 This establishes the monotonicity of γ' with respect to β . □

408 B Twenty Questions Experiment

409 The last experiment is a 20Q game (Karbasi et al., 2012) with 20 animals. The agent is represented
410 by an LLM and it is optimized against a user, which is also represented by an LLM. The animals are
411 listed in Figure 2 and the horizon of the game is $n = 6$.

412 Both the agent and user are implemented using Llama-3.1-8B. The role of the agent is

413 *You try to guess an animal. Respond with up to 6 words.*

414 The question of the agent is generated using prompt

415 *Ask a question.*

416 It is conditioned on the history of the conversation. The role of the user is

Question	Answer	Reward
Does it live on land?	Yes	0.100
Does it have four legs?	Yes	0.200
Does it have a tail?	Yes	0.200
Does it primarily eat plants?	No	0.600
Does it have sharp claws?	Yes	0.600
Is it a carnivorous mammal?	Yes	0.600

Figure 3: One 20Q game between the user and agent. The target animal is dog.

Bluetooth Speaker Phone Charger Air Fryer Yoga Mat
 Water Bottle Ring Doorbell Echo Dot Wireless Earbuds
 Protein Powder LED Strip Lights Portable Power Bank
 Coffee Maker Weighted Blanket Desk Lamp Wireless Mouse
 Reusable Straws Robot Vacuum Shower Curtain
 Cast Iron Skillet Kindle Paperwhite

Figure 4: Products in the Amazon selection game.

417 *Answer with Yes or No. No period.*

418 The answer of the user is generated by prompt

419 *You think of [animal]. You are asked: [question]*

420 where [animal] is replaced by the target animal name from Figure 2 and [question] is replaced by the
 421 last question of the agent. The reward is the fraction of eliminated animals. The animal is eliminated
 422 if at least one property of the animal disagrees with at least one answer of the user. One conversation
 423 between the user and agent is shown in Figure 3.

424 C Amazon Product Selection Experiment

425 The last experiment is a product selection game on a set of 20 Amazon products. The agent is tasked
 426 with narrowing down to a specific product by asking yes/no questions. The products are listed in
 427 Figure 4 and the horizon of the game is $n = 4$.

428 Both the agent and user are implemented using Llama-3.1-8B. The agent is provided with the system
 429 message:

430 *You are playing a 20 Questions game to guess an Amazon product from this list: [list of products].*
 431 *Ask clear yes/no questions to efficiently narrow down the possibilities. Keep questions concise*
 432 *(ideally under 10 words). The user will only respond with Yes or No.*

433 The question of the agent is generated using prompt:

434 *Ask a question.*

435 It is conditioned on the history of the conversation. The user’s response is generated by prompt:

436 *You think of [product]. You are asked: [question]*

437 where [product] is replaced by the target product name from Figure 4 and [question] is replaced
 438 by the last question of the agent. The reward is the fraction of eliminated products. A product is
 439 eliminated if its response to a question differs from the target product’s response to the same question.
 440 This reward calculation creates a natural submodular structure as questions eliminate overlapping
 441 subsets of products.

Question	Answer	Reward
Is the product electronic?	Yes	0.350
Can the product be held in your hand?	Yes	0.550
Does the product plug into a wall outlet?	No	0.850
Does the product require charging?	Yes	0.850

Figure 5: Example run of the Amazon product selection game.

442 The baseline model achieves a reward of 0.837 ± 0.005 . We compare this with various configurations
443 of our methods: (1) **KL-SubPO** with standardized trajectory rewards achieves 0.841 ± 0.004 ; (2)
444 **KL-SubPO** with clipped rewards from below at 0 achieves 0.858 ± 0.004 ; (3) **KL-PO** with standardized
445 per-step gains achieves 0.828 ± 0.005 ; (4) **KL-PO** with clipped rewards from below at 0 achieves
446 0.847 ± 0.004 .