
Contextual Bandits for Large-Scale Structured Discrete Constrained Optimization Problems

Pavithra Harsha^{1*} Chitra Subramanian¹ Naoki Abe¹ Shivaram Subramanian¹
Amadou Ba¹ Kevin Arturo Fernández Román² Mauricio Longinos Garrido²
Miao Liu¹ Aurélie Lozano¹ Chandrasekhar Narayanaswami¹
¹IBM Research ²IBM CIO *Contact Author: pharsha@us.ibm.com

Abstract

1 We study contextual bandits in high-dimensional combinatorial action spaces arising
2 in structured constrained optimization problems, such as IT resource allocation
3 and retail assortment pricing. Key quantities in the objective or constraints must
4 be estimated from data during the trial sequence of actions. In [12], we propose
5 a novel, practical, and transparent approach based on general-purpose regression
6 oracles with Inverse Gap Weighting (IGW) for seamless integration within an
7 optimization framework. IGW sampling is efficiently managed by: (a) a column-
8 generation reformulation of the underlying Mixed Integer Programming (MIP)
9 model which allows for flexible lower-level predictors, causal coherence, and efficient
10 representation of large action spaces; (b) a diverse solution pool generation to
11 balance the exploration-exploitation trade-off in large-action spaces. To address
12 non-smooth rewards induced by constraints, we introduce a risk-averse phased
13 learning strategy. Experiments on an IT auto-scaling task demonstrate substantial
14 reductions in cumulative regret, with added gains from risk-averse methods that
15 effectively manage constraint violations. This submission summarizes [12] and
16 sketches our extensions towards a full theoretical regret analysis.

17 1 Introduction

18 In a contextual bandit setting a learning/decision-making agent, repeatedly receives a context, selects
19 an action, based on past observations, and receives a reward. The learner’s goal is to maximize
20 the cumulative reward of an employed policy, striking the right balance between exploration and
21 exploitation. In [12] we propose the design of a practical contextual bandit algorithm for scenarios
22 with a high-dimensional action space involving a structured discrete constrained optimization prob-
23 lem. These problems frequently arise in the operations management domain such as IT resource
24 allocation [16, 4] and retail assortment pricing [7, 13]. Rather than treating the whole action space as
25 an opaque entity admitting a global reward regression function, as is done in past works on contextual
26 bandits, here we take the *transparent white box* approach in which we explicitly leverage the rich
27 structure that is inherent in the reward optimization problem.

28 First of all, in many real-world settings, the overall reward can be naturally formulated as a function
29 in terms of lower-level estimable functionals which capture the causal structure necessary to represent
30 the action-reward relationship. These lower-level functional relationships can be modeled effectively
31 using machine learning methods, by leveraging the inherent problem structure that restricts the causal
32 dependencies to smaller dimensions. The response variables of these lower-level models are predicted
33 for incoming contexts and used in a combinatorial white-box optimization formulation. These
34 formulations are generally more tractable than their black-box counterparts due to their structured
35 nature. Leveraging a white-box approach allows us to ensure that the underlying lower-level machine
36 learning models are causally adequate, namely, confounders are controlled for and other causal

37 constraints can be enforced (e.g., demand decreases with price, or latency with resources). This also
 38 improves transparency and interpretability and supports predictor reuse across use cases.

39 The second key characteristic of these applications is the presence of operational constraints such as
 40 latency targets, sales goals, or capacity limits. While such constraints can be added as penalty terms
 41 in generic reward learning frameworks, it undermines the reward smoothness along the constraint
 42 boundaries, challenging much of the guarantees on contextual bandit methods. It makes it difficult to
 43 satisfy requirements or attain high rewards due to causal inadequacies and constraint violations.

44 Our novel solution to this problem setting leverages cutting edge advances in two fields: contextual
 45 bandits and discrete constrained optimization. On the bandits side, recent work have given rise to the
 46 development of general, near-optimal algorithms that are based on ‘regression oracles’ [8, 9, 17, 5].
 47 A key advantage of this type of approach in our context is its ability to accommodate complex, non-
 48 standard reward models, including those we aim to consider, in which lower-level response models are
 49 composed to derive the overall reward function. The core method that is relevant to our work is the
 50 probabilistic action selection strategy referred to as ‘Inverse Gap Weighting’ (IGW) [9, 2, 1]. Existing
 51 extensions for large action spaces either assume linear reward embeddings [18] or reward/action-
 52 space smoothness [19] or hierarchical action structure [15] and are not applicable to our setting. We
 53 build and extend the notion of ‘diversity’ in the IGW sampling introduced by [18]. On the discrete
 54 constrained optimization side, we leverage the advanced technique of Column Generation (CG)
 55 in our context and develop a master-subproblem reformulation of the underlying problem which
 56 is tractably solvable by mixed-integer programming (MIP) for many practical applications. This
 57 CG approach accommodates flexible, general-purpose predictors as input, ensures causal coherence
 58 and exploits the inherent lower-dimensional problem structure for a more efficient representation of
 59 the large action space. Our extensions to effectively handle optimization constraints highlight the
 60 dichotomy between extensive early exploration in traditional bandits and more conservative, risk
 61 averse exploration under constraints, with limited violations until boundaries are better understood.

62 2 Problem Setup

63 We study the following discrete structured constrained optimization problem:

$$P(x) : \min_{\mathbf{a} \in \mathcal{A}} \sum_{i \in \mathcal{I}} c_i(x, \mathbf{a}_i), \quad \text{subject to } L_i(x, \mathbf{a}_i) \leq 0 \quad \forall i \in \mathcal{I}. \quad (1)$$

64 Let $x \in \mathcal{X}$ and $\mathbf{a} \in \mathcal{A}$ refer to the context and action vector in their respective spaces. The
 65 action set \mathcal{A} restricts each action dimension $j \in J$ to a discrete set \mathcal{A}^j , so $\mathbf{a} = (a^j)_{j \in J}$ where
 66 $a^j \in \mathcal{A}^j$. The objective components c_i and constraint terms L_i for $i \in \mathcal{I}$ are the context- and
 67 action-dependent response variables, but each depends only on a lower-dimensional partial/sub action
 68 vector $\mathbf{a}_i = (a^j)_{j \in J_i} \in \mathcal{A}_i$ where $J_i \subset J$ and $|J_i| \ll |J|$, typically at most 3-4. These partial
 69 action vectors may overlap across i , preventing trivial decomposition. The constraints L_i capture
 70 performance relative to a user-specified threshold. As the L_i ’s are learned over time, we also use a
 71 soft-penalized formulation with user-defined $\lambda_i \forall i \in \mathcal{I}$:

$$P^\lambda(x) : \min_{\mathbf{a} \in \mathcal{A}} Z^\lambda(x, \mathbf{a}), \quad \text{where } Z^\lambda(x, \mathbf{a}) = \sum_{i \in \mathcal{I}} c_i(x, \mathbf{a}_i) + \lambda_i L_i(x, \mathbf{a}_i)^+. \quad (2)$$

72 A contextual bandit protocol of the decision problem proceeds as follows in T rounds. In each
 73 round $t \in [T]$, the learner receives a context $x_t \in \mathcal{X}$, selects an action $\mathbf{a}_t \in \mathcal{A}$ and observes the
 74 response variables $\{c'_{ti}(\mathbf{a}_{ti}), L'_{ti}(\mathbf{a}_{ti})\} \forall i \in \mathcal{I}$ and based on that realizes penalized cost $Z_t^\lambda(\mathbf{a}_t) =$
 75 $\sum_{i \in \mathcal{I}} c'_{ti}(\mathbf{a}_{ti}) + \lambda_i L'_{ti}(\mathbf{a}_{ti})^+$. We assume that in each round t , although the contexts are chosen
 76 arbitrarily, the response variables conditioned on the context are sampled from fixed but unknown
 77 distributions $\mathbb{P}_{c_i}(\cdot|x)$ and $\mathbb{P}_{L_i}(\cdot|x) \forall i \in \mathcal{I}$, respectively. We assume the learner has access to a class
 78 of regression functions or predictors $\mathcal{F}_{c_i}, \mathcal{F}_{L_i}^+$ for each $i \in \mathcal{I}$ with bounded outputs, where $L_i^+ =$
 79 $\max\{L_i, 0\}$ denotes the constraint violation. We work with the predictors for constraint violation,
 80 L^+ , rather than L , in order to recover the true expected cost, $E[Z_t^\lambda(\mathbf{a}_{ti})] = \sum_{i \in \mathcal{I}} E[c'_{ti}(\mathbf{a}_i)] +$
 81 $\lambda_i E[L'_{ti}(\mathbf{a}_{ti})^+]$. We make the following realizability assumption [3, 8, 9, 17], where the predictors
 82 α_i^* aim to approximate the *true cost and constraint violation functions*, respectively.

83 **Assumption 2.1 (Realizability):** *There exists a regression function $\alpha_i^* \in \mathcal{F}_{\alpha_i}$ such that*
 84 *$E[\alpha'_{ti}(\mathbf{a}_i)|x_t = x] = \alpha_i^*(x, \mathbf{a}_i)$ for all $\mathbf{a}_i \in \mathcal{A}_i$ and $t \in [T]$ where $\alpha_i \in \{c_i, L_i^+\}$, $\alpha'_{ti} \in$*
 85 *$\{c'_{ti}, L'_{ti}^+\} \forall i \in \mathcal{I}$.*

86 **Regret:** Given the set of predictors, $f = \{\hat{\alpha}_i \in \mathcal{F}_{\alpha_i} | \alpha_i \in \{c_i, L_i^+\} \ i \in \mathcal{I}\}$, we define the
 87 induced policy as $\pi_f(x) := \arg \min_{\mathbf{a} \in \mathcal{A}} Z_f^\lambda(x, \mathbf{a})$, where Z_f^λ is obtained from Eq.(2) by substituting
 88 predictors f . The aim of the learner is to minimize the regret with respect to the optimal policy
 89 $\pi^* := \pi_{f^*}$ defined as: $\text{Reg}_{\text{CB}} = \sum_{t=1}^T Z_t^\lambda(a_t) - Z_t^\lambda(\pi^*(x_t))$.

90 3 Proposed Solution

91 Our proposed solution integrates several key ideas described briefly below (refer to [12] for details).

92 **Generating lower-level predictors:** We assume access to regression oracles for lower-level
 93 predictors $\alpha_i \in \{c_i, L_i^+\} \ \forall i \in \mathcal{I}$ that are trained on sequence of context–action–response data.
 94 Using only the lower-dimensional action vector as input to the regression oracles, we effectively
 95 incorporate the inherent problem structure. The predictors can be general parametric and non-
 96 parametric models, such as neural networks, regression trees, random forests, and kernels, optionally
 97 with causal constraints, to generalize effectively across contexts and actions.

98 **Optimizing MIP via column generation (CG):** We employ CG, a decomposition technique compat-
 99 ible with any state-of-art regression model for lower-level functions, to efficiently solve the discrete
 100 optimization problem to near-optimality. Our CG approach exploits the sparsity in the correlations
 101 between decisions to decompose the model into a high-dimensional ‘restricted’ master program
 102 (RMP) and several low-dimensional subproblems. CG iterates between these levels till convergence
 103 and produces an RMP that can be solved directly as a MIP using off-the-shelf optimization packages
 104 such as CPLEX and Gurobi to obtain not just an optimal solution but a range of feasible solutions.

105 **Generating a diverse solution pool:** MIP solvers can generate and store multiple solutions beyond
 106 the optimal one in a configurable solution pool while traversing the combinatorial search tree. Users
 107 can specify the pool capacity, optimality gap (relative or absolute), and a diversity-based replacement
 108 strategy to maintain varied solution pool, all without affecting solver performance (see [6, 11]).

109 **IGW sampling:** The original sampling procedure targeted small, fixed action spaces, with regret
 110 scaling linearly in action-space size, making it unsuitable for large spaces. For large spaces [18]
 111 introduced reward-function approximations and diverse action pool-based IGW sampling. We devise
 112 a heuristic IGW method that samples only from the diverse solution pool, improving the exploration
 113 and exploitation balance in large spaces. The underlying equation are given in the proposed algorithm.

114 **Proposed algorithm - OptCB:** We now introduce an initial version of our proposed algorithm,
 115 referred to as *OptCB* (Optimization-based Contextual Bandit), with the details outlined in Algorithm
 116 1. It seamlessly integrates key components: (a) general offline regression oracles for learning (steps
 117 4,5,10), (b) MIP optimization with CG at the core and built-in diverse pool generation (step 6), and (c)
 118 IGW sampling (steps 7 and 8). Its practical design enables easy integration into standard workflows
 119 that already combine ML and discrete optimization.

120 **Risk-averse learning with constraints:** The exploration–exploitation principle, supported by
 121 theory [17, 18], suggests extensive early exploration followed by gradual reduction. Under con-
 122 straints, however, violations heavily impact the penalized objective $Z^\lambda(x, \mathbf{a})$, and risk considerations
 123 discourage excessive early exploration due to cost discontinuities. While some boundary violations
 124 are necessary for learning, large violations are prohibitively costly and best avoided. This motivates
 125 a risk-averse strategy: prioritize exploration within the “known” boundary early on, then gradually
 126 expand around its vicinity as it is better understood. Algorithm 2 illustrates a simple two-phase
 127 variant (see [12] for another variant), where the transition is triggered by a significant change in the
 128 gradient of the observed cumulative squared error series $Z_t^\lambda(x, \mathbf{a})$, a proxy for regression regret.

129 4 Empirical Validation

130 We validate our approach in synthetic environments using real-world IT resource management data.
 131 Our method reduces cumulative regret by 80% over the vanilla IGW baseline in small action spaces
 132 (256 actions) and 90% in large spaces (10,000 actions), with further gains from risk-averse strategies
 133 (8–10% and 57–65%, respectively). See plots in appendix. At scale (100M actions), we achieve
 134 1000× speedups (35 min/iteration vs. 2 sec). Statistical significance tests and ablation studies confirm
 135 that each component of our design is critical to maximizing cumulative rewards.

Algorithm 1 OptCB

1: **Input:**
Epoch schedule $0 = \tau_0 < \tau_1 < \tau_2 < \dots < \tau_{T'} = T$.
MIP constraints target penalty vector: λ .
MIP pool parameters (gap, size, capacity): g, s, l .
IGW sampling parameters: γ_0, ρ .

2: **for** epoch $m = 0, 1, 2, \dots, T' - 1$ **do**
3: **for** round $t = \tau_m + 1, \dots, \tau_{m+1}$ **do**
4: Observe context x_t .
5: Receive predictions \hat{c}_i, \hat{L}_i^+ from regression oracles.
6: Solve the MIP reformulation $P^\lambda(x_t)$ with pool parameters (g, s, l) to obtain a diverse solution pool S .
7: Set IGW sampling probability $p_t(\cdot)$: (i) Evaluate penalized objective $Z^\lambda(\mathbf{a}) \forall \mathbf{a} \in S$. (ii) Define $S^* = \operatorname{argmin}_{\mathbf{a} \in S} Z^\lambda(\mathbf{a})$. (iii) Let Z^{λ^*} be the corresponding optimal objective. (iv) Set $\gamma_t = \gamma_0 \tau_m^\rho$ (see [12] for alternative approaches). (v) Determine sampling probability as follows:
$$p_t(\mathbf{a}) = \frac{1}{|S| + \gamma_t (Z^\lambda(\mathbf{a}) - Z^{\lambda^*})} \forall \mathbf{a} \in S \setminus S^*, \text{ and } p_t(\mathbf{a}) = \frac{[1 - \sum_{\mathbf{a} \in S \setminus S^*} p_t(\mathbf{a})]}{|S^*|} \forall \mathbf{a} \in S^*,$$

8: Sample $\mathbf{a}_t \sim p_t(\cdot)$, execute \mathbf{a}_t and observe resultant response variables $c_{it}, L_{it} \forall i \in I$
9: **end for**
10: Re-train regression predictors \hat{c}_i, \hat{L}_i^+ with all data including that from epoch m using offline least squares.
11: **end for**

Algorithm 2 Two-phase RiskAverse OptCB

Assumption: Monotonic Predictors $L_i(x, \cdot)$ *Same as OptCB barring following changes*

1: **Input (additional):**
Initialize power-law exponent: $\theta_0 > 1$.
Target phase transition power-law exponent: $\bar{\theta}$.
Additional IGW sampling parameter: $\bar{\rho} > \rho$.

2: In step 6, solve reformulation $P(x_t)$ with hard constraints if $\theta_{\tau_m} > \bar{\theta}$ and its soft counterpart $P^\lambda(x_t)$ otherwise. To ensure hard problem is not infeasible, if for any $i \in \mathcal{I}$, $\min L_i(x_t, \cdot) > 0$, we modify $L_i(x_t, \cdot)$ to $L_i(x_t, \cdot) - \min L_i(x_t, \cdot)$ to obtain a solution in round t .

3: In step 7(iv), set $\gamma_t = \gamma_0 \tau_m^\rho$ if $\theta_{\tau_m} > \bar{\theta}$, else set $\gamma_t = \gamma_0 \tau_m^{\bar{\rho}}$ (see [12] for alternative approaches).

4: After step 8, compute regression regret $\operatorname{Reg}_{\text{sq}}(\cdot)$ or its proxy, cumulative squared error of $Z^\lambda(x, \mathbf{a})$.

5: After step 10, perform log-log regression on all the computed $\operatorname{Reg}_{\text{sq}}(\cdot)$ values to obtain the power-law coefficient $\theta_{\tau_{m+1}}$.

136 **5 Regret Analysis**

137 IGW sampling, combined with an estimation–decision meta-algorithm, unifies several prior regret
138 guarantees in contextual bandits through the decision-to-estimation (DEC) coefficient [10]. Bounded
139 DEC, together with bounded squared-loss regression regret under the realizability assumption, yields
140 optimal regret guarantees in contextual information settings with finite, linear and certain continuous
141 action settings [9, 10]. However, existing works for bounding DEC, including that in [18], do not
142 extend naturally to the discrete, large-scale combinatorial settings considered here.

143 We introduce a refined DEC notion that captures the interplay between the decomposed estimation
144 error and its impact on decision quality in our formulation. Recall that our approach leverages the
145 MIP solver’s ability to generate diverse near-optimal solutions during branch-and-bound search while
146 maintaining performance and ensuring the solutions are within a user-specified optimality gap. The
147 IGW sampling scheme over the diverse solution pool allows us to visit the ϵ -independent dimensions
148 in finite time t_ϵ , with the number of such dimensions upper bounded by the Eluder dimension [14].
149 Using the optimality gap together with the Eluder dimension, we upper bound the refined DEC, in the
150 penalized reward setting, to obtain \sqrt{T} regret when T exceeds t_ϵ . On the other hand when $T < t_\epsilon$ the
151 regret is controlled by the radius of the ball covering the action space leveraging the diverse solution
152 pool. Extending this theoretical treatment to capture the trade-off between early exploration and
153 the risk of large constraint violations, and to develop alternative risk-aware estimation-to-decision
154 algorithms, remains an open and promising direction for future research.

155 **A Appendix: Regret Plots**

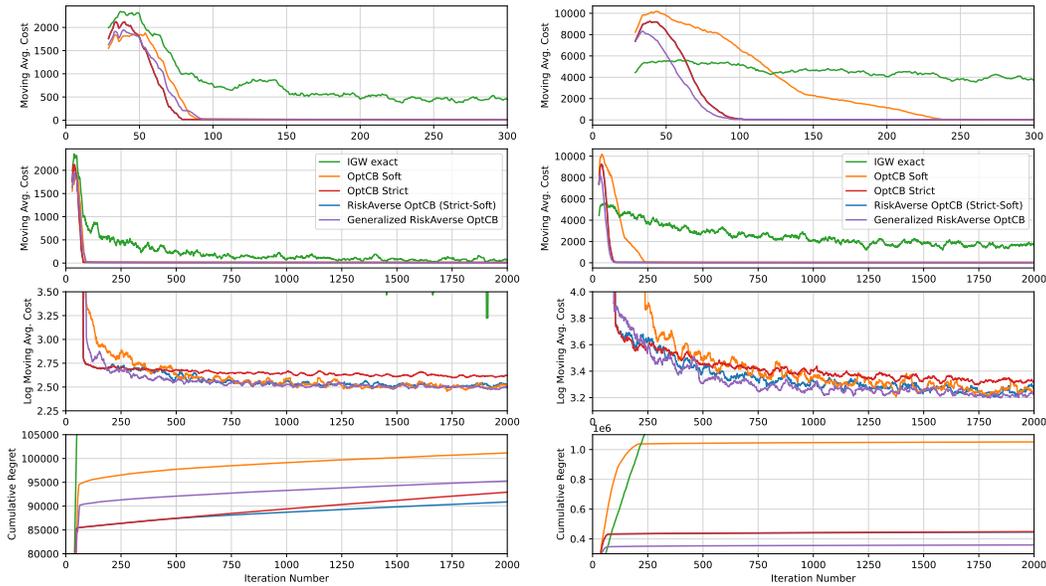


Figure 1: Instantaneous cost and cumulative regret metrics for: (left plot) small-scale environment with 256 combinatorial actions; and (right plot) large-scale environment with 10K combinatorial actions.

156 **References**

157 [1] Naoki Abe, Alan W Biermann, and Philip M Long. Reinforcement learning with immediate
 158 rewards and linear hypotheses. *Algorithmica*, 37:263–293, 2003.

159 [2] Naoki Abe and Philip M Long. Associative reinforcement learning using linear probabilistic
 160 concepts. In *ICML*, pages 3–11. Citeseer, 1999.

161 [3] Alekh Agarwal, Miroslav Dudík, Satyen Kale, John Langford, and Robert Schapire. Contextual
 162 bandit learning with predictable rewards. In *Artificial Intelligence and Statistics*, pages 19–26.
 163 PMLR, 2012.

164 [4] Romil Bhardwaj, Kirthevasan Kandasamy, Asim Biswal, Wenshuo Guo, Benjamin Hindman,
 165 Joseph Gonzalez, Michael Jordan, and Ion Stoica. Cilantro: Performance-Aware resource
 166 allocation for general objectives via online feedback. In *17th USENIX Symposium on Operating
 167 Systems Design and Implementation (OSDI 23)*, pages 623–643, 2023.

168 [5] Alberto Bietti, Alekh Agarwal, and John Langford. A contextual bandit bake-off. *Journal of
 169 Machine Learning Research*, 22(133):1–49, 2021.

170 [6] CPLEX. Cplex solution pool. <https://www.ibm.com/docs/en/icos/22.1.2?topic=solutions-what-is-solution-pool>.
 171

172 [7] Kris Johnson Ferreira, Bin Hong Alex Lee, and David Simchi-Levi. Analytics for an online
 173 retailer: Demand forecasting and price optimization. *Manufacturing & service operations
 174 management*, 18(1):69–88, 2016.

175 [8] Dylan Foster, Alekh Agarwal, Miroslav Dudík, Haipeng Luo, and Robert Schapire. Practical
 176 contextual bandits with regression oracles. In *International Conference on Machine Learning*,
 177 pages 1539–1548. PMLR, 2018.

- 178 [9] Dylan Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits
179 with regression oracles. In *International Conference on Machine Learning*, pages 3199–3210.
180 PMLR, 2020.
- 181 [10] Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity
182 of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- 183 [11] Gurobi. Gurobi solution pool. [https://docs.gurobi.com/projects/optimizer/en/
184 current/features/solutionpool.html#solution-pool](https://docs.gurobi.com/projects/optimizer/en/current/features/solutionpool.html#solution-pool).
- 185 [12] Pavithra Harsha, Chitra Subramanian, Naoki Abe, Shivaram Subramanian, Amadou
186 Ba, Kevin Arturo Fernández Román, Mauricio Longinos Garrido, and Chandrasekhar
187 Narayanaswami. Practical contextual bandits for large-scale structured discrete constrained
188 optimization problems. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge
189 Discovery and Data Mining V.2*, KDD '25, page 850–859, 2025.
- 190 [13] Pavithra Harsha, Shivaram Subramanian, and Markus Ettl. A practical price optimization
191 approach for omnichannel retailing. *INFORMS Journal on Optimization*, 1(3):241–264, 2019.
- 192 [14] Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic
193 exploration. *Advances in Neural Information Processing Systems*, 26, 2013.
- 194 [15] Rajat Sen, Alexander Rakhlin, Lexing Ying, Rahul Kidambi, Dean Foster, Daniel N Hill, and
195 Inderjit S Dhillon. Top-k extreme contextual bandits with arm hierarchy. In *International
196 Conference on Machine Learning*, pages 9422–9433. PMLR, 2021.
- 197 [16] Jiuchen Shi, Hang Zhang, Zhixin Tong, Quan Chen, Kaihua Fu, and Minyi Guo. Nodens:
198 Enabling resource efficient and fast {QoS} recovery of dynamic microservice applications in
199 datacenters. In *2023 USENIX Annual Technical Conference (USENIX ATC 23)*, pages 403–417,
200 2023.
- 201 [17] David Simchi-Levi and Yunzong Xu. Bypassing the monster: A faster and simpler optimal
202 algorithm for contextual bandits under realizability. *Mathematics of Operations Research*,
203 47(3):1904–1931, 2022.
- 204 [18] Yinglun Zhu, Dylan J Foster, John Langford, and Paul Mineiro. Contextual bandits with
205 large action spaces: Made practical. In *International Conference on Machine Learning*, pages
206 27428–27453. PMLR, 2022.
- 207 [19] Yinglun Zhu and Paul Mineiro. Contextual bandits with smooth regret: Efficient learning in
208 continuous action spaces. In *International Conference on Machine Learning*, pages 27574–
209 27590. PMLR, 2022.