

LARGE-SCALE MDPs WITH KL CONTROL COST AND ITS APPLICATION TO CROWDSOURCING

YASIN ABBASI-YADKORI, PETER L. BARTLETT, XI CHEN, ALAN MALEK

MOTIVATION

- Markov decision process: modeling sequential decisions
- E.g. queueing network, robot planning
- Dynamic Programming can solve for small state problems
- Applications can have *large* state spaces
- Here: in the KL-cost setting, can efficiently do large state spaces

NOTATION

An MDP is defined by:

- State space $\mathcal{X} = \{1, \dots, X\}$
- Action space \mathcal{A}
- Transition Kernel $K : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_{\mathcal{X}}$
- Loss function $\ell : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$

The problem is to:

- Policy $\pi : \mathcal{A} \rightarrow \Delta_{\mathcal{A}}$
- Find policy to minimize

$$J_{\pi}(x) = \mathbb{E} \left[\sum_{t=0}^{\infty} \ell(X_t, \pi) | X_0 = x \right]$$

Aim for optimality *within a restricted family of policies*.

EXTENDING TO LARGE STATE SPACES

- Parametric class of policies π_{θ} for $\theta \in \Theta$ with losses J_{θ}
- Bellman operator:

$$(\mathbf{L}J_{\theta})(x) = \min_{a \in \mathcal{A}} (\ell(x, a) + \mathbb{E}_{K(x, a)}[J_{\theta}(x') | x])$$

- Optimal policy has $\mathbf{L}J_{\theta} = J_{\theta}$.
- Linear Programming formulation:

$$\begin{aligned} \min_{\theta} c^{\top} J_{\theta}, & \quad (\text{low cost}) \\ \text{s.t. } \mathbf{L}J_{\theta} \leq J_{\theta} & \quad (J_{\theta} \text{ is feasible}) \end{aligned}$$

- Look for efficient relaxations, e.g.

$$\min_{\theta} c^{\top} J_{\theta} + \|\mathbf{L}J_{\theta} - J_{\theta}\|$$

- Previous ADP techniques used $J_{\theta} = \Psi\theta$

LINEARLY SOLVABLE MDPs FROM [TODOROV]

- $\mathcal{A} = \Delta_{\mathcal{X}}$
- Loss: $\ell(x, P(\cdot|x)) = q(x) + D_{KL}(P(\cdot|x) || P_0(\cdot|x))$
 - state loss $q(x)$, base dynamics P_0
 - infinite loss unless $P \ll P_0$
- Terminal state z : $q(z) = 0$ and $P_0(z|z) = 1$
- Total cost of policy P

$$J_P(x) = \mathbb{E} \left[\sum_{t=0}^{\infty} \ell(X_t, P) | X_0 = x \right]$$

- Greedy action is:

$$P_J(\cdot|x) = \arg \min_{p \in \Delta_{\mathcal{X}}} J_p(x) \propto \frac{1}{Z(x)} P_0(x'|x) e^{-J_P(x')}$$

- Bellman's operator becomes linear in $g(x) = e^{-J(x)}$:

$$e^{\mathbf{L}J(x)} = e^{q(x)} \sum_{x'} P_0(x, x') e^{-J(x')}$$

LARGE STATE SPACES FOR LMDPs

- Intuition: take $J_{\theta} = -\log(\Psi\theta)$ so $e^{\mathbf{L}J_{\theta}}$ is linear in θ
- Approximate unconstrained optimization:

$$\min_{\theta} c^{\top} J_{\theta} + H \|\mathbf{L}J_{\theta} - J_{\theta}\|$$

- $\|\mathbf{L}J_{\theta} - J_{\theta}\|$ not convex in θ , but

$$e^{-\max\{\mathbf{L}J_{\theta}, J_{\theta}\}} \|\mathbf{L}J_{\theta} - J_{\theta}\| \leq \|e^{-\mathbf{L}J_{\theta}} - e^{-J_{\theta}}\|$$

- Optimization relaxed to:

$$\min_{\theta} -c^{\top} \log(\Psi\theta) + H \|e^{-q} P_0 \Psi\theta - \Psi\theta\|$$

OUR ALGORITHM FOR TOTAL COST

Input: x_1, N, H , step sizes $(\eta_t), v$.
Initialize $\theta_1 = \mathbf{0}$.
for $t := 1, 2, \dots, N$ **do**
 Sample trajectory $(x_1, a_1, \dots, z) \sim v$.
 Compute the stochastic subgradient r_t .
 Update $\theta_{t+1} = \Pi_{\mathcal{W}}(\theta_t - \eta_t r_t)$.
end for
 $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \theta_t$.
Return policy $P_{J_{\hat{\theta}_T}}$

PERFORMANCE BOUND

Theorem 1. Choose $H \geq e^{\|q\| + \log\|1/\Psi\|}$. Let $\hat{\theta}$ be an ϵ -optimal solution. Then, for any $\theta \in \Theta$ with $l_{\theta} = \min(J_{\theta}, \mathbf{L}J_{\theta})$,

$$\begin{aligned} J_{P_{J_{\hat{\theta}}}}(x_1) \leq & \inf_{J_{\theta} \in \mathcal{J}} \{J_{P_{J_{\theta}}}(x_1) + \mathcal{E}(J_{\theta})\} + \epsilon \\ & + \|P_{J_{\hat{\theta}}} - Q\|_1 \max_{T \in \mathcal{T}} \sum_{x \in T} |J_{\hat{\theta}}(x) - \mathbf{L}J_{\hat{\theta}}(x)| \end{aligned}$$

The penalty function

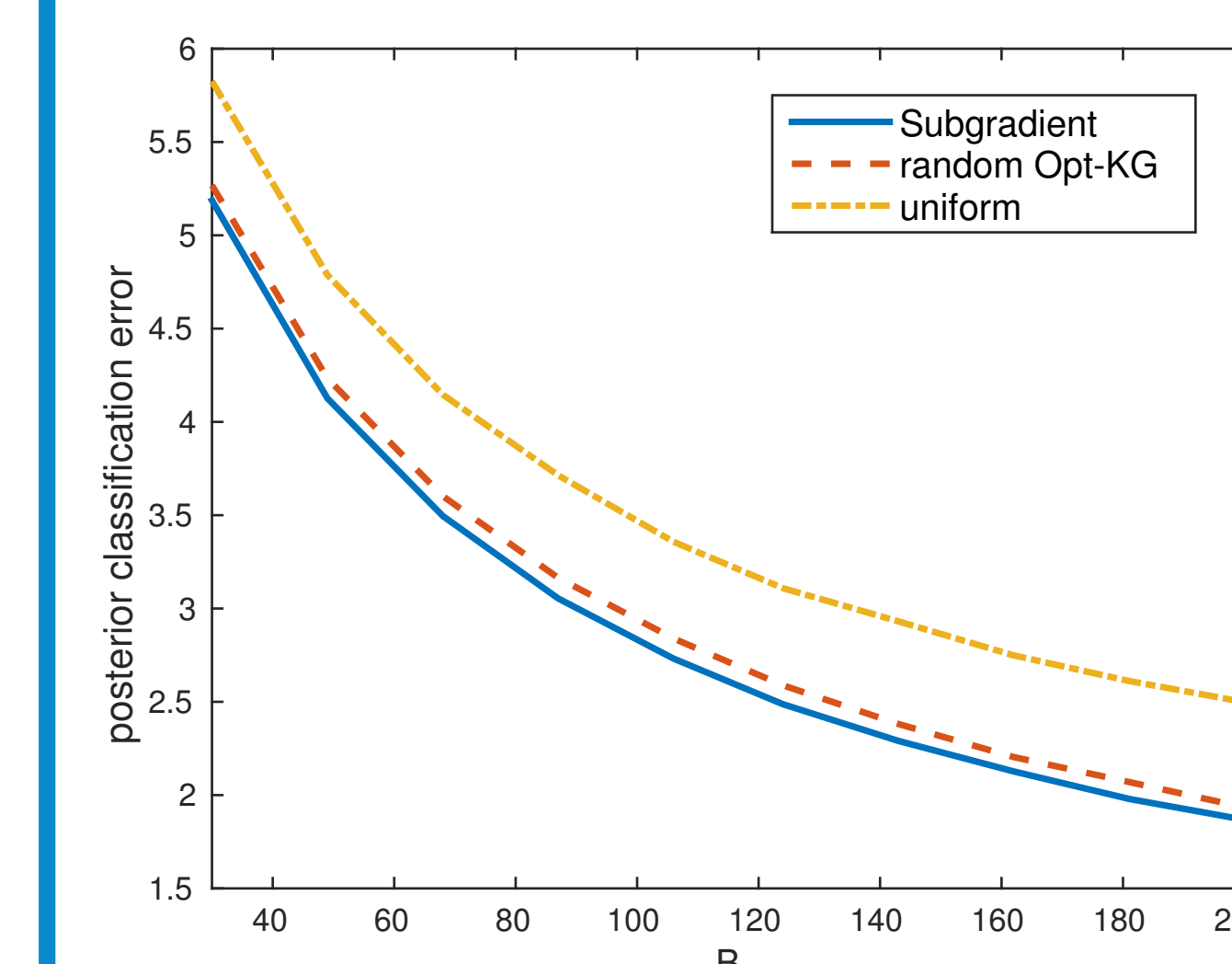
$$\mathcal{E}(J_{\theta}) = \sum_{T \in \mathcal{T}} \sum_{x \in T} (HQ(T)e^{-l_{\theta}(x)} + P_{J_{\theta}}(T)) |J_{\theta}(x) - \mathbf{L}J_{\theta}(x)|$$

is related to how far J_{θ} is from the optimal value function.

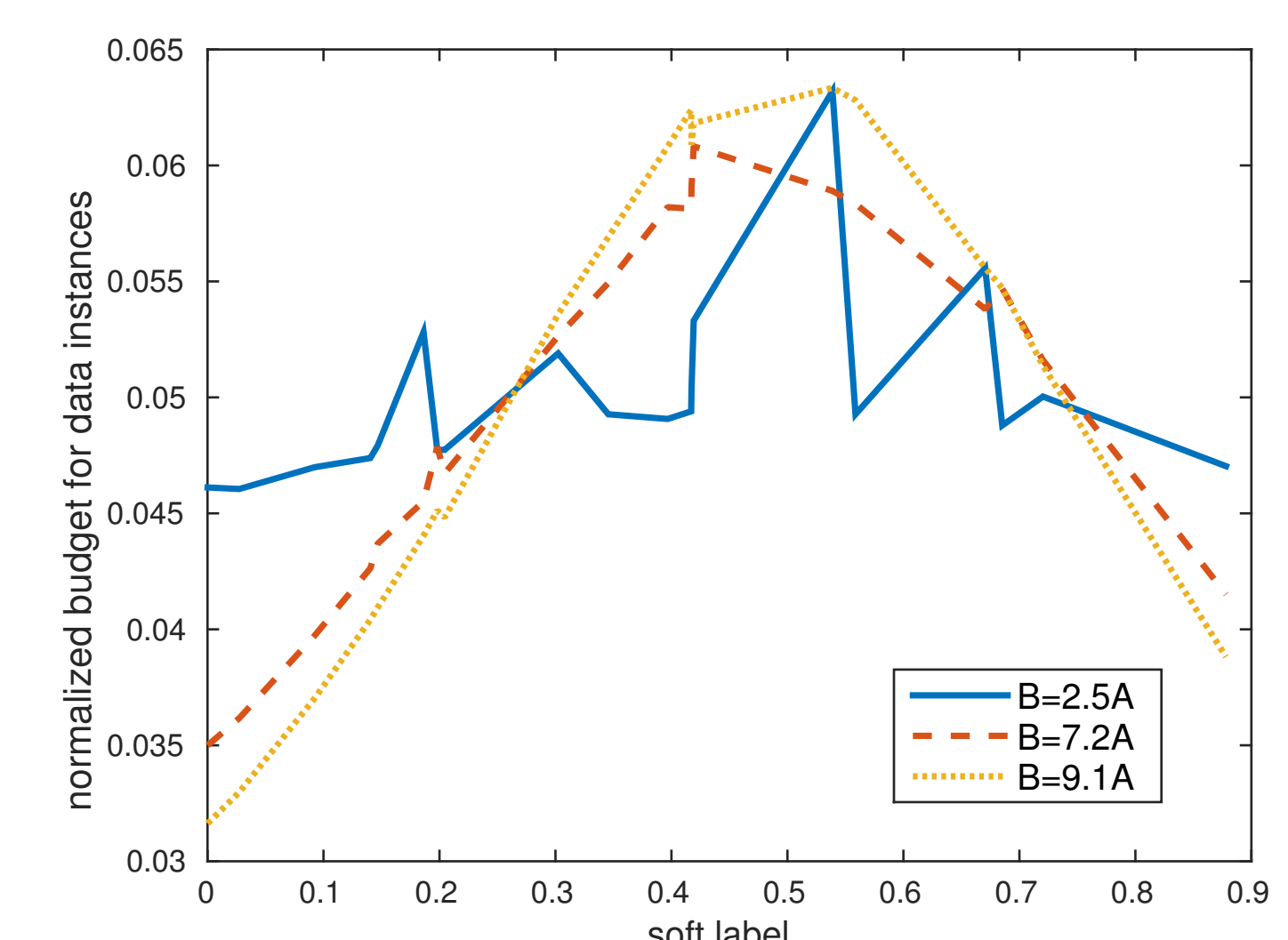
CROWDSOURCING

- Need to label A items.
- Each item has soft label $\mu_i \in [0, 1]$
- Guess if $\{\mu_i \geq \frac{1}{2}\}$ for as many i as we can
- For $t = 1, \dots, T$:
 - Pick $i \in \{1, \dots, A\}$
 - Receive $X_t \sim \text{Bern}(\mu_i)$
- Use Beta prior \Rightarrow MDP dynamics equivalent to Bayesian updates
- P_0 limits transitions,
- $q(x)$ rewards correct labels

EXPERIMENTAL RESULTS



Average error of three policies. Our method requires 10% fewer samples for same accuracy



Portion of budget vs. soft label. Harder soft labels receive more budget, and the difference grows with B .