

Bandits with Preference Feedback: A Stackelberg Game Perspective

Barna Pasztor*, Parnian Kassraie*, Andreas Krause

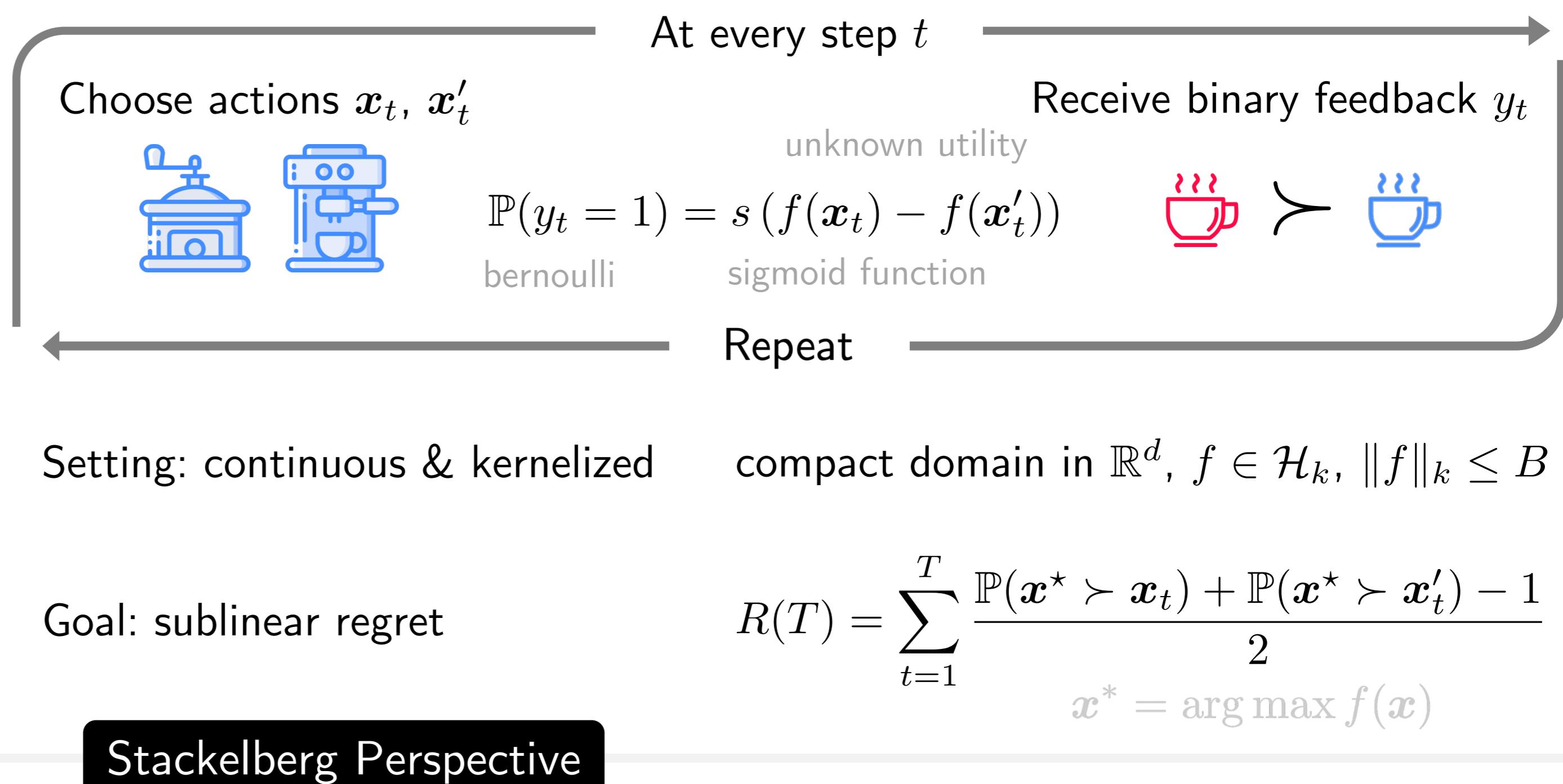


Challenges

- Continuous action space
- Qualitative preference feedback
- Costly sampling
- Complexity of exploration & exploitation

Contributions

- Stackelberg Game formulation
- Practical confidence bounds for kernelized utilities
- No-regret guarantee
- Very promising performance



View actions as players in a Stackelberg Game

- With objective $\mathbb{P}(\mathbf{x} \succ \mathbf{x}')$, both players choose \mathbf{x}^* via backward induction
- True preference is unknown
- Approximate it with a lower-bound

MaxMinLCB Acquisition Function

$$\begin{aligned} \mathbf{x}_t &= \arg \max_{\mathbf{x}} \text{LCB}_t(\mathbf{x} \succ \omega(\mathbf{x})) \\ \text{s.t. } \omega(\mathbf{x}) &= \arg \min_{\mathbf{x}'} \text{LCB}_t(\mathbf{x} \succ \mathbf{x}') \\ \mathbf{x}'_t &= \omega(\mathbf{x}_t) \end{aligned}$$

Leader

Follower

Organically balances exploration & exploitation

- What's the role of the Leader?
- What's the role of the Follower?

Theorem (Regret - Informal)

With an appropriate choice of β_t , MaxMinLCB satisfies

$$\mathbb{P}\left(\forall T \geq 1 : R(T) \leq C_1 (\gamma_T + \log 1/\delta) \sqrt{T}\right) \geq 1 - \delta$$

Preference-based inference is *equivalent* to learning with direct feedback, up to choice of kernel.

$$\tilde{k}((x_1, x'_1), (x_2, x'_2)) = k(x_1, x_2) + k(x'_1, x'_2) - k(x_1, x'_2) - k(x'_1, x_2)$$

$$h_t \leftarrow \arg \min \frac{\lambda}{2} \|h\|_{\tilde{k}}^2 + \sum_{\tau=1}^t -y_\tau \log [s(h(\mathbf{x}_\tau, \mathbf{x}'_\tau))] - (1 - y_\tau) \log [1 - s(h(\mathbf{x}_\tau, \mathbf{x}'_\tau))]$$

To construct a lower-bound for

- $h_t(\mathbf{x}, \mathbf{x}')$ estimates the utility gap $f(\mathbf{x}) - f(\mathbf{x}')$
- $\sigma_t(\mathbf{x}, \mathbf{x}')$ quantifies the estimation uncertainty

$$\text{LCB}_t(\mathbf{x}, \mathbf{x}') = s(h_t(\mathbf{x}, \mathbf{x}')) - \beta_t \sigma_t(\mathbf{x}, \mathbf{x}')$$

Theorem (Anytime Preference-based Conf Seq)

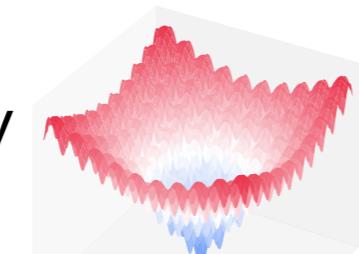
Choosing $\beta_t \asymp \gamma_t + \log(1/\delta)$ satisfies

$$\forall t \geq 1, \mathbf{x}, \mathbf{x}' \in \mathcal{X} : |\mathbb{P}(\mathbf{x} \succ \mathbf{x}') - s(h_t(\mathbf{x}, \mathbf{x}'))| \leq \beta_t \sigma_t(\mathbf{x}, \mathbf{x}')$$

with probability greater than $1 - \delta$.

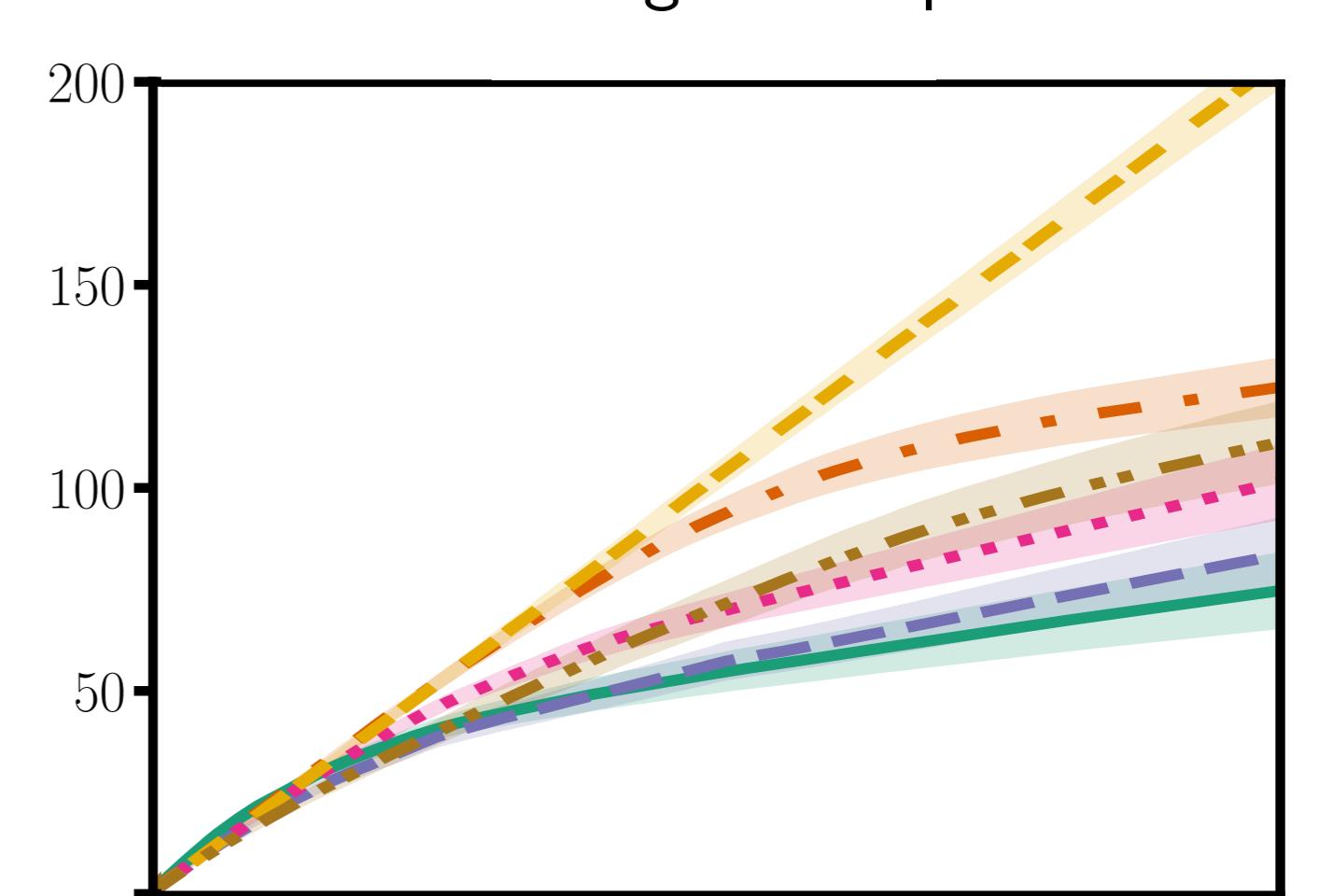
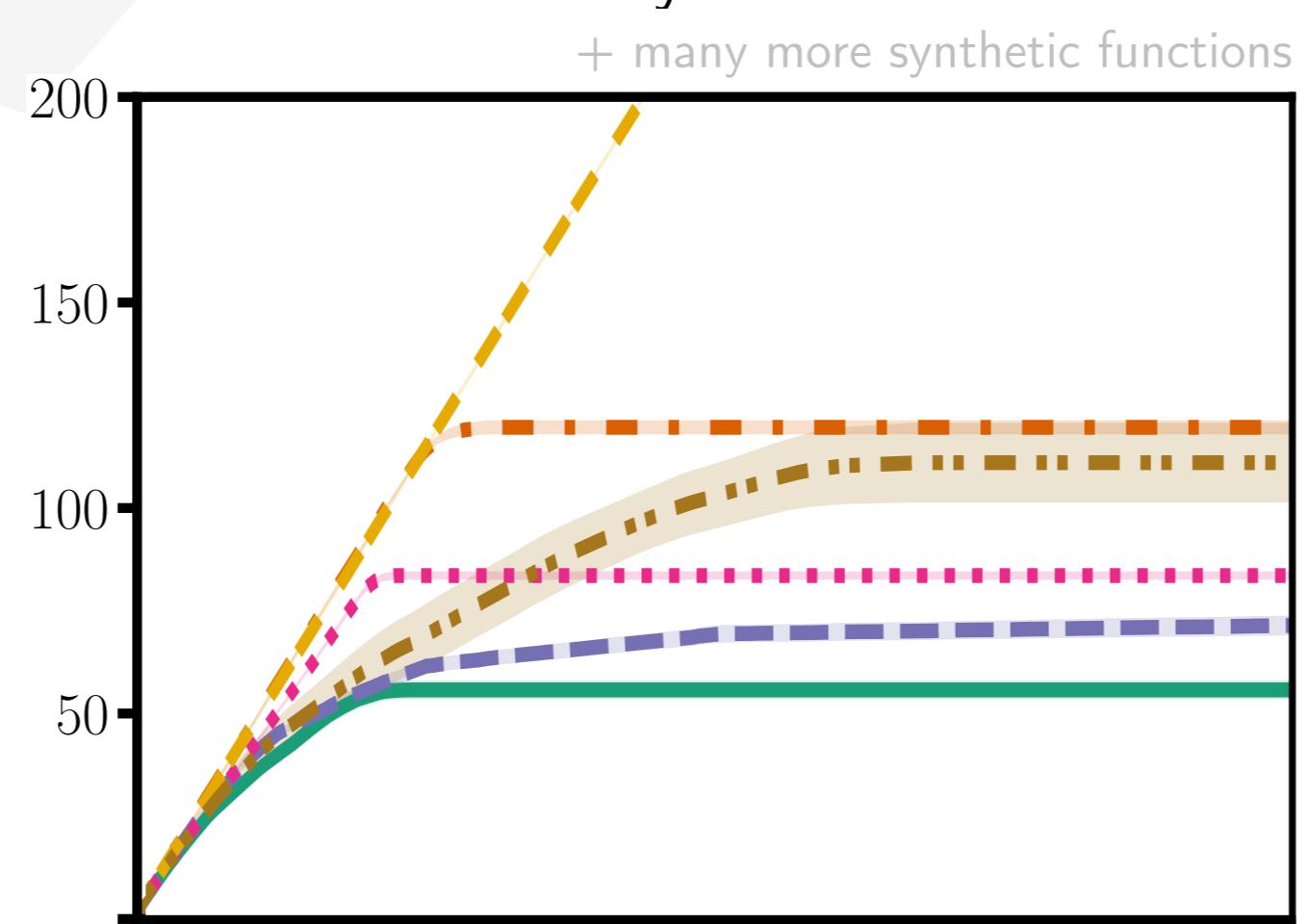
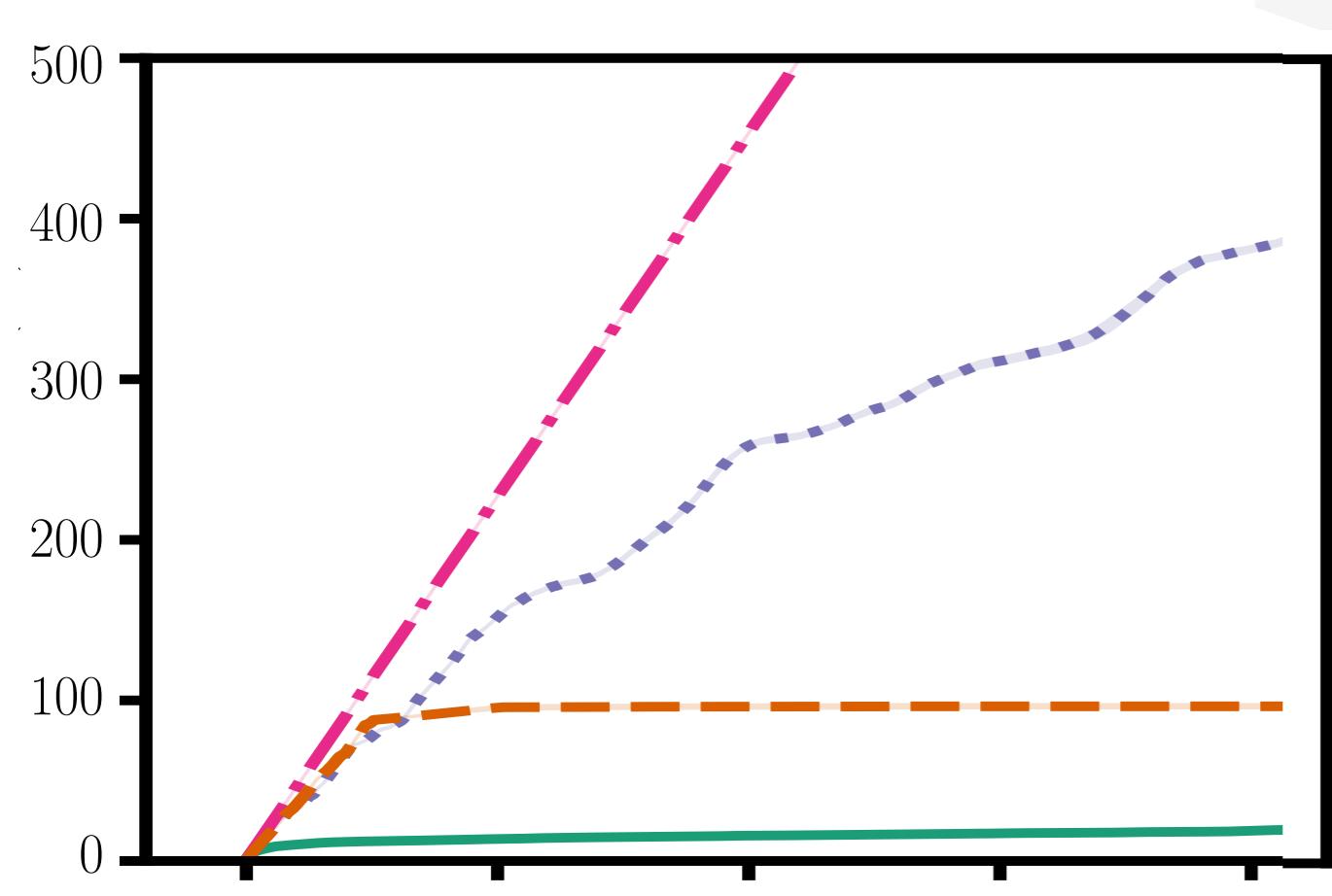
Result

regret of logistic bandit on Ackley reward using different conf. seqs.



pref-based bandit benchmark on the 2D Ackley function

pref-based bandit benchmark on text embeddings of Yelp reviews



Up Next

Applications in RLHF
adaptive fine-tuning of LLMs to niche domains, personalized & pluralist usage

Learning w Finite Recall
choosing an action from recent history to improve costs & feedback quality

Welfare Maximization
accepting feedback from multiple sources and aggregating the preference