
Online Decision Making with Generative Action Sets

Jianyu Xu

Carnegie Mellon University
Pittsburgh, PA 15213
jianyux@andrew.cmu.edu

Vidhi Jain

Carnegie Mellon University
Pittsburgh, PA 15213
vidhij2@andrew.cmu.edu

Bryan Wilder

Carnegie Mellon University
Pittsburgh, PA 15213
bwilder@andrew.cmu.edu

Aarti Singh

Carnegie Mellon University
Pittsburgh, PA 15213
aarti@andrew.cmu.edu

Abstract

With advances in generative AI, decision-making agents can now dynamically create new actions during online learning, but action generation typically incurs costs that must be balanced against potential benefits. We study an online learning problem where an agent can generate new actions at any time step by paying a one-time cost, with these actions becoming permanently available for future use. The challenge lies in learning the optimal sequence of two-fold decisions: which action to take and when to generate new ones, further complicated by the triangular tradeoffs among exploitation, exploration and *creation*. To solve this problem, we propose a doubly-optimistic algorithm that employs Lower Confidence Bounds (LCB) for action selection and Upper Confidence Bounds (UCB) for action generation. Empirical evaluation on healthcare question-answering datasets demonstrates that our approach achieves favorable generation-quality tradeoffs compared to baseline strategies. From theoretical perspectives, we prove that our algorithm achieves the optimal regret of $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$, providing the first sublinear regret bound for online learning with expanding action spaces.

1 Introduction

Sequential decision-making problems involve agents repeatedly selecting actions from a candidate set to maximize cumulative reward. Traditional approaches assume a fixed set of available actions, focusing on the exploration-exploitation tradeoffs: balancing empirically high-reward actions (exploitation) against less-tested alternatives (exploration). However, advances in generative AI have introduced a new paradigm where contemporary systems can dynamically *expand* their action spaces by *creating* novel actions over time. This capability introduces an additional strategic dimension that agents should also balance immediate performance with strategic investments in future capabilities enabled by new actions. Consider the following motivating scenarios:

Example 1.1 (Healthcare Question-Answering Systems). *AI-powered healthcare platforms must decide between reusing existing vetted responses from their FAQ libraries or investing in creating new, tailored responses for novel patient inquiries. Each custom response requires costly expert review and validation (potentially hundreds of dollars when accounting for clinical expertise). However, once created and vetted, these responses become reusable assets. When a patient in a given region asks “What are healthy meals during pregnancy?”, the system faces a critical choice: provide a generic response about pregnancy nutrition, or invest in creating a new response more specific to typical foods in that region, benefiting hundreds of future expectant mothers in similar settings.*

Example 1.2 (Personalized Advertisement). *An advertising platform may initially start with a finite set of ad templates for different user contexts. Over time, the platform observes new user segments and decides to design specialized ads (with initial design and production costs) perfectly customized to the new user subgroups. Once created, these specialized ads become available for future targeting at no additional cost.*

In both scenarios, the agent must decide at each time step whether to select an existing action or pay a one-time cost to instantiate a new action perfectly suited to the observed context. This introduces a novel *create-to-reuse* problem that goes beyond traditional exploration-exploitation tradeoffs.

Problem formulation. In this work, we study a contextual bandit problem with an *actively expanding* action space. At each time t , the agent first observes a context x_t . Then it can either

- (a) Pull an *existing* arm at no cost but incur some loss, or
- (b) Pay a fixed one-time cost to generate a *new* arm and incur zero(0) loss.

A detailed problem description is shown as follows:

Initialization: Context-to-action oracle $\mathcal{A}(\cdot)$. A library $S_1 = \{f, \mathcal{A}(f)\}$ with context keys f and vetted custom actions $\mathcal{A}(f)$.
 For $t = 1, 2, \dots, T$:

1. Observe $x_t \in \mathbb{R}^d$ (patient question arrives).
2. The algorithm decides whether to create a customized response to x_t . If YES, then
 - (i) Generation oracle produces and deploys $a_t = \mathcal{A}(x_t)$ (custom response to x_t).
 - (ii) Receive a fixed loss c (creation cost).
 - (iii) Update $S_{t+1} := S_t \cup \{x_t : a_t\}$ (add new context-action pair to the library).
3. If NO, then
 - (i) Select an existing context key $f_t \in S_t$ and retrieve $a_t = S_t(f_t)$.
 - (ii) Receive a loss $l_t := d(x_t, f_t) + N_t$ (noisy mismatch penalty).
 - (iii) Update $S_{t+1} := S_t$ (library unchanged).

Here $d(x, a) := (x - a)^\top W(x - a)$ is a quadratic distance function for $x, a \in \mathbb{R}^d$, with an *unknown* $W \in \mathbb{S}_+^d$. N_t is an i.i.d. σ -subGaussian noise. Please refer to Appendix B for a rigorous problem setup and all technical assumptions.

We highlight two important features of this formulation. First, step (b) is notable in that the agent accesses the action space only through an oracle that is prompted with the context x_t . By contrast, applications of previous bandit formulations would operate directly in a separate action space. In our motivating settings though, this action space might be very complicated: the space of all possible texts, all possible antibodies, and so on. In an increasing number of applications, practitioners deal with such action spaces through calling a separate model that returns a tailored response to a given x_t (e.g. an LLM). The bandit algorithm thus operates as a decision-making layer on top of this oracle. The second important feature of the problem formulation is that, once generated, the new arm can be reused in future rounds without incurring additional expense. The key is to judiciously decide when to pay the cost of adding such a specialized action and when to rely on existing arms.

This setting presents two fundamental challenges that prevent existing bandits and online learning methods from solving this problem. First, we face a triangular trade-off among three competing objectives: exploitation (using known good arms), exploration (learning about uncertain arms), and *creation* (which lies between exploration and exploitation, as it satisfies the immediate need at current time while also enriching the action set for the future). Second, we have no prior experience with potential new actions nor unlimited freedom to generate arbitrary ones – each creation must be specifically tailored to the current context.

Summary of Contributions Our main contributions are fourfold:

1. **Problem Modeling:** We establish a new problem formulation that allows for costly expansion of the action space in online learning, formalizing the create-to-reuse framework.

Algorithm 1 Doubly-Optimistic Algorithm

```
1: Initialization  $\Sigma_0 = \lambda \cdot I_{d^2}, b_0 = \vec{0}_{d^2}, S_1 = \{\vec{1}_d\}$ , hyper parameter  $\alpha$ .
2: for  $t = 1, 2, \dots, T$  do
3:   Observe  $x_t \in \mathbb{R}^d$ 
4:   for  $\forall f \in S_t$  do
5:     Denote  $\phi(x, f) := \text{Vec}[(x - f)(x - f)^\top] \in \mathbb{R}^{d^2}$  and
      
$$\begin{aligned} \Delta_t(x, f) &:= \alpha \cdot \sqrt{\phi(x_t, f)^\top \Sigma_{t-1}^{-1} \phi(x_t, f)}, \quad \bar{d}_t(x, f) := \phi(x, f)^\top \Sigma_{t-1}^{-1} b_{t-1} \\ \hat{d}_t(x, f) &:= \bar{d}_t(x, f) + \Delta_t(x, f), \quad \check{d}_t(x, f) := \bar{d}_t(x, f) - \Delta_t(x, f) \end{aligned} \quad (1)$$

6:   end for
7:   Select  $f_t := \arg\min_{f \in S_t} \check{d}_t(x_t, f)$ .
8:   if  $Z_t == 1$  with  $Z_t \sim \text{Ber}(\min\{1, \frac{1}{c} \cdot \hat{d}_t(x_t, f_t)\})$  as an i.i.d. Bernoulli random variable
     then
9:     Take action  $a_t = x_t$  at a cost  $c$ .
10:    Receive loss  $l_t = 0$ .
11:    Update action set  $S_{t+1} = S_t \cup \{a_t\}$ .
12:    Keep  $\Sigma_t := \Sigma_{t-1}$  and  $b_t := b_{t-1}$  without updating.
13:   else
14:     Take action  $a_t = f_t$  at no cost.
15:     Receive loss  $l_t = d(x_t, f_t) + N_t$ .
16:     Update action set  $S_{t+1} = S_t$ .
17:     Update parameters
      
$$\Sigma_t := \Sigma_{t-1} + \phi(x_t, a_t)\phi(x_t, a_t)^\top, \quad b_t := b_{t-1} + l_t \cdot \phi(x_t, a_t). \quad (2)$$

18:   end if
19: end for
```

2. **Algorithmic Framework:** We propose a *doubly-optimistic* algorithm that uses Lower Confidence Bounds (LCB) when selecting among existing actions, and Upper Confidence Bounds (UCB) when deciding whether to generate new actions. This design simultaneously exploits near-optimal actions and enables creation without excessive hesitation.
3. **Empirical Validation:** We conduct experiments on real-world healthcare question-answering datasets, demonstrating that our approach achieves favorable generation-quality tradeoffs compared to baselines. Our results show the method gracefully interpolates between pure reuse and always-create policies while maintaining superior performance.
4. **Optimal Regret Guarantees:** Under a semi-parametric loss model, our algorithm achieves $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$ expected regret, where T is the time horizon and d is the dimension of covariates. We prove this rate is optimal by establishing a matching $\Omega(T^{\frac{d}{d+2}})$ information-theoretic lower bound.

Technical Novelty. The crux of our approach is a **double optimism** principle, which resolves the unique challenge of balancing creation with exploration/exploitation. Among existing actions, we rely on their *LCB* comparisons to both exploit high-performing actions and continue exploring uncertain ones. When evaluating creation decisions, we compare the *UCB* loss of the best existing action against the fixed generative cost, triggering creation with appropriate probability. This double optimism perspective naturally maximizes the long-term value of new actions while tightly controlling worst-case regret.

Related Works. Our work is closely related to multi-armed bandits, bandits with constraints, and facility location problems. We present a detailed discussion in Appendix A.

2 Algorithm

To solve the contextual bandits problem with expanding action space, we propose our “Doubly-Optimistic” algorithm. In this section, we present the algorithm design and highlight its properties. We will analyze and bound its cumulative regret in the next section.

The pseudocode of our algorithm is displayed as Algorithm 1. At each time t , it inherits the linear-regression parameters Σ_{t-1} , b_{t-1} an action set S_t from $(t-1)$, and receives a context vector x_t from the nature. With Σ_{t-1} and b_{t-1} , it estimates the empirical loss of each existing action $f \in S_t$ as $\bar{d}_t(x_t, f)$, along with an uncertainty bound $\Delta_t(x_t, f)$. Then we take the following two steps to figure out the action a_t to take.

- (i) **Lower Confidence Bound (LCB) loss on existing actions.** For each action f , we calculate the LCB loss as $\check{d}_t(x_t, f) = \bar{d}_t(x_t, f) - \Delta_t(x_t, f)$. Then we select f_t as the arg-minimum of all $\check{d}_t(x_t, f)$ over all f . Note: we do not propose f_t immediately.
- (ii) **Upper Confidence Bound (UCB) chance to create a new action.** After retrieving f_t as the argmin of LCB losses, we turn to believe in its UCB loss $\hat{d}_t(x_t, f_t) = \bar{d}_t(x_t, f_t) + \Delta_t(x_t, f_t)$ while contrasting to the fixed arm-adding cost c . With a probability of $\min\{1, \frac{\hat{d}_t(x_t, f_t)}{c}\}$, we create and take the new action $a_t = x_t$, and update the action set $S_{t+1} = S_t \cup \{x_t\}$ accordingly. Otherwise, we take the existing action $a_t = f_t$, receive a random loss l_t , and update the parameters Σ_t and b_t accordingly.

As the argmin of LCB loss, f_t represents the least possible loss to suffer as an optimist, which balances exploration versus exploitation under uncertainties we possess from history. Similar methods are applied in a broad group of contextual bandits literature such as [Chu et al. \(2011\)](#).

As the probability $\frac{\hat{d}_t(x_t, f_t)}{c}$ induced by the UCB loss of f_t , it increases the chance of creating a new arm at x_t to the most (within a risk Δ_t we can tolerate). This design enables us to estimate the “necessity” of creation, bounding the total expected loss over a group of x_t ’s *before* the first action being created among them. We will explain this later in Lemma [G.5](#).

Computational complexity Algorithm 1 incurs a time complexity at $O(d^4 T \cdot K_{\max}) \leq O(d^4 T^2)$, as it compute matrix-to-vector products of d^2 -dimension for every action $f \in S_t$ at each round t , and there are at most T arms. By noting that the expected number of newly created arms is on the order of $O(T^{\frac{d}{d+2}})$, we can refine the *expected* complexity to $O(T^{\frac{2d+2}{d+2}})$. Moreover, the key step of updating the inverse covariance Σ_t^{-1} can be carried out in $O(d^4)$ via *Woodbury matrix identity* which states $(A + xx^T)^{-1} = A^{-1} - \frac{1}{1+x^T A^{-1} x} (A^{-1} x)(A^{-1} x)^T$. Although this non-linear complexity is necessary for achieving our optimal regret guarantees, we note that several standard techniques could potentially improve computational performance in practice, including low-rank approximation and randomized sketching methods.

Regret Bounds. We sequentially present our theoretical guarantees on the regret upper and lower bounds, as the following two theorems. (The regret analysis is deferred to Appendix [C](#)).

Theorem 2.1 (Regret upper bound). *With assumptions made in Appendix [C.1](#), the expected regret of our Algorithm 1 is upper bounded by $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$.*

Theorem 2.2 (Regret lower bound). *For any online learning algorithm, there exists an instance of problem setting presented in Appendix [B](#), such that the regret is at least $\Omega(T^{\frac{d}{d+2}})$ with respect to T .*

Numerical Results. We conduct numerical experiments to validate our method’s performance. We first run the original algorithm on low-dimensional synthetic data to demonstrate the regret dependence on T . Then we adapt our algorithm to real-world healthcare Q&A scenarios and show better tradeoffs between generation cost and mismatching loss compared to baselines. A detailed presentation of numerical results can be found in Appendix [D](#).

Conclusion. In this paper, we introduced an online decision-making problem where new actions can be generated on the fly, at a fixed cost, and then reused indefinitely. To address the balance among exploitation, exploration, and creation, we proposed a doubly-optimistic algorithm that achieves $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$ optimal regret (validated in theory and simulations). We also implemented our algorithm on a real-world healthcare Q&A dataset to make decisions on generating new answers v.s. applying an FAQ. Our results open up new avenues for optimizing creation decisions in online learning, with potential extensions to broader loss models and flexible creation costs.

References

- Agarwal, A., Hsu, D., Kale, S., Langford, J., Li, L., and Schapire, R. (2014). Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning (ICML-14)*, pages 1638–1646.
- Agrawal, S. and Devanur, N. (2016). Linear contextual bandits with knapsacks. *Advances in neural information processing systems*, 29.
- Angluin, D. (1988). Queries and concept learning. *Machine learning*, 2(4):319–342.
- Auer, P. (2002). Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine Learning Research*, 3(Nov):397–422.
- Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002). Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256.
- Badanidiyuru, A., Kleinberg, R., and Slivkins, A. (2013). Bandits with knapsacks. In *Annual Symposium on Foundations of Computer Science (FOCS-13)*, pages 207–216. IEEE.
- Bellman, R. (1958). Studies in the mathematical theory of inventory and production.
- Chen, B., Chao, X., and Ahn, H.-S. (2019). Coordinating pricing and inventory replenishment with nonparametric demand learning. *Operations Research*, 67(4):1035–1052.
- Cheng, H.-T., Koc, L., Harmsen, J., Shaked, T., Chandra, T., Aradhye, H., Anderson, G., Corrado, G., Chai, W., Ispir, M., et al. (2016). Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*, pages 7–10.
- Chu, W., Li, L., Reyzin, L., and Schapire, R. (2011). Contextual bandits with linear payoff functions. In *International Conference on Artificial Intelligence and Statistics (AISTATS-11)*, pages 208–214.
- Farquhar, G., Gustafson, L., Lin, Z., Whiteson, S., Usunier, N., and Synnaeve, G. (2020). Growing action spaces. In *International Conference on Machine Learning*, pages 3040–3051. PMLR.
- Fotakis, D. (2008). On the competitive ratio for online facility location. *Algorithmica*, 50(1):1–57.
- Guo, X., Kulkarni, J., Li, S., and Xian, J. (2020). On the facility location problem in online and dynamic models. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2020)*, pages 42–1. Schloss Dagstuhl–Leibniz-Zentrum für Informatik.
- He, X., Liao, L., Zhang, H., Nie, L., Hu, X., and Chua, T.-S. (2017). Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*, pages 173–182.
- Huang, S.-J., Jin, R., and Zhou, Z.-H. (2010). Active learning by querying informative and representative examples. *Advances in neural information processing systems*, 23.
- Immorlica, N., Sankararaman, K. A., Schapire, R., and Slivkins, A. (2019). Adversarial bandits with knapsacks. In *Annual Symposium on Foundations of Computer Science (FOCS-19)*, pages 202–219. IEEE.
- Jin, C., Netrapalli, P., Ge, R., Kakade, S. M., and Jordan, M. I. (2019). A short note on concentration inequalities for random vectors with subgaussian norm. *arXiv preprint arXiv:1902.03736*.
- Kaplan, H., Naori, D., and Raz, D. (2023). Almost tight bounds for online facility location in the random-order model. In *Proceedings of the 2023 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1523–1544. SIAM.
- Kothawade, S., Chopra, S., Ghosh, S., and Iyer, R. (2022). Active data discovery: Mining unknown data using submodular information measures. In *ICML workshop on Adaptive Experimental Design and Active Learning in the Real World*.
- Lai, T. L. and Robbins, H. (1985). Asymptotically efficient adaptive allocation rules. *Advances in applied mathematics*, 6(1):4–22.

- Lee, H., Im, J., Jang, S., Cho, H., and Chung, S. (2019). Melu: Meta-learned user preference estimator for cold-start recommendation. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1073–1082.
- Li, L., Chu, W., Langford, J., and Schapire, R. E. (2010). A contextual-bandit approach to personalized news article recommendation. In *Proceedings of the 19th international conference on World wide web*, pages 661–670.
- Liu, S., Jiang, J., and Li, X. (2022). Non-stationary bandits with knapsacks. *Advances in Neural Information Processing Systems*, 35:16522–16532.
- Liu, Y., Logan, B., Liu, N., Xu, Z., Tang, J., and Wang, Y. (2017). Deep reinforcement learning for dynamic treatment regimes on medical registry data. In *2017 IEEE international conference on healthcare informatics (ICHI)*, pages 380–385. IEEE.
- Meyerson, A. (2001). Online facility location. In *Proceedings 42nd IEEE Symposium on Foundations of Computer Science*, pages 426–431. IEEE.
- Mieghem, J. A. V. and Rudi, N. (2002). Newsvendor networks: Inventory management and capacity investment with discretionary activities. *Manufacturing & Service Operations Management*, 4(4):313–335.
- Moor, M., Banerjee, O., Abad, Z. S. H., Krumholz, H. M., Leskovec, J., Topol, E. J., and Rajpurkar, P. (2023). Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265.
- Rajendran, P. T., Espinoza, H., Delaborde, A., and Mraidha, C. (2023). Unsupervised unknown unknown detection in active learning. In *The IJCAI-2023 AISafety and SafeRL Joint Workshop*.
- Rajpurkar, P., Chen, E., Banerjee, O., and Topol, E. J. (2022). Ai in health and medicine. *Nature medicine*, 28(1):31–38.
- Settles, B. (2009). Active learning literature survey.
- Seung, H. S., Oppor, M., and Sompolinsky, H. (1992). Query by committee. In *Proceedings of the fifth annual workshop on Computational learning theory*, pages 287–294.
- Slivkins, A. et al. (2019). Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning*, 12(1-2):1–286.
- Xu, J., Wang, X., Wang, Y.-X., and Jiang, J. (2025a). Joint pricing and resource allocation: An optimal online-learning approach. *arXiv preprint arXiv:2501.18049*.
- Xu, J., Wang, Y., Chen, X., and Wang, Y.-X. (2025b). Dynamic pricing with adversarially-censored demands. *arXiv preprint arXiv:2502.06168*.
- Zador, P. L. (1964). *Development and evaluation of procedures for quantizing multivariate distributions*. Stanford University.
- Zhao, P., Shan, J.-W., Zhang, Y.-J., and Zhou, Z.-H. (2024). Exploratory machine learning with unknown unknowns. *Artificial Intelligence*, 327:104059.

Appendix

A Related Works

Here we discuss related literature on the most relevant topics in online decision making, as well as on broader fields including active learning, digital healthcare, recommendation system, and inventory management.

Multi-Armed and Contextual Bandits. The multi-armed bandit (MAB) problem has been extensively studied since [Lai and Robbins \(1985\)](#). The classic framework ([Auer et al., 2002](#); [Agarwal et al., 2014](#)), that a decision-maker repeatedly selects from a fixed set of arms, was extended to contextual bandits ([Li et al., 2010](#); [Chu et al., 2011](#)) where rewards depend on observable contexts. The crux is to balance exploration and exploitation with the goal of *regret* minimization. Please refer to [Slivkins et al. \(2019\)](#) for a comprehensive discussion.

Online Facility Location. Online facility location (OFL), studied by [Meyerson \(2001\)](#), [Fotakis \(2008\)](#), and [Guo et al. \(2020\)](#), is closely related to our formulation. In OFL, algorithms decide whether to open new facilities or assign requests to existing ones, minimizing facility costs plus assignment distances. While structurally similar to our problem, there are crucial differences. First, OFL assumes *known* distance metrics, while we must learn *unknown* parameters defining distances. Second, OFL *automatically* assigns points to nearest facilities, while we must *actively* select actions under uncertainty. Therefore, OFL involves a *two-way* trade-offs between immediate costs and future benefits, whereas our problem requires a *three-way* balance between exploitation, exploration, and creation, necessitating our novel algorithmic approach.

Online Learning with Resource Constraints. Another line of related research studies resource-limited bandits, such as “bandits with knapsack (BwK)” ([Badanidiyuru et al., 2013](#)) and its versions ([Agrawal and Devanur, 2016](#); [Immorlica et al., 2019](#); [Liu et al., 2022](#)). In these scenarios, each arm-pulling consumes some portion of a finite resource (e.g., budget, time, or capacity), and the algorithm aims to optimize the cumulative reward before resources run out. However, these approaches cannot be directly applied to our problem because of a key difference in resource consumption patterns. In BwK, resource consumption only affects the current period’s decision-making. In contrast, our setting involves a one-time cost for creating new arms that provides benefits across all future periods through expansion of the action space. Besides, BwK mostly assumes a *hard* constraint on budgets, while we adopt a *soft* constraint as an additional cost in our problem setting.

Active Learning Active learning frameworks fundamentally embody the exploration-exploitation-creation paradigm by allowing algorithms to strategically choose their training data, thereby naturally connecting to sequential decision-making with expanding action spaces. [Settles \(2009\)](#) established the theoretical foundations for query selection strategies, while membership query synthesis approaches ([Angluin, 1988](#)) demonstrated how active learners can create entirely new query types rather than merely selecting from existing unlabeled data pools. Query-by-Committee methods ([Seung et al., 1992](#)) and extended through frameworks like QUIRE by [Huang et al. \(2010\)](#) show how multiple learning strategies can be combined to create adaptive query selection policies that balance informativeness and representativeness. Closer work on meta-active learning and the “Growing Action Spaces” framework by [Farquhar et al. \(2020\)](#) directly address expanding action spaces through curriculum learning approaches that progressively grow query complexity. The create-to-reuse framework maps directly onto active learning’s core mechanisms: systems invest computational effort in synthesizing new query types, developing committee-based strategies, and learning meta-policies for query selection, creating reusable query generation mechanisms and adaptive selection strategies that can be applied across different datasets, domains, and learning tasks, while continuously expanding their query capabilities as they encounter new data distributions and learning scenarios.

Exploratory Learning for Unknown Unknowns. Another notable progress is the exploratory machine learning (ExML) framework ([Zhao et al., 2024](#)). The authors introduced a novel and insightful approach to address unexpected unknown unknowns by exploring additional feature information through environmental interactions within a budget constraint, where an optimal bandit identification strategy is proposed to guide the feature exploration. There are several follow-up developments ([Kothawade et al., 2022](#); [Rajendran et al., 2023](#)). Compared to our create-to-use framework, there are two main differences: On the one hand, their work addresses the strategic choices of “create” while the current exploratory decisions would not be “in use” of future decisions.

On the other hand, their cost serves as a budget consumption instead of a tradeoffs with cumulative utilities, analogous to the divergence between regret minimization and best-arm identification.

Digital Healthcare and Clinical Decision Support Digital healthcare and clinical decision support systems (CDSS) represent a rapidly evolving field where AI-powered systems must continuously balance the utilization of established medical knowledge with the creation of novel, patient-specific treatment protocols. Foundational work by [Rajpurkar et al. \(2022\)](#) on diagnostic AI systems and the comprehensive framework established by [Moor et al. \(2023\)](#) demonstrate how modern medical AI systems expand beyond narrow, single-task applications to flexible models capable of diverse medical reasoning tasks. Reinforcement learning approaches in critical care, particularly the systematic review by [Liu et al. \(2017\)](#) covering 21 RL applications in intensive care units, illustrate how these systems extend from discrete medication dosing decisions to continuous, multi-dimensional treatment optimization spaces. The create-to-reuse paradigm is particularly evident in precision medicine applications, where systems invest computational resources in developing personalized treatment protocols that can subsequently be applied to patients with similar phenotypic characteristics, effectively creating reusable clinical knowledge that scales across patient populations while maintaining individualized care quality.

Recommendation Systems and Personalization Recommendation systems research has evolved from static collaborative filtering approaches to sophisticated frameworks that dynamically balance the exploitation of existing user preferences with the creation of new personalized recommendation strategies. Neural Collaborative Filtering by [He et al. \(2017\)](#) and the Wide & Deep Learning framework by [Cheng et al. \(2016\)](#) established the foundation for deep learning approaches that can capture complex user-item interactions beyond traditional matrix factorization methods. Meta-learning approaches, particularly by [Lee et al. \(2019\)](#) demonstrate how recommendation systems can treat each user as a distinct learning task, creating personalized model parameters that generalize across different applications and contexts. It is worth mentioning that the multi-armed bandit approaches in recommendation systems ([Li et al., 2010](#)) naturally embody the exploration-exploitation-creation tradeoffs by continuously balancing known user preferences with the discovery of new content types and recommendation strategies. Our create-to-reuse framework directly parallels these systems’ core functionality: recommendation systems routinely invest computational resources in creating personalized embeddings, meta-learned initialization parameters, and graph neural network representations that serve as reusable templates for rapid adaptation to new users, items, and interaction modalities, while continuously expanding their action spaces through dynamic catalog growth and emerging user behavior patterns.

Inventory Management Inventory management and supply chain systems represent a mature operations research domain where organizations continuously face fundamental tradeoffs between optimizing existing supply chain capabilities and investing in new suppliers, products, or distribution channels. [Bellman \(1958\)](#) established the mathematical foundations of inventory theory, while dynamic capacity expansion models ([Mieghem and Rudi, 2002](#)) demonstrate how firms balance existing capacity utilization with flexible resource investments that create new operational capabilities. The problem of inventory management often coexists with revenue management ([Chen et al., 2019](#)), resource allocation ([Xu et al., 2025a](#)), and adversarial online learning ([Xu et al., 2025b](#)) that occurs frequently in modern supply chains. The create-to-reuse framework aligns naturally with supply chain decision-making: organizations invest upfront in new suppliers, products, or distribution capabilities that become reusable assets for future deployment across different demand scenarios.

B Problem Setup

We now formalize the problem of creating-to-reuse as an online decision-making framework. In order to demonstrate the problem setting, we start with the healthcare Q&A scenario described in Example 1.1. As an abstraction, each arriving patient question is represented as a d -dimensional *context* vector x_t in a learned semantic embedding space. The system maintains a context library S_t of vetted FAQ entries, implemented as a *hash table* where each context that has been previously added serves as a key to its corresponding custom respond (or generally the *action*) generated by an oracle $\mathcal{A}(\cdot)$. Crucially, the algorithm operates only in the context representation space by searching through context keys in S_t . When a new question x_t arrives, the algorithm makes decisions based on estimated losses and can either:

- (a) Decide to create a new custom response by paying a fixed cost c and adding context x_t as a new key to the library. The generation oracle $\mathcal{A}(\cdot)$ then automatically produces the tailored action $a_t = \mathcal{A}(x_t)$, and the pair (x_t, a_t) becomes permanently available for future reuse. *Or*
- (b) Select an existing context key $f \in S_t$ from the library. The system automatically retrieves the corresponding action $a_t = S_t(f) = \mathcal{A}(f)$ and deploys it for context x_t , incurring a mismatch loss $d(x_t, f)$ that reflects the difference between (1) the custom response to context x_t versus (2) the action tailored for another context f .

Technically, we consider the following problem setting.

Initialization: Context-to-action oracle $\mathcal{A}(\cdot)$. A library $S_1 = \{f, \mathcal{A}(f)\}$ with context keys f and vetted custom actions $\mathcal{A}(f)$.
 For $t = 1, 2, \dots, T$:

1. Observe $x_t \in \mathbb{R}^d$ (patient question arrives).
2. The algorithm decides whether to create a customized response to x_t . If YES, then
 - (i) Generation oracle produces and deploys $a_t = \mathcal{A}(x_t)$ (custom response to x_t).
 - (ii) Receive a fixed loss c (creation cost).
 - (iii) Update $S_{t+1} := S_t \cup \{x_t : a_t\}$ (add new context-action pair to the library).
3. If NO, then
 - (i) Select an existing context key $f_t \in S_t$ and retrieve $a_t = S_t(f_t)$.
 - (ii) Receive a loss $l_t := d(x_t, f_t) + N_t$ (noisy mismatch penalty).
 - (iii) Update $S_{t+1} := S_t$ (library unchanged).

In this formulation, $d(x_t, f_t)$ captures the expected mismatch loss when deploying an action originally designed for context f_t to serve context x_t . While this fundamentally reflects the difference between $\mathcal{A}(x_t)$ and $\mathcal{A}(f_t)$ in the action space, the algorithm can only estimate this through context-space relationships since it lacks direct access to $\mathcal{A}(x_t)$ (actions having not been generated yet).

For theoretical analysis, our main modeling assumption is that this mismatch can be captured by a *squared distance* function in the context space. In experiments, we consider other forms for the mismatch distance.

Assumption B.1 (Quadratic parametric loss). We assume the distance function satisfies

$$d(x, f) := (x - f)^\top W (x - f) \quad (3)$$

where $W \in \mathbb{S}_+^d$ is an **unknown** positive semi-definite $d \times d$ matrix. Accordingly, denote

$$\begin{aligned} w &:= \text{Vec}(W) \in \mathbb{R}^{d^2} \\ \phi(x, f) &:= \text{Vec}[(x - f)(x - f)^\top] \in \mathbb{R}^{d^2}, \end{aligned} \quad (4)$$

and we have an equivalent definition as $d(x, f) := \phi(x, f)^\top w$.

Why we assume a quadratic parametric loss? The motivation is that contexts are embedded in a space where different dimensions capture semantically relevant information. The cost of reusing an action designed for one context when serving another can be modeled as a distance on this representation space, although the exact importance weighting of different semantic dimensions (captured by matrix W) is unknown to the learner. Since $d(x_t, x_t) = 0$, our formulation measures the *excess* cost due to not generating a custom action for each x_t . This fits scenarios where the algorithm interacts with complex action spaces through oracle $\mathcal{A}(x_t)$ (human expert or generative model); our aim is to achieve good performance relative to this oracle’s capabilities. Modeling $d(\cdot, \cdot)$ as a squared distance function captures more structure than linear parametric choices while remaining more tractable than nonparametric formulations. Furthermore, the empirical results our algorithm performs on real-world Healthcare Q&A datasets validate the robustness of this modeling.

Goal of Algorithm Design. Our goal is to minimize the expected *total loss*. We will rigorously define the performance metric and technical assumptions at the beginning of Appendix C.

C Regret Analysis

In this section, we provide a regret analysis of our algorithm. We first state the performance metric and necessary technical assumptions. Then we present the main theorem on the algorithmic regret upper bound. Finally, we provide a corresponding lower bound that matches the leading term of the upper bound with respect to T .

C.1 Definitions and Assumptions

As we have stated by the end of Appendix B, our goal is to minimize the *total loss*. In order to measure the performance, we adopt the expected *regret* as the loss metric, which is defined as follows:

Definition C.1 (Optimal and Regret). Denote the minimal expected loss¹ that is *achievable* in hindsight as OPT_h , which equals:

$$OPT_h := \min_{S := \{S_1, S_2, \dots, S_T, S_{T+1} | S_{t+1} \setminus S_t \subseteq \{x_t\}\}} c \cdot |S_{T+1}| + \sum_{t=1}^T \min_{f \in S_{t+1}} d(x_t, f). \quad (5)$$

There also exists a *non-achievable* minimal loss denoted as OPT_o , which is only accessible by an omniscient oracle that knows $\{x_t\}_{t=1}^T$ and selects an optimal option set ahead of time:

$$OPT_o := \min_S c \cdot |S| + \sum_{t=1}^T \min_{f \in S} d(x_t, f). \quad (6)$$

From the definition, we know that $OPT_o \leq OPT_h$. Also, denote the expected loss obtained by our algorithm as ALG , which equals:

$$ALG := c \cdot |S_{T+1}| + \sum_{t=1}^T \min_{f \in S_{t+1}} d(x_t, f). \quad (7)$$

Define the regret REG as the expected loss difference² between OPT_h and ALG .

$$REG := \mathbb{E}[ALG - OPT_h] \quad (8)$$

We then make two distributional assumptions on the covariates and the noises, respectively.

Assumption C.2 (Covariate distribution and norm bound). Assume $x_t \in \mathbb{R}^d, t = 1, 2, \dots, T$ are drawn from independent and identical distributions (i.i.d.), with $d \geq 2$. Also, assume a norm bound as $\|x_t\|_2 \leq 1$.

Assumption C.2 is necessary for us to effectively learn the metric matrix W through online linear regression. For the same reason, we assume a subGaussian noise on the observations as follows:

Assumption C.3 (Noise distribution). Assume that $N_t \in \mathbb{R}, t = 1, 2, \dots, T$ are drawn from σ -subGaussian i.i.d., where σ is a universal constant.

C.2 Regret Bounds

In this subsection, we sequentially present our theoretical guarantees on the regret upper and lower bounds, as the following two theorems.

Theorem C.4 (Regret upper bound). *With assumptions made in Appendix C.1, the expected regret of our Algorithm 1 is upper bounded by $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$.*

Proof Sketch. We prove Theorem C.4 in the following sequence:

1. (Lemma G.1) We upper bound the non-achievable minimal loss as $OPT_o = O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}})$. This is proved by a fine-grid covering of the space.
2. (Lemma G.3) We upper bound the algorithmic loss ALG within a constant competitive ratio of OPT_o adding cumulative prediction errors: $\mathbb{E}[ALG] = O(\mathbb{E}[OPT_o] + \sum_{t=1}^T \Delta_t(x_t, f_t))$. To prove this, we divide $\{x_t\}$'s into "good" and "bad" groups, and bound their excess loss respectively.

¹Expectation taken over observation noises only. Same for the definition of OPT_o .

²Expectation taken over the $\{x_t\}_{t=1}^T$ series.

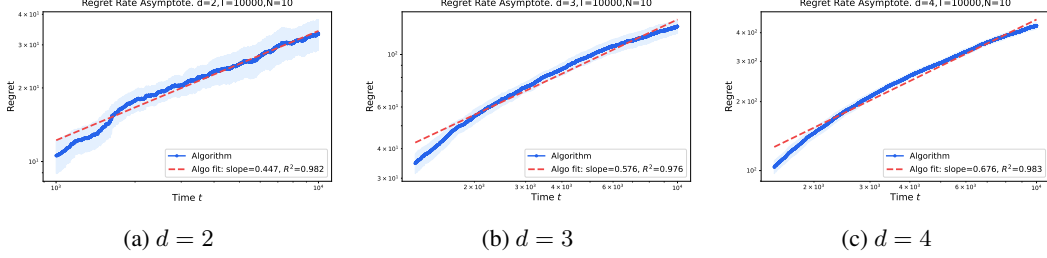


Figure 1: Regret curves for $T = 10000$ and $d = 2, 3, 4$ in log-log scales, repeated by $N = 10$ epochs. The slope of the linear asymptote under log-log diagram indicates the power dependence of regret on T , which should be $\frac{d}{d+2}$.

3. (Lemma G.8) We upper bound the excess risk $\mathbb{E}[\sum_{t=1}^T \Delta_t(x_t, f_t)] = O(d\sqrt{T \log T})$ by standard online linear regression (similar to Chu et al. (2011) by replacing d with d^2).
4. Finally, we derive the regret upper bound as $REG = \mathbb{E}[ALG - OPT_h] = O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$ according to the three steps above.

Please refer to Appendix G for all technical details of this proof, including rigorous statements of lemmas and derivations of inequalities. ■

To show the optimality of the regret upper bound proposed above, we present the information-theoretic lower regret bound.

Theorem C.5 (Regret lower bound). *For any online learning algorithm, there exists an instance of problem setting presented in Appendix B, such that the regret is at least $\Omega(T^{\frac{d}{d+2}})$ with respect to T (despite the dependence on d).*

We defer the proof to Appendix G.6. The main idea is to apply the $\Omega(K^{-\frac{2}{d}})$ lower bound for the K-nearest-neighbors (K-NN) problem, along with an optimal choice of K that balance this term with $c \cdot K$. Theorem C.5 indicates that our algorithm achieves an optimal regret with respect to T .

D Empirical Performance

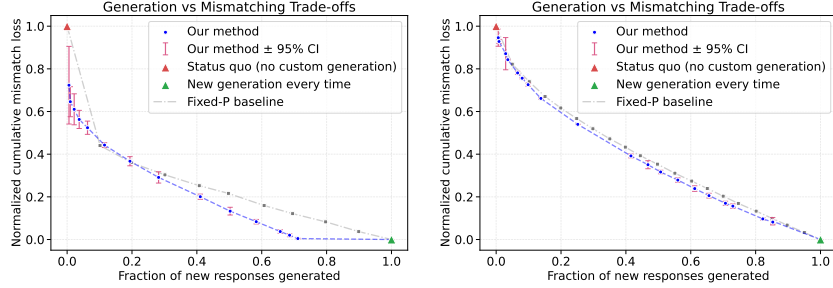
In this section, we conduct numerical experiments to validate our method’s performance. We first run the original algorithm on low-dimensional synthetic data to demonstrate the regret dependence on T . Then we adapt our algorithm to real-world healthcare Q&A scenarios and show better tradeoffs between generation cost and mismatching loss compared to baselines.

D.1 Regret Validation on Synthetic Data

We evaluate our doubly-optimistic algorithm on synthetic data with dimensions $d = 2, 3, 4$ over time horizon $T = 10,000$, repeated for $N = 10$ epochs. Context vectors x_t are drawn from L_2 -normalized uniform distributions, with noise $N_t \sim \mathcal{N}(0, 0.05)$. We calculate regret by comparing the algorithmic loss against OPT_o (defined in Eq. (6)), approximated by randomized K-means++ with Lloyd iterations over potentially optimal values of K . We do not apply OPT_h as its computational cost is exponentially dependent on T .

Figure 1 presents the regret curves in log-log scale to reveal the power dependence of regret on T . Our method exhibits empirical slopes of 0.447, 0.576, 0.676 for $d = 2, 3, 4$ respectively, aligning closely with the theoretical rates which should be $\frac{d}{d+2}$ according to Theorem C.4. These results validate our theoretical analysis in synthetic environments.

Note: We restrict experiments to low-dimensional settings due to the computational cost of OPT_o (a necessary component of regret) in high dimensions, where K-means++ becomes ineffective and the underlying nearest neighbor problem is NP-hard. Despite these computational limitations, the synthetic validation confirms that our approach achieves the predicted theoretical regret rates, providing confidence in its performance for moderate-dimensional real-world applications.



(a) Numerical results on a maternal health Q&A dataset from Nivi . Inc. (b) Numerical results on the public Medical Q&A Dataset.

Figure 2: Tradeoffs between normalized generation costs (x-axis) and normalized mismatching loss (y-axis) on two healthcare Q&A datasets. A lower/left curve indicates a better performance. Each blue point represents the (generation cost, mismatching loss) pair caused by a choice of c . In both cases, our algorithm outperforms the baseline that randomly generates custom responses with a variety of fixed probabilities p (each gray point represents a choice of p).

D.2 Generation-Quality Tradeoffs Analysis on Healthcare Q&A Datasets

We evaluate our algorithm on two real-world healthcare Q&A datasets to demonstrate its practical effectiveness:

1. **Nivi’s Maternal Health Dataset:** A dataset containing 839 user queries, with 12 pre-written FAQs for pregnant women, provided by Nivi.Inc, a company that provides healthcare chatbot services on WhatsApp.
2. **Medical Q&A Dataset:** A public collection of 47,457 medical question-answer pairs curated from 12 NIH websites (<https://www.kaggle.com/datasets/gvaldenebro/cancer-q-and-a-dataset>).

Experimental Setup. Our experimental framework models the create-to-reuse decision process operating entirely in the context representation space. All questions are mapped to embeddings using OpenAI’s pre-trained `text-embedding-3-small` model, creating a semantic representation space where the algorithm makes decisions. For each arriving question context x_t , the algorithm decides whether to select an existing context key f from the FAQ library or add x_t as a new context key and then invoke the custom answer generation oracle $\mathcal{A}(\cdot)$.

Custom answer generation differs across datasets to reflect their nature. For the maternal health dataset, custom answers are generated by GPT-5 with carefully designed prompts including safety guardrails and emergency detection protocols appropriate for healthcare contexts. For the Medical Q&A dataset, custom answers are directly retrieved from the pre-existing responses associated with each question entry.

Crucially, the mismatch loss feedback occurs in the action space rather than the context space. For current question context x_t and an existing context f in the FAQ library, the loss is calculated as $(1 - \text{cosine similarity})$ between x_t ’s **custom answer** and f ’s **custom answer**. This reflects our core assumption that the algorithm operates in context space while true loss manifests in action space, accessible only through the generation oracle $\mathcal{A}(\cdot)$.

As we also mentioned in Section 2, to maintain computational tractability, we model the estimated loss function as $\bar{d}(x, f) = (\theta^\top (x - f))^2$ (on the maternal health dataset) or adopt a neural network $d(x, f; \Theta)$ (on Medical Q&A Dataset).

We evaluate our doubly-optimistic algorithm against a fixed-probability baseline strategy. This baseline makes i.i.d. Bernoulli decisions $\sim \text{Ber}(p)$ at each time step: with probability p , generate a custom response; otherwise, select the most similar existing context from the library based on cosine similarity between question embeddings (note that it has no access to the custom answer before generation). To comprehensively evaluate performance across different cost-accuracy preferences, we vary the probability parameter p uniformly across $[0, 1]$ for the baseline. Meanwhile, we also vary

the creation cost parameter c from 0 to 100 for our algorithm, generating complete tradeoff curves for both approaches.

The numerical results are depicted in Figure 2, where points and curves closer to the bottom-left indicate superior performance. We plot cumulative generation costs against cumulative mismatch losses, with both metrics normalized separately to $[0, 1]$ scale for interpretability. Generation costs are normalized by the total cost of the always-generate strategy, while mismatch losses are normalized by the *status quo* strategy that never generates custom responses. Note that these represent the two components of total loss in our formulation, depicted separately for clearer analysis. Each gray point represents a different choice of p for the baseline, forming a curve that represents the best possible performance achievable by any fixed-probability strategy. Each experiment runs $N = 10$ epochs with 95% confidence intervals computed using Wald’s test.

Results on Nivi’s Maternal Health Dataset. Figure 2a presents the generation-quality tradeoffs. Starting with 12 pre-written FAQs, our algorithm demonstrates several key advantages:

1. **Context Clustering:** Compared with the always-generating strategy (green triangle), approximately 30% of user questions exhibit sufficient similarity to existing FAQs, as evidenced by the algorithm achieving near-zero mismatch loss when generating responses for 70% of queries.
2. **Efficiency Gains:** Compared with *status quo* (red triangle), strategic addition of just a few targeted FAQs reduces mismatch loss by approximately 25% (as evidenced by the algorithmic curve approaching the point $(0, 0.75)$), highlighting the value of adaptive creation decisions over static policies.
3. **Pareto Optimality:** Our algorithm consistently outperforms fixed-probability baselines throughout the entire generation spectrum, with statistical significance demonstrated by 95% confidence intervals. The doubly-optimistic approach effectively pushes the performance frontier toward Pareto optimality.

Results on Medical Q&A Dataset. Figure 2b presents results on the public Medical Q&A dataset. We establish the initial FAQ library by prompting GPT-5 to classify all questions into 32 categories by topic, then randomly sampling 10 question-answer pairs from each category to create a generic response (a total of 32 FAQs).

Compared with the always-generating and FAQ-only baselines respectively, our algorithm can reduce about 60% generation cost and about 60% mismatch loss, leading to positive-sum tradeoffs (indicated by the convex curve). Also, it achieves statistically significant improvements over fixed-probability baselines across nearly the entire generation spectrum, as confirmed by 95% confidence intervals. However, the performance gains are notably smaller than those observed on the private maternal health dataset. We attribute this difference to the greater diversity in the Medical Q&A dataset, spanning 37 question types across 32 medical topics. In contrast, Nivi’s dataset focuses specifically on maternal health with more concentrated topics and frequently recurring keywords, producing clearer semantic connections and stronger correlations between context and action similarities that enable more effective learning.

The results validate our theoretical framework in practice, demonstrating that principled confidence bound approaches for creation decisions significantly outperform heuristic alternatives in real-world healthcare applications where both response quality and resource efficiency are critical.

E Discussions

Dynamic and Context-Dependent Creation Costs. Our current framework assumes a fixed creation cost c across all time steps and contexts. A natural extension would allow time-varying costs c_t or context-dependent costs $c(x_t)$ that reflect realistic scenarios where creation difficulty varies with problem complexity or resource availability. This generalization would better capture applications like drug discovery, where synthesis costs depend on molecular complexity, or content generation, where review costs vary with topic sensitivity. However, this extension introduces significant algorithmic challenges, as evidenced by the substantially worse competitive ratios in variant-cost online facility location problems, where even achieving constant competitive ratios becomes impossible under adversarial sequences.

Non-Parametric and Neural Function Approximation. While our theoretical analysis focuses on parametric quadratic loss functions $d(x, f)$, our empirical experiments demonstrate promising results when replacing the distance function with neural networks and using LLM-as-a-judge for feedback evaluation. Extending the theoretical guarantees to broader function classes, particularly neural networks or kernel methods, would significantly broaden the applicability of our framework. The key challenge lies in controlling the complexity of the function class while maintaining meaningful regret bounds, potentially requiring techniques from neural tangent kernels (NTK) or Bayesian optimization (BO) to handle the high-dimensional hypothesis space.

F Conclusion

In this paper, we introduced an online decision-making problem where new actions can be generated on the fly, at a fixed cost, and then reused indefinitely. To address the balance among exploitation, exploration, and creation, we proposed a doubly-optimistic algorithm that achieves $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}} + d\sqrt{T \log T})$ regret. This regret rate was proved optimal with a matching lower bound, and was validated through simulations. We also implemented our algorithm on a real-world healthcare Q&A dataset to make decisions on generating new answers v.s. applying an FAQ. Our results open up new avenues for optimizing creation decisions in online learning, with potential extensions to broader loss models and flexible creation costs.

G Proof Details

Here we extend the proof sketch of Theorem C.4 provided in Appendix C. According to the roadmap depicted, to validate Theorem C.4, we only need to prove the following Lemmas G.1, G.3 and G.8. We first propose the lemma that bounds OPT_o .

Lemma G.1 (OPT_o upper bound). *We have $OPT_o = O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}})$.*

Proof sketch. We propose a context set (library) \tilde{S} such that $c \cdot |\tilde{S}| + \sum_{t=1}^T \min_{f \in \tilde{S}} d(x_t, f) = O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}})$. Specifically, we let $\tilde{S} := \{[N_1, N_2, \dots, N_d]^\top | N_i \in [\frac{1}{\Delta}], i = 1, 2, \dots, d\}$ as a Δ -covering set over the context space of $[0, 1]^d$. On the one hand, the cumulative mismatch loss due to discretization of the context space is $O(T \cdot \Delta^2 d)$. On the other hand, the total cost of adding new contexts to the set is $O((\frac{1}{\Delta})^d)$. Let $\Delta = T^{-\frac{1}{d+2}} d^{-\frac{1}{d+2}}$ and the total loss is $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}})$. Please kindly find a detailed proof in Appendix G.1. ■

Before getting into the main lemma that upper bounds ALG , we present another lemma showing the concentration of $\bar{d}_t(x_t, f)$ within $\Delta_t(x_t, f)$.

Lemma G.2 (Δ_t as estimation error). *The estimation error of $|\bar{d}_t(x_t, f) - d(x_t, f)|$ is upper bounded by $\Delta_t(x_t, f)$ with high probability. As a consequence, we have $d(x_t, f) - 2\Delta_t(x_t, f) \leq \bar{d}_t(x_t, f) \leq d(x_t, f) \leq \hat{d}_t(x_t, f) \leq d(x_t, f) + 2\Delta_t(x_t, f)$.*

The proof of Lemma G.2 is deferred to Appendix G.2. In the following, we state the lemma that upper bounds the algorithmic loss by a constant competitive ratio over OPT_o adding estimation errors. According to

Lemma G.3 (Constant competitive ratio). *We have $ALG \leq 60OPT_o + 54 \sum_{t=1}^T \Delta_t(x_t, f_t)$.*

Proof. Before starting the proof, we emphasize that all operations we make in this proof are made in the *context* space. As we frequently mention in this paper, the actions are only accessible through the oracle $\mathcal{A}(x)$ for some context x . Therefore, the context library S_t is sometimes referred as a “set” without causing misunderstandings.

First of all, we note that the following two $\{x_t\}_{t=1}^T$ series have identical joint distributions:

- (a) Sample a sequence of x_1, x_2, \dots, x_T independently from an identical distribution \mathbb{D}_X . (iid)
- (b) Sample a set of $Z := \{z_1, z_2, \dots, z_T\}$ independently from the identical distribution \mathbb{D}_X , and then sample $\{x_t\}_{t=1}^T$ as a uniformly random permutation of Z , i.e. $\{x_t\}_{t=1}^T \sim U(\sigma(Z))$. Here $\sigma(Z)$ denotes the permutation set of Z . (iid + permutation)

Given this property, we assume that $\exists Z = \{z_1, z_2, \dots, z_T\}, z_t \stackrel{\text{i.i.d.}}{\sim} \mathbb{D}_X, \{x_t\}_{t=1}^T \sim U(\sigma(Z))$. In the following, we will keep using the notations of $\{x_t\}_{t=1}^T$ and Z accordingly.

Consider the optimal offline solution S^* such that

$$\begin{aligned} OPT_o &= c \cdot |S^*(x_1, x_2, \dots, x_T)| + \sum_{t=1}^T \min_{f \in S^*} d(x_t, f) \\ &= c \cdot |S^*(z_1, z_2, \dots, z_T)| + \sum_{t=1}^T \min_{f \in S^*} d(x_t, f). \end{aligned} \quad (9)$$

Here we denote $S^*(x_1, x_2, \dots, x_T)$ and $S^*(z_1, z_2, \dots, z_T)$ differently to show that the offline solution is not dependent on the permutation, with slight abuse of notation. Denote $S^* =: \{c_1^*, c_2^*, \dots, c_K^*\}$. For each $c_i^*, i = 1, 2, \dots, K$, denote a subset of $\{x_t\}$ as C_i^* such that $\min_{f \in S^*} d(x_t, f) = d(x_t, c_i^*), \forall x_t \in C_i^*$. In other words, C_i^* consists of all x_t 's that are assigned to c_i^* in the optimal solution S^* . Denote $A_i^* := \sum_{t: x_t \in C_i^*} d(x_t, c_i^*)$ as the total optimal cost associated with c_i^* , and $a_i^* := \frac{A_i^*}{|C_i^*|}$ as the average cost in C_i^* .

Now, we define C_i^g and C_i^b as separated GOOD and BAD subsets of C_i^* , respectively, such that

$$\begin{aligned} C_i^g \subset C_i^*, C_i^b \subset C_i^*, |C_i^g| = |C_i^b| = \frac{|C_i^*|}{2} \\ d(x_g, c_i^*) \leq d(x_b, c_i^*), \forall x_g \in C_i^g, x_b \in C_i^b. \end{aligned} \quad (10)$$

In other words, C_i^g and C_i^b represent the nearest half and the farthest half of x_t 's in the set C_i^* , in terms of distance to c_i^* . Note that the sets C_i^g and C_i^b are determined by Z and not relevant to the permutation. Therefore, once Z is realized, the random sequence $\{x_t\}_{t=1}^T$ does not affect C_i^g and C_i^b .

Given these notations, we present and prove the following two lemmas: a Lemma G.4 bounding the *total* loss of GOOD x_t 's, and a Lemma G.6 bounding the *individual* loss of each BAD x_t 's.

Lemma G.4. *The total loss caused by all $x_t \in C_i^g$ is upper bounded as*

$$\sum_{t: x_t \in C_i^g} \mathbb{E}[l_t | \{x_t\}_{t=1}^T] \leq 3c + 4A_i^* + 4 \sum_{x_t \in C_i^g} d(x_t, c_i^*) + 6 \sum_{t=1}^T \Delta_t(x_t, f_t). \quad (11)$$

Proof of Lemma G.4. Denote the context set (library) sequence as $\{S_t\}_{t=1}^T$. Also, denote $\Delta_t := \Delta_t(x_t, f_t)$ and $d_t^* := d(x_t, c_i^*)$ for simplicity. In fact, any $x_t \in C_i^g$ falls in one of the following two cases:

(I) When $\exists e_i \in S_T$ such that $d(e_i, c_i^*) \leq 2a_i^*$, we further categorize x_t into three sub-cases:

I.(a). At time t , we select context e_i and deploy $a_t = \mathcal{A}(e_i)$ (i.e., x_t is matched to context e_i). We have

$$d(x_t, e_i) \leq 2(d(x_t, c_i^*) + d(c_i^*, e_i)) \leq 2(d_t^* + 2a_i^*). \quad (12)$$

The first inequality is due to

$$d(a, b) + d(b, c) \geq \frac{1}{2}d(a, c), \forall a, b, c \in \mathbb{R}^{d^2}. \quad (13)$$

as a quadratic form. Hence

$$\mathbb{E}[l_t | \{x_t\}_{t=1}^T] \leq 2d(x_t, e_i) + 2\Delta_t(x_t, f_t) \leq 4(d_t^* + 2a_i^*) + 2\Delta_t. \quad (14)$$

I.(b). At time t , $e_i \in S_t$ but $a_t \neq \mathcal{A}(e_i)$, i.e. x_t is matched to some other context f_t even with the existence of e_i . Now we have

$$d(x_t, f_t) - 2\Delta_t \leq \check{d}_t(x_t, f_t) \leq \check{d}_t(x_t, e_i) \leq d(x_t, e_i). \quad (15)$$

The second inequality comes from the arg-minimum definition of f_t , and the first and third inequalities is from Lemma G.2. Therefore, we have

$$d(x_t, f_t) \leq d_t(x_t, e_i) + 2\Delta_t \leq 2(d(x_t, c_i^*) + d(c_i^*, e_i)) + 2\Delta_t \leq 2(d_t^* + 2a_i^*) + 2\Delta_t \quad (16)$$

Hence we have

$$\mathbb{E}[l_t | \{x_t\}_{t=1}^T] \leq 2d(x_t, e_i) + 2\Delta_t \leq 4(d_t^* + 2a_i^*) + 6\Delta_t. \quad (17)$$

I.(c). $e_i \notin S_t$ at time t , i.e. x_t is matched to some f_t before any close-enough context e_i being added. In this case, we propose the following lemma that provides an *overall* loss bound for any group of $\{x_t\}$'s, on which no new actions have been created.

Lemma G.5 (Constant loss bound before a new action being generated). *Denote $Q := \{x_{t_i}, i = 1, 2, \dots, n | 1 \leq t_1 \leq \dots \leq t_n \leq T\}$ as a subsequence of $\{x_t\}_{t=1}^T$. Also, denote t_k as the first time in Q such that a new action is generated, i.e. $a_{t_k} = \mathcal{A}(x_{t_k})$ and $a_{t_i} \neq \mathcal{A}(x_{t_i}), i \leq k-1$. We have*

$$\mathbb{E}[\sum_{i=1}^{k-1} l_{t_i} | \{x_t\}_{t=1}^T] \leq c. \quad (18)$$

We defer the proof of Lemma G.5 to Appendix G.3, where we will prove a generalized claim. According to Lemma G.5, the total expected loss for all x_t in this case can be bounded by c .

(II) When $\forall e \in S_T$ satisfies $d(e, c_i^*) > 2a_i^*$, we know that no new action are generated at time t , $\forall t : x_t \in C_i^g$. Then we again apply Lemma G.5 and upper bound the expected total loss by c .

Combining Case I (a,b,c) and Case II, along with a separate cost c of adding e_i , we have an upper bound on the expected total loss for all $t : x_t \in C_i^g$ as follows:

$$\begin{aligned} \mathbb{E}[\sum_{t:x_t \in C_i^g} l_t] &\leq 4 \sum_{t:x_t \in C_i^g} d_t^* + 8 \sum_{t:x_t \in C_i^g} a_i^* + 6 \sum_{t:x_t \in C_i^g} \Delta_t + 3c \\ &= 4 \sum_{t:x_t \in C_i^g} d_t^* + 4A_i^* + 6 \sum_{t:x_t \in C_i^g} \Delta_t + 3c. \end{aligned} \quad (19)$$

Here the last line comes from $|C_i^g| = \frac{|C_i^*|}{2}$. This proves Lemma G.4. \blacksquare

The previous lemma bounds the *total* loss of GOOD x_t 's, while the following lemma will bound the *individual* loss of BAD x_t 's".

Lemma G.6. *For each individual $x_t \in C_i^b$, the expected loss is upper bounded as*

$$\mathbb{E}[l_t | Z] \leq 4d(x_t, c_i^*) + 4\Delta_t(x_t, f_t) + \frac{2}{|C_i^*|} \cdot (c + 8 \sum_{s:x_s \in C_i^g} \mathbb{E}[l_s | Z] + 8 \sum_{s:x_s \in C_i^g} d(x_s, c_i^*)). \quad (20)$$

Proof sketch of Lemma G.6. Intuitively, later-arrived x_t 's should be facing a *better* situation as there are more action candidates. Therefore, for any $x_t \in C_i^b$, if there exists a good point x_g that emerges before the occurrence of x_t , we can upper bound $\mathbb{E}[l_t]$ with $\mathbb{E}[l_g]$ adding $d(x_t, c_i^*)$. This is because we can at least match x_t to the existing in-library context that x_g was matched to. Denote f_g as the existing context whose custom action x_g was assigned to. According to the "triangular inequality" shown as Equation (13) (up to constant coefficient), we have: $\mathbb{E}[l_t] \leq O(d(x_t, f_g)) \leq O(d(x_t, c_i^*) + d(c_i^*, f_g)) \leq O(d(x_t, c_i^*) + d(c_i^*, x_g) + d(x_g, f_g)) = O(d(x_t, c_i^*) + d(x_g, c_i^*) + \mathbb{E}[l_g])$. If there is no such a x_g (with very small probability), then we upper bound the expected loss by c . For a detailed proof of Lemma G.6, please kindly refer to Appendix G.4. \blacksquare

Combining Lemma G.4 and Lemma G.6 above, we have

$$\begin{aligned}
& \mathbb{E} \left[\sum_{t: x_t \in C_i^*} l_t \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[\sum_{t: x_t \in C_i^*} l_t \middle| Z \right] \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[\sum_{s: x_s \in C_i^g} l_s \middle| Z \right] + \mathbb{E} \left[\sum_{r: x_r \in C_i^b} l_r \middle| Z \right] \right] \\
&= \mathbb{E} \left[\mathbb{E} \left[\sum_{t: x_t \in C_i^*} l_t \middle| \{x_t\}_{t=1}^T \right] + \mathbb{E} \left[\sum_{r: x_r \in C_i^b} l_r \middle| Z \right] \right] \\
&\leq \mathbb{E} [3c + 4A_i^* + 4 \sum_{s: x_s \in C_i^g} d(x_s, c_i^*) + 6 \sum_{s: x_s \in C_i^g} \Delta_s(x_s, f_s)] \\
&\quad + \mathbb{E} [4 \sum_{r: x_r \in C_i^b} d(x_r, c_i^*) + 4 \sum_{r: x_r \in C_i^b} \Delta_r(x_r, f_r) \\
&\quad + \frac{|C_i^*|}{2} \cdot \frac{2}{|C_i^*|} \cdot (c + 8 \sum_{s: x_s \in C_i^g} \mathbb{E}[l_s] + 8 \sum_{s: x_s \in C_i^g} d(x_s, c_i^*))] \\
&\leq \mathbb{E} [28c + 40A_i^* + 40 \sum_{s: x_s \in C_i^g} d(x_s, c_i^*) + 54 \sum_{t: x_t \in C_i^*} \Delta_t(x_t, f_t)] \\
&\leq \mathbb{E} [28c + 60A_i^* + 54 \sum_{t: x_t \in C_i^*} \Delta_t(x_t, f_t)].
\end{aligned} \tag{21}$$

Here the last inequality is because $\sum_{s: x_s \in C_i^g} d(x_s, c_i^*) \leq \frac{\sum_{s: x_s \in C_i^g} + \sum_{r: x_r \in C_i^b}}{2} = \frac{A_i^2}{2}$. On the other hand, the sum of losses in OPT_o that are associated to c_i^* equals $c + A_i^*$. Therefore, we have $ALG \leq 60OPT_o + 54 \sum_{t=1}^T \Delta_t(x_t, f_t)$. This ends the proof of Lemma G.3. \blacksquare

Remark G.7. The reason for us to divide $\{x_t\}$'s into GOOD and BAD subsets is twofold.

- (1) We can upper-bound the *total* loss of all GOOD points, mainly because we have Lemma G.5 such that the Case I(c) and Case II hold. Lemma G.5 states that for any group of $\{x_t\}$'s, the expected cost before a new action being created (i.e. before a new context is added to the library) among them is no more than c . Therefore, if there does not exist an e_i close enough to c_i^* , we know that no new actions have been created among GOOD $\{x_t\}$'s (since any GOOD x_t satisfies $d(x_t, c_i^*) \leq 2a_i^*$ and therefore is a qualified candidate e_i once being added to the existing context library). However, this does not hold for BAD points, as they may still trigger new action generations although their contexts are faraway from c_i^* .
- (2) We can only upper-bound the *individual* loss of each BAD x_t due to the reason in (1) above. The individual upper bound for a BAD point is applicable for a GOOD point, but this would introduce a linear dependence on $T \cdot c$ in the overall loss instead of a constant ratio.

Now we propose the lemma where we upper bound the cumulative estimation error.

Lemma G.8 (Linear regression excess risk). *The cumulative absolute error of online linear regression with least-square estimator satisfies $\sum_{t=1}^T \Delta_t = O(\sqrt{d^2 T \log T})$.*

We defer the proof of Lemma G.8 to Appendix G.5 as a standard result from linear regression. According to what we stated earlier, this completes the proof of Theorem C.4.

In the following subsections, we present the proof details of lemmas proposed above.

G.1 Proof of Lemma G.1

Proof. Let $\tilde{S} = \{[N_1, N_2, \dots, N_d]^\top | N_i \in [\frac{1}{\Delta}], i = 1, 2, \dots, d\}$. On the one hand, for any context $x = [x_1, x_2, \dots, x_d]^\top \in \mathbb{R}^d, \|x\|_2 \leq 1$, consider $f_x := [\lfloor \frac{x_1}{\Delta} \rfloor \cdot \Delta, \lfloor \frac{x_2}{\Delta} \rfloor \cdot \Delta, \dots, \lfloor \frac{x_d}{\Delta} \rfloor \cdot \Delta]$.

Due to the definition of \tilde{S} , we know that $f_x \in \tilde{S}$. Also we have $d(x, f_x) = \|x - f_x\|_W^2 < \|[\Delta, \Delta, \dots, \Delta]^\top\|_W^2 \leq \lambda_{\max}(W) \cdot \Delta^2 d$. On the other hand, we have $|\tilde{S}| = (\frac{1}{\Delta})^d$. Denote S^* as the solution to OPT_o (as defined in Eq. (9)), we have

$$\begin{aligned}
OPT_o &= c \cdot |S^*| + \sum_{t=1}^T \min_{f \in S^*} d(x_t, f) \\
&\leq c \cdot |\tilde{S}| + \sum_{t=1}^T \min_{f \in \tilde{S}} d(x_t, f) \\
&\leq c \cdot \left(\frac{1}{\Delta}\right)^d + \sum_{t=1}^T d(x_t, f_{x_t}) \\
&\leq c \left(\frac{1}{\Delta}\right)^d + T \cdot \lambda_{\max}(W) \cdot \Delta^2 d \\
&= O\left(\frac{1}{\Delta^d} + T \Delta^2 d\right),
\end{aligned} \tag{22}$$

and we let $\Delta = T^{-\frac{1}{d+2}} d^{-\frac{1}{d+2}}$ to make the RHS = $O(T^{\frac{d}{d+2}} d^{\frac{d}{d+2}})$. This proves the lemma. \blacksquare

G.2 Proof of Lemma G.2

Proof. Here we prove a more general result on ridge regression:

Lemma G.9. *Let $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ are d -dimension vectors, and $y_i := x_i^\top \theta^* + N_t$, where $\theta^* \in \mathbb{R}^d$ is a fixed unknown vector such that $\|\theta^*\|_2 \leq 1$, and N_t is a martingale difference sequence subject to σ -subGaussian distributions. Denote $X = [x_1, x_2, \dots, x_n]^\top \in \mathbb{R}^{n \times d}$ and $Y = [y_1, y_2, \dots, y_n]^\top \in \mathbb{R}^n$. Let the ridge regression estimator*

$$\hat{\theta} := (X^\top X + I_d)^{-1} X^\top Y$$

where I_d is the $d \times d$ identity matrix. Then with probability $\Pr \geq 1 - \delta$, we have

$$|x^\top (\theta^* - \hat{\theta})| \leq O\left((1 + \sqrt{\log \frac{2}{\delta}}) \sqrt{x^\top (X^\top X + I_d)^{-1} x}\right) \tag{23}$$

holds for any $\delta > 0$ and $x \in \mathbb{R}^d$.

Proof of Lemma G.9. Denote $N := [N_1, N_2, \dots, N_n]^\top \in \mathbb{R}^n$ as the vector of noises in the labels. Then we have

$$\hat{\theta} = (X^\top X + I_d)^{-1} X^\top X \theta^* + (X^\top X + I_d)^{-1} X^\top N. \tag{24}$$

Therefore, the difference between θ^* and $\hat{\theta}$ can be characterized as

$$\begin{aligned}
\theta^* - \hat{\theta} &= \theta^* - (X^\top X + I_d)^{-1} X^\top X \theta^* - (X^\top X + I_d)^{-1} X^\top N \\
&= (X^\top X + I_d)^{-1} (X^\top X + I_d) \theta^* - (X^\top X + I_d)^{-1} X^\top X \theta^* - (X^\top X + I_d)^{-1} X^\top N \\
&= (X^\top X + I_d)^{-1} \theta^* - (X^\top X + I_d)^{-1} X^\top N \\
&= (X^\top X + I_d)^{-1} (\theta^* - X^\top N).
\end{aligned} \tag{25}$$

As a result, we have

$$\begin{aligned}
|x^\top (\theta^* - \hat{\theta})| &= |x^\top (X^\top X + I_d)^{-1} (\theta^* - X^\top N)| \\
&\leq |x^\top (X^\top X + I_d)^{-1} \theta^*| + |x^\top (X^\top X + I_d)^{-1} X^\top N|.
\end{aligned} \tag{26}$$

For the simplicity of notation, denote $A := (X^\top X + I_d)^{-1}$, then we have $|x^\top (\theta^* - \hat{\theta})| \leq \|x^\top A \theta^*\|_2 + \|x^\top A X^\top N\|$. On the one hand, for the first term we have

$$\begin{aligned}
|x^\top A \theta^*| &\leq \|A^\top x\|_2 \cdot \|\theta^*\|_2 \\
&\leq \sqrt{x^\top A A^\top x} \cdot 1 \\
&\leq \sqrt{x^\top A x}.
\end{aligned} \tag{27}$$

The second line is because $\|\theta^*\| \leq 1$, and the last inequality is because $A = A^\top = (X^\top X + I_d)^{-1} \prec I_d$.

On the other hand, for the second term, recall that we set $A := (X^\top X + I_d)^{-1}$ and $\theta^* - \hat{\theta} = A(\theta^* - X^\top N)$. We consider the random variable $x^\top A X^\top N = \sum_{t=1}^n \alpha_t N_t$, where the deterministic coefficients $\alpha_t := (x^\top A X^\top)_t$, $t = 1, \dots, n$.

Notice that $\{N_t\}$ is a martingale difference sequence with subGaussian tails. According to Jin et al. (2019, Proposition 7), which is a subGaussian version of Azuma–Hoeffding’s Inequality, let $d = 1$ and there exists a constant C_J such that

$$\left| \sum_{t=1}^n \alpha_t N_t \right| \leq C_J \cdot \sqrt{\sum_{t=1}^n \alpha_t^2 \log \frac{2}{\delta}}. \quad (28)$$

with probability $\Pr \geq 1 - \delta$. Here $\|\alpha\|_2^2 = \sum_{t=1}^n \alpha_t^2 = x^\top A X^\top X A x \leq x^\top A x$ because $X^\top X \preceq X^\top X + I_d$.

Therefore, with probability at least $1 - \delta$,

$$|x^\top A X^\top N| \leq C_J \sqrt{x^\top A x \log \left(\frac{2}{\delta} \right)}. \quad (29)$$

Returning to $|x^\top(\theta^* - \hat{\theta})| \leq |x^\top A \theta^*| + |x^\top A X^\top N|$, we already established (using $\|\theta^*\|_2 \leq 1$) that $|x^\top A \theta^*| \leq \sqrt{x^\top A x}$. Combining this with Eq. (29) as the martingale tail bound, we get

$$\begin{aligned} |x^\top(\theta^* - \hat{\theta})| &\leq \sqrt{x^\top A x} + C_J \cdot \sqrt{x^\top A x \log \left(\frac{2}{\delta} \right)} \\ &= \left(1 + C_J \cdot \sqrt{\log \left(\frac{2}{\delta} \right)} \right) \sqrt{x^\top (X^\top X + I_d)^{-1} x}. \end{aligned} \quad (30)$$

This ends the proof of Lemma G.9. ■

Now let us go back to the proof of Lemma G.2. We apply this lemma for $6T$ times: in the proof of Lemma G.4 as Case I(a), I(b), I(c) (or Lemma 5.6) and II, and in the proof of Lemma G.6 as Case I and Case II, in each of which we adopt this concentration bound for each existing context $f \in S_t$, which is at most T . Therefore, we let $\delta \leftarrow \frac{1}{6T^2} \delta$ and let $\lambda = 1$, $\alpha = (1 + C_J \cdot \sqrt{\log \frac{12T^2}{\delta}}) \cdot \|W\|_F$. According to Lemma G.9, we prove that $\bar{d}(x_t, f) - \Delta_t(x_t, f) \leq d(x_t, f) \leq \bar{d}(x_t, f) + \Delta_t(x_t, f)$ holds for any $f \in S_t$ and $\forall t = 1, 2, \dots, T$, with probability $\Pr \geq 1 - \delta$. Therefore, we have proved Lemma G.2. ■

G.3 Proof of Lemma G.5

Proof. Notice that at each time t_k , with probability $\Pr = \frac{\hat{d}_{t_k}(x_{t_k}, f_{t_k})}{c}$ we terminate this stochastic process, and with the rest $\Pr = 1 - \frac{\hat{d}_{t_k}(x_{t_k}, f_{t_k})}{c}$ we add $d_{t_k}(x_{t_k}, f_{t_k})$ to our cumulative expected loss. Since $\hat{d}_{t_k}(x_{t_k}, f_{t_k}) \geq d_{t_k}(x_{t_k}, f_{t_k})$, $\forall k \in [n]$, we may instead prove a generalized version of this lemma.

Lemma G.10. *Consider an infinite sequence $\{p_1, p_2, \dots, p_k, \dots\}$ where $p_k \in [0, 1]$. The initial sum $S = 0$. At each time k , with probability p_k we stop this stochastic process, otherwise we add p_k to the sum S . We show that $\mathbb{E}[S] \leq 1$.*

Lemma G.10 is a generalization of Lemma G.5 since we add $d_{t_k}(x_{t_k}, f_{t_k}) \leq \hat{d}_{t_k}(x_{t_k}, f_{t_k})$ at each time k in the latter setting.

Denote a random variable I_k as follows: $I_k = 1$ if the stochastic process has not stopped by the end of time k , and $I_k = 0$ otherwise. In the case when $I_k = 1$, we add p_k to the sum S . Therefore, we have

$$\mathbb{E}[S] = \sum_{k=1}^{\infty} p_k I_k$$

. Also, we know that the probability that $I_k = 1$ is $\Pr[I_k = 1] = \prod_{i=1}^k (1 - p_i)$. As a result, we have

$$\begin{aligned} \mathbb{E}[S] &= \mathbb{E}\left[\sum_{k=1}^{\infty} p_k \cdot I_k\right] \\ &= \sum_{k=1}^{\infty} p_k \prod_{i=1}^k (1 - p_i). \end{aligned} \tag{31}$$

In the following, we show that $\sum_{k=1}^{\infty} p_k \prod_{i=1}^k (1 - p_i) \leq 1$. We first consider $p_k \in (0, 1)$. Denote $Q_0 := 1$ and $Q_k := \Pr[I_k = 1] = \prod_{i=1}^k (1 - p_i)$, and we know $Q_k = (1 - p_k)Q_{k-1} \leq Q_{k-1}$. Also, we have $p_k Q_{k-1} = (1 - (1 - p_k))Q_{k-1} = Q_{k-1} - Q_k$.

For the rigorousness of the proof, we first show that $\sum_{k=1}^{\infty} p_k Q_k$ is finite. Denote

$$T_n := \sum_{k=1}^n p_k Q_k, \tag{32}$$

and we have

$$\begin{aligned} T_n &\leq \sum_{k=1}^n p_k Q_{k-1} \\ &= \sum_{k=1}^n Q_{k-1} - Q_k \\ &= Q_0 - Q_n \\ &< Q_0 = 1 \end{aligned} \tag{33}$$

As $T_n < 1$ and $T_{n+1} \geq T_n, \forall n \geq 1$, we have

$$\lim_{n \rightarrow \infty} T_n \leq 1 \tag{34}$$

according to the Monotone Convergence Theorem. Then we slightly generalize the results above from $p_k \in (0, 1)$ to $p_k \in [0, 1]$, i.e. incorporating 0 and 1. In fact, if $p_k = 0$, then we may skip this $p_k Q_k$ term. Otherwise if $p_k = 1$, consider the first m s.t. $p_m = 1$, and then we still have $E[S] = \sum_{k=1}^{m-1} p_k I_k = T_{m-1} < 1$ and $I_m = I_M = 0$ for any $M \geq m, M \in \mathbb{Z}^+$.

Therefore, we have

$$\begin{aligned} \mathbb{E}[S] &= \sum_{k=1}^{\infty} p_k Q_k \\ &\leq \sum_{k=1}^{\infty} p_k Q_{k-1} \\ &= \sum_{k=1}^{\infty} Q_{k-1} - Q_k \\ &= Q_0 - \lim_{k \rightarrow \infty} Q_k \\ &\leq 1. \end{aligned} \tag{35}$$

This ends the proof of Lemma G.10 and therefore proves Lemma G.5. ■

G.4 Proof of Lemma G.6

Proof. Consider the moment when a $x_t \in C_i^b$ arrives, and denote s as the most recent moment ($s < t$) such that $x_s \in C_i^g$. According to the uniform permutation assumption from Z to $\{x_t\}_{t=1}^T$, this x_s can be any $z \in C_i^g$ with equal probability as $\Pr = \frac{1}{|C_i^g|} = \frac{2}{|C_i^*|}$. In the following, we analyze the expected loss $\mathbb{E}[l_t]$ by two cases:

- (I) If $x_s \in C_i^g$ does exist before x_t occurs. Denote $f_t^* := \operatorname{argmin}_{f \in S_t} d(c_i^*, f)$ as the closest context to c_i^* existed by the time t . Then we have:

$$\begin{aligned} \mathbb{E}[l_t | \{x_t\}_{t=1}^T] &= c \cdot \frac{\hat{d}_t(x_t, f_t)}{c} + d(x_t, f_t) \cdot (1 - \frac{\hat{d}_t(x_t, f_t)}{c}) \\ &\leq \hat{d}_t(x_t, f_t) + d(x_t, f_t) \\ &\leq \check{d}_t(x_t, f_t) + 2\Delta_t + \check{d}_t(x_t, f_t) + 2\Delta_t \\ &\leq 2\check{d}_t(x_t, f_t^*) + 4\Delta_t \\ &\leq 2d(x_t, f_t^*) + 4\Delta_t \\ &\leq 4d(x_t, c_i^*) + 4d(c_i^*, f_t^*) + 4\Delta_t. \end{aligned} \quad (36)$$

Also, denote $\hat{f}_s := \operatorname{argmin}_{f \in S_s} d(x_s, f)$ as the best existing context that can be matched to x_s by the time s . Then we know that

$$\begin{aligned} \mathbb{E}[l_s | \{x_t\}_{t=1}^T] &\geq d(x_s, \hat{f}_s) \\ &\geq \frac{1}{2}d(c_i^*, \hat{f}_s) - d(x_s, c_i^*) \\ &\geq \frac{1}{2}d(c_i^*, f_s^*) - d(x_s, c_i^*). \end{aligned} \quad (37)$$

Combining Eq. (36) with Eq. (37), we have

$$\begin{aligned} \mathbb{E}[l_t | \{x_t\}_{t=1}^T] &\leq 4d(x_t, c_i^*) + 4\Delta_t + 4 \cdot 2(\mathbb{E}[l_s | \{x_t\}_{t=1}^T] + d(x_s, c_i^*)) \\ &\leq 4d_t^* + 2\Delta_t + 8 \cdot \frac{2}{|C_i^*|} \cdot \left(\sum_{s: x_s \in C_i^g} \mathbb{E}[l_s | \{x_t\}_{t=1}^T] + d_s^* \right). \end{aligned} \quad (38)$$

Again, the last line of Eq. (38) comes from the i.i.d. assumption of x_s .

- (II) If $x_s \in C_i^g$ does not exist before x_t occurs, i.e. $x_r \in C_i^b, \forall r \leq t-1$. According to the uniform permutation from Z to $\{x_t\}_{t=1}^T$, this event happens with probability $\frac{2}{|C_i^*|}$. In this case, if $\hat{d}_t(x_t, f_t) \geq c$, then we suffer a cost c at time t ; otherwise $\hat{d}_t(x_t, f_t) < c$, and we either generate a new action (with cost c) or suffer an expected loss at $d(x_t, f_t) \leq \hat{d}_t(x_t, f_t) < c$. In a nutshell, the expected loss does not exceed c .

Combining with Case I and Case II above, we immediately get Eq. (20). ■

G.5 Proof of Lemma G.8

Proof. Denote $\Delta_t := \Delta_t(x_t, f_t)$ for simplicity. In the following, we first reduce the summation of estimation error Δ_t to the regret of a $K(\leq T)$ -arm linear bandit problem, up to constant factors. In fact, according to Lemma G.2, we know that $d(x_t, f) \leq \check{d}_t(x_t, f) \leq d(x_t, f)$. Since we select $f_t = \operatorname{argmin}_{f \in S_t} \check{d}_t(x_t, f)$, we have

$$\begin{aligned} d(x_t, f_t^*) - 2\Delta_t &\leq d(x_t, f_t) - 2\Delta_t \\ &\leq \check{d}_t(x_t, f_t) \\ &\leq \check{d}_t(x_t, f_t^*) \\ &\leq d(x_t, f_t^*), \end{aligned} \quad (39)$$

where $f_t^* := \operatorname{argmin}_{f \in S_t} d(x_t, f)$ as the best existing context for x in the current context library at time t . Therefore, the performance gap between f_t and f_t^* can be bounded as $d(x_t, f_t) - d(x_t, f_t^*) \leq$

$2\Delta_t$. On the other hand, since $d(x_t, f_t) = \langle w, \phi(x_t, f_t) \rangle$ is a linear loss function, we consider $\phi(x_t, f)$ as the “context”³ of each arm $f \in S_t$, and then we form a linear contextual bandit problem setting. Recall that $\Delta_t = \alpha \cdot \sqrt{\phi(x_t, f_t)^\top \Sigma_{t-1}^{-1} \phi(x_t, f_t)}$. According to [Chu et al. \(2011\)](#) Lemma 3 (which originates from [Auer \(2002\)](#) Lemma 13), we have

$$\begin{aligned} \sum_{t=1}^T \Delta_t &= \sum_{t=1}^T \alpha \cdot \sqrt{\phi(x_t, f_t)^\top \Sigma_{t-1}^{-1} \phi(x_t, f_t)} \\ &\leq \alpha \cdot 5 \sqrt{(d^2) |\Psi_{T+1}| \log |\Psi_{T+1}|} \\ &\leq 5\alpha \sqrt{d^2 T \log T}. \end{aligned} \tag{40}$$

Here the second line is because the dimension of contexts are d^2 as $\phi(x_t, f) = \text{Vec}((x_t - f)(x_t - f)^\top) \in \mathbb{R}^{d^2}$, and the third line comes from the original definition of Ψ_t as a subset of $[t-1]$. ■

G.6 Proof of Lower Bound (Theorem C.5)

Proof. In order to prove the lower bound, we show the following facts

1. $OPT_o = \Omega(T^{\frac{d}{d+2}})$ according to the K -nearest-neighbors(K-NN) lower bound.
2. Any online facility location algorithm suffers at least $(2 - o(1))$ -competitive-ratio, i.e. $ALG \geq (2 - o(1))OPT_h$.

In the following, we present two lemmas corresponding to the facts above.

Lemma G.11 (OPT_o lower bound). *We have $OPT_h \geq OPT_o \geq \Omega(T^{\frac{d}{d+2}})$.*

Proof. Denote

$$\begin{aligned} OPT_o(K) &:= \min_{S: |S|=K} c \cdot |S| + \sum_{t=1}^T \min_{f \in S} d(x_t, f) \\ &= K + T \cdot \min_{S: |S|=K} \frac{1}{T} \sum_{t=1}^T \min_{f \in S} d(x_t, f). \end{aligned} \tag{41}$$

This equals T times K -nearest-neighbors (K-NN) loss plus K . According to [Zador \(1964\)](#) (i.e. Zador’s Theorem in coding theory), the mean squared distance to the nearest codebook center in \mathbb{R}^d space in L_r -norm is lower bounded by $\Omega(K^{-\frac{r}{d}})$. This is directly applicable to K-NN which effectively partitions points by their nearest neighbors. Hence, the quantization lower bound established by Zador’s Theorem translates into a lower bound on K-NN’s average squared loss. Therefore, we let $r = 2$ to fit in our setting, and then have

$$OPT_o = \min_{K \in [T]} OPT_o(K) = \Omega(c \cdot K + T \cdot K^{-\frac{2}{d}}) = \Omega(T^{\frac{d}{d+2}}), \tag{42}$$

where the last line is an application of Hölder’s Inequality that $K + T \cdot K^{-\frac{2}{d}} \geq K^{\frac{\frac{2}{d}}{1+\frac{2}{d}}} (T \cdot K^{-\frac{2}{d}})^{\frac{1}{1+\frac{2}{d}}} = T^{\frac{d}{d+2}}$, and the equality holds at $K = T^{\frac{d}{d+2}}$. ■

Lemma G.12 (Theorem 5.1 in [Kaplan et al. \(2023\)](#)). *Let \mathcal{A} be an algorithm for online facility location in the i.i.d. model, then, the competitive ratio of \mathcal{A} is at least $2 - o(1)$.*

Combining Lemma G.11 and Lemma G.12, we know that $REG = ALG - OPT_h \geq (2 - o(1) - 1)OPT_h \geq 0.5OPT_o = \Omega(T^{\frac{d}{d+2}})$. This proves Theorem C.5 ■

³Here we denote this covariate as the *context* as it serves as an environmental description in the contextual bandits. We denote $f \in S_t$, which was denoted as a context in the library, an *arm* of this contextual bandits.