

# Planning in Large-Scale Markov Decision Problems

Yasin Abbasi-Yadkori <sup>1</sup>, Peter Bartlett <sup>1,2</sup>, Xi Chen <sup>3</sup>, **Alan Malek**<sup>2</sup>

<sup>1</sup>Queensland University of Technology

<sup>2</sup>University of California, Berkeley

<sup>3</sup>NYU Stern School of Business

June 27th, 2016

# Map

- Problem: planning in an MDP with large state space
- Goal: find near-optimal policy in low dimensional family
- Average Cost
  - ▶ Parameterize dual LP
  - ▶ Obtain “agnostic” guarantee
  - ▶ Queueing network
- KL-cost
  - ▶ Exploit Linearly Solvable MDPs
  - ▶ Parameterize log of loss function
  - ▶ Crowdsourcing

# Motivation

- Markov decision process: modeling sequential decisions
- Decouple learning and planning, e.g. [?]
- E.g. queueing network, robot planning
- Can solve for small state spaces
- Large state spaces: “curse of dimensionality”

# Outline

- 1 Motivation
- 2 Linearly Solvable MDPs
- 3 Extending to large dimensions
- 4 Experiments

# MDPs

A Markov Decision Process is specified by:

- State space  $\mathcal{X} = \{1, \dots, X\}$
- Action space  $\mathcal{A}$
- Transition Kernel  $P : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_{\mathcal{X}}$
- Loss function  $\ell : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}^+$

Planning problem:

- Find policy  $\pi : \mathcal{X} \rightarrow \Delta_{\mathcal{A}}$  to minimize value function

$$J_{\pi}^{\gamma}(x) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \ell(\mathbf{X}_t, \pi) \mid \mathbf{X}_0 = x \right] \quad (\text{discounted cost})$$

$$J_{\pi}(x) = \lim_{n \rightarrow \infty} \mathbb{E} \left[ \frac{1}{n} \sum_{t=0}^n \ell(\mathbf{X}_t, \pi) \mid \mathbf{X}_0 = x \right] \quad (\text{average cost})$$

$$J_{\pi}(x) = \mathbb{E} \left[ \sum_{t=0}^{\infty} \ell(\mathbf{X}_t, \pi) \mid \mathbf{X}_0 = x \right] \quad (\text{total cost})$$

# Discounted Cost

- Earliest attempt to solve planning problem [?]
- Define the *Bellman Operator*  $L^\gamma$

$$(L^\gamma J)(x) = \min_a \ell(x, a) + \gamma \mathbb{E}[J(X_1) | X_0 = x, A_0 = a]$$

- $L^\gamma$  is an  $L_\infty$  contraction:  $\|J - J'\|_\infty \leq \gamma \|L^\gamma J - L^\gamma J'\|_\infty$
- If  $\gamma < 1$ , there is a unique solution  $J^* = \lim_{k \rightarrow \infty} L^{\gamma k} J$
- $J$  is optimal iff  $L^\gamma J = J$
- Optimal policy is greedy:

$$\pi^*(a|x) = \mathbb{I}\{a = \arg \min_a \ell(x, a) + \gamma \mathbb{E}[J^*(X_1) | X_0 = x, A_0 = a]\}$$

- Unfortunately, Bellman iteration is  $O(X^2A)$

## Average Cost

- More complicated:  $L^1$  not a contraction
- Need to measure w.r.t. the average cost,  $\lambda \in \mathbb{R}$  and rely on Markov Chain stationarity
- Define the differential cost function  $h \in \mathbb{R}^X$  and Bellman operator

$$Lh(x) := \min_{a \in \mathcal{A}} \left[ \ell(x, a) + \sum_y P(y|x, a)h(y) \right]$$

- Bellman optimality:  $Lh = h + \lambda \mathbf{1}$
- [?] Thm. 8.4.1: Suppose  $\lambda$  and  $h$  satisfy  $Lh \geq h + \lambda \mathbf{1}$ . Then  $\lambda \leq \lambda^*$ .
- Motivates exact average-cost LP [?]

$$\begin{aligned} & \max_{\lambda, h} \lambda, \\ & \text{s.t. } h + \lambda \mathbf{1} \leq Lh \end{aligned}$$

- Always has a solution [Thm. 8.4.3] for recurrent chains

# Linear Programming Formulation

- Define  $B \in \mathbb{R}^{XA \times X}$  by  $B_{(x,a),y} = \{x = y\}$ . We can write:

$$h + \lambda \mathbf{1} \leq Lh = \min_a \left[ \ell(x, a) + \sum_y P(y|x, a)h(y) \right]$$

$\Leftrightarrow$

$$h(x) + \lambda \leq \ell(x, a) + \sum_y P(y|x, a)h(y) \quad \forall x, a$$

$\Leftrightarrow$

$$B(\lambda \mathbf{1} + h) \leq \ell + Ph$$

- Average-cost LP equivalent too

$$\max_{\lambda, h} \lambda,$$

$$\text{s.t. } B(\lambda \mathbf{1} + h) \leq \ell + Ph$$

- Dimension  $X$ , number of constraints  $O(XA)$ . Intractable!



# Large state space

- Parametric class of value functions  $J_\theta$  or  $h_\theta$  for  $\theta \in \Theta \subset \mathbb{R}^d$
- For any value function  $J$  or  $h$ , there is a greedy policy  $\pi_J$  or  $\pi_h$  (the argmax in  $L^\gamma$ )

## Problem (Large-State Planning Problem)

Assume:

- $\mathcal{X}$  is very large
- We have entrywise access to  $P$  and  $\ell$
- Goal: find  $\theta$  to minimize

$$J_{\pi_{J_\theta}}$$

Greedy policy using value function  $J_\theta$

Value function of running this greedy policy

# Approximate solutions

- Approximate Dynamic programming
  - ▶ Attempt to minimize  $\theta$  directly, e.g. OGD
  - ▶ Approximate policy iteration; e.g. LSPI [?]
- Approximate Linear program
  - ▶ For a feature matrix  $\Psi \in \mathbb{R}^{X \times d}$  for some  $d \ll X$

$$\max_{\lambda, h} \lambda,$$

$$\text{s.t. } B(\lambda \mathbf{1} + \Psi \theta) \leq \ell + P\Psi \theta$$

# Previous work

- Approximate Dynamic Programming (linear approximation of the value function): [??]
- Approximate Linear Programming: (approximately solving LP) [?????????].
- Solving LMDPs (with no theoretical guarantees): [?] and [??]

## Previous work: average cost

- Average cost suffers from a new set of problems
  - ▶ State-relevance vectors are not in average cost LP
  - ▶ Lyapunov function ideas are hard to extend
- First algorithms studied: [?]
  - ▶ Awkward: had one LP to estimate  $\lambda$  and a second to estimate  $h^*$
  - ▶ Requires feasibility
- [?] first looked at minimizing the dual LP, but provided no performance bounds (described versions of DP algorithms in the dual space)

# Outline

1 Motivation

**2 Linearly Solvable MDPs**

3 Extending to large dimensions

4 Experiments

# The Dual

- Recall: average cost LP

$$\begin{aligned} \max_{\lambda, h} \quad & \lambda, \\ \text{s.t.} \quad & B(\lambda \mathbf{1} + h) \leq \ell + Ph \end{aligned}$$

- Dual is

$$\begin{aligned} \min_{\mu \in \mathbb{R}^{XA}} \quad & \ell^\top \mu, \\ \text{s.t.} \quad & \mathbf{1}^\top \mu = 1, \mu \geq \mathbf{0}, (P - B)^\top \mu = \mathbf{0}. \end{aligned}$$

- $\mu$  is a stationary distribution of  $P^\pi$  for  $\pi(a|x) \propto \mu_{x,a}$

# The Dual ALP

- Feature matrix  $\Phi \in \mathbb{R}^{XA \times d}$ ; constrain  $\mu = \Phi\theta$ ,  $\theta \in B_2(0, S)$

$$\begin{aligned} \min_{\theta \in B_2(0, S)} \quad & \ell^T \Phi \theta, \\ \text{s.t.} \quad & \mathbf{1}^T \Phi \theta = 1, \Phi \theta \geq \mathbf{0}, (\mathbf{P} - \mathbf{B})^T \Phi \theta = \mathbf{0}. \end{aligned}$$

- Policy  $\pi_\theta(a|x) \propto [(\Phi\theta)(x, a)]_+$
- $\mu_\theta$  is the stationary distribution of  $P^{\pi_\theta}$
- Intuition:  $\mu_\theta \approx \Phi\theta$ , so  $\min_\theta \ell^T \mu_\theta \approx \min_\theta \ell^T \Phi\theta$

# Interpretation

- $\Phi_\theta$  is an approximate stationary distributions
- The primal of the dual ALP is:

$$\begin{aligned} \min_{\lambda, h} \lambda \\ \text{s.t. } \Phi^\top(\ell + (P - B)h - \lambda \mathbf{1}) \in \Phi^+ \end{aligned}$$

where  $\Phi^+ = \{x \in \mathbb{R}^d \mid \exists \nu \geq 0 \text{ s.t. } x = \Phi^\top \nu\}$

- Similar to weighted constraint aggregation



# Reducing Constraints

- Still intractable:  $d$ -dimensional problem but  $O(XA)$  constraints
- Form the convex cost function:

$$\begin{aligned}c(\theta) &= \ell^\top \Phi \theta + H \|\lceil \Phi \theta \rceil_-\|_1 + H \|(\mathbf{P} - \mathbf{B})^\top \Phi \theta\|_1 \\ &= \ell^\top \Phi \theta + H \sum_{(x,a)} |[\Phi_{(x,a),:} \theta]_-| + H \sum_{x'} |(\Phi \theta)^\top (\mathbf{P} - \mathbf{B})_{:,x'}|\end{aligned}$$

- Sample  $(x_t, a_t) \sim q_1$  and  $y_t \sim q_2$
- Unbiased subgradient estimate:

$$\begin{aligned}g_t(\theta) &= \ell^\top \Phi - H \frac{\Phi_{(x_t, a_t),:}}{q_1(x_t, a_t)} \mathbb{I}\{\Phi_{(x_t, a_t),:} \theta < 0\} \\ &\quad + H \frac{(\Phi^\top (\mathbf{P} - \mathbf{B})_{:, y_t})^\top}{q_2(y_t)} \operatorname{sgn}((\Phi \theta)^\top (\mathbf{P} - \mathbf{B})_{:, y_t})\end{aligned}$$

# The Stochastic Subgradient Method for MDPs

**Input:** Constants  $S, H > 0$ , number of rounds  $T$ .  
Let  $\Pi_{\Theta}$  be the Euclidean projection onto  $S$ -radius 2-norm ball.  
Initialize  $\theta_1 \propto 1$ .  
**for**  $t := 1, 2, \dots, T$  **do**  
    Sample  $(x_t, a_t) \sim q_1$  and  $y_t \sim q_2$ .  
    Compute subgradient estimate  $g_t$   
    Update  $\theta_{t+1} = \Pi_{\Theta}(\theta_t - \eta_t g_t)$ .  
**end for**  
 $\hat{\theta}_T = \frac{1}{T} \sum_{t=1}^T \theta_t$ .  
Return policy  $\pi_{\hat{\theta}_T}$ .

## Theorem

Given some  $\epsilon > 0$ , the  $\hat{\theta}_T$  produced by the stochastic subgradient method after  $T = 1/\epsilon^4$  steps satisfies

$$\ell^\top \mu_{\hat{\theta}_T} \leq \min_{\theta \in B(0, S)} \left( \ell^\top \mu_\theta + \frac{V(\theta)}{\epsilon} + O(\epsilon) \right)$$

with probability at least  $1 - \delta$ , where  $V = O(V_1 + V_2)$  is a violation function defined by

$$V_1(\theta) = \|\lceil \Phi \theta \rceil - \theta\|_1$$

$$V_2(\theta) = \|(\mathbf{P} - \mathbf{B})^\top \Phi \theta\|_1.$$

The big-O notation hides polynomials in  $S$ ,  $d$ ,  $C_1$ ,  $C_2$ , and  $\log(1/\delta)$ .

# Discussion

- Previous bounds were of the form  $\inf_r \|h^* - \Psi r\|$
- Can remove the awkward  $V(\theta)/\epsilon + O(\epsilon)$  by taking a grid of  $\epsilon$
- Constants:

$$C_1 = \underbrace{\max_{x,a} \frac{\|\Phi(x,a)_{:,x}\|}{q_1(x,a)}}_{\text{Control via } \Phi \text{ and } q_1}, \quad C_2 = \underbrace{\max_x \frac{\|\Phi^T(P-B)_{:,x}\|}{q_2(x)}}_{\text{control via structure of } P}$$

- $V(\theta^*)$  measures the difficulty of the problem
- Assume fast mixing: for every policy  $\pi$ ,  $\exists \tau(\pi) > 0$  s.t.  $\forall d, d' \in \Delta_{\mathcal{X}}$ ,

$$\|dP^\pi - d'P^\pi\|_1 \leq e^{-1/\tau(\pi)} \|d - d'\|_1$$

# Proof Outline

- First,

$$\mathbf{1}^\top \mu = 1, \|\mu\|_1 \leq 1 + \epsilon_1, \|\mu^\top (\mathbf{P} - \mathbf{B})\|_1 \leq \epsilon_2 \Rightarrow$$
$$\left\| \mu_{\pi_{\mu^+}} - \mu \right\|_1 \leq \tau(\mu_{\mu}) \log(1/\epsilon_1) O(\epsilon_1 + \epsilon_2)$$

- SGD theorem

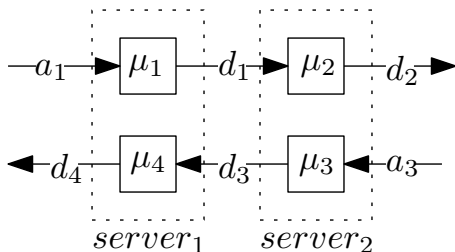
$$\ell^\top \Phi \hat{\theta}_T + H(V(\hat{\theta})) \leq \ell^\top \Phi \theta^* + H(V(\theta^*)) + O\left(\frac{SH(C_1 + C_2)}{\sqrt{T}}\right)$$

- Use  $\Phi \theta \approx \mu_\theta$

$$\ell^\top \mu_{\hat{\theta}_T} - \ell^\top \mu_{\theta^*} \leq HO(V_1(\theta^*) + V_2(\theta^*)) + O\left(\frac{H(C_1 + C_2)}{\sqrt{T}}\right)$$

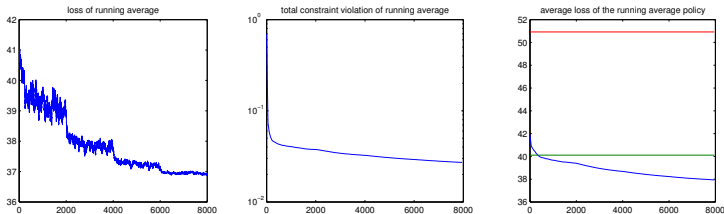
- Optimize  $H$  and  $T$

## Queueing network example (Rybko-Stolyar)



- Customers arrive at  $\mu_1/\mu_3$  then move to  $\mu_2/\mu_4$
- Server 1 processes  $\mu_1$  or  $\mu_4$ , server 2 processes  $\mu_2$  or  $\mu_3$
- Features: indicators of sub-blocks in state-action space, stationary distribution of LONGER and LBSF heuristics
- Loss is the total queue size
- $a_1 = a_3 = .08$ ,  $d_1 = d_2 = .12$ , and  $d_3 = d_4 = .28$ ,  $X = 902500$

# Experiments: Results



- The left plot: linear objective of the running average, i.e.  $\ell^T \Phi \hat{\theta}_t$ .
- The center plot: sum of the two constraint violations of  $\hat{\theta}_t$
- The right plot:  $\ell^T \mu_{\hat{\theta}_t}$ . The two horizontal lines correspond to the loss of two heuristics, LONGER and LBFS.

# Dual ALP Summary

- Presented an algorithm to solve average-cost large-scale MDPs
  - ▶ Restricted the dual LP to a subspace to reduce dimension
  - ▶ Used Stochastic Gradient Descent to sample constraints
- Presented oracle inequality guaranteeing we perform well w.r.t. best policy in the subspace.
- Demonstrated algorithm on a queueing network



- 1 Motivation
- 2 Linearly Solvable MDPs
- 3 Extending to large dimensions**
- 4 Experiments

# KL-cost

- Introduced in [?]
- $\mathcal{A} = \Delta_{\mathcal{X}}$ : we are playing polcies
- Loss:  $\ell(x, P) = q(x) + D_{KL}(\underbrace{P}_{\text{Learner's action}} \parallel \underbrace{P_0(\cdot|x)}_{\text{Base dynamics}})$
- Infinite loss unless  $P \ll P_0$
- Terminal state  $z$ :  $q(z) = 0$  and  $P_0(z|z) = 1$
- Obejective is total cost:

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

## Details

- $LJ_P(x) = \min_P \left[ \ell(x, P) + \sum_y P(y|x) J_P(y) \right]$
- Greedy action is:

$$P_J(\cdot|x) = \arg \min_{p \in \Delta_{\mathcal{X}}} \sum_y p(y) \log \frac{p(y)}{P_0(y|x) e^{-J(y)}} = \frac{P_0(x'|x) e^{-J(x')}}{Z(x)}$$

with  $Z = P_0 e^{-J}$

- This implies  $LJ = q - \log(Z)$
- Value function is the solution to:

$$J = LJ \Leftrightarrow J = q - \log(P_0 e^{-J})$$

- Exponentiating:  $LJ = J \Leftrightarrow e^{-q} P_0 e^{-J} = e^{-J}$

## Parameterizing $J_\theta$

- Previous approaches:  $J_\theta = \Psi\theta$
- Instead:  $J_\theta = -\log(\Psi\theta)$
- Surrogate optimization:

$$\min_{\theta} c^T J_\theta + \underbrace{\|LJ_\theta - J_\theta\|}_{\text{Bellman error}} \quad (1)$$

- $\|LJ_\theta - J_\theta\|$  not convex in  $\theta$ , but

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

- Plugging  $\Psi\theta = e^{-J_\theta}$  into (??):

$$\min_{\theta} -c^T \log(\Psi\theta) + \|e^{-qP_0} \Psi\theta - \Psi\theta\|$$

# Parameterizing $J_\theta$

- Previous approaches:  $J_\theta = \Psi\theta$
- Instead:  $J_\theta = -\log(\Psi\theta)$
- Surrogate optimization:

$$\min_{\theta} c^T J_\theta + \underbrace{\|LJ_\theta - J_\theta\|}_{\text{Bellman error}} \quad (1)$$

- $\|LJ_\theta - J_\theta\|$  not convex in  $\theta$ , but

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

- Plugging  $\Psi\theta = e^{-J_\theta}$  into (??):

$$\min_{\theta} -c^T \log(\Psi\theta) + \underbrace{\|e^{-qP_0} \Psi\theta - \Psi\theta\|}_{\text{Bellman operator}}$$

# Our algorithm

- Let  $\mathcal{T}$  be the set of trajectories with  $x_1 \sim c$  with distribution  $Q(\cdot)$
- Recall relaxed optimization:

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

- Optimization is equal to:

$$\min_{\theta} -c^T \log(\Psi\theta) + \sum_{T \in \mathcal{T}} Q(T) \sum_{x \in T} \left| e^{-q(x)} P_0 \Psi\theta(x) - \Psi\theta(x) \right|$$

- Use stochastic gradient descent by sampling trajectories

## Theorem

Let  $\hat{\theta}$  be an  $\epsilon$ -optimal solution returned by SGD. Then,

$$\begin{aligned} J_{P_{J_{\hat{\theta}}}}(x_1) &\leq \inf_{\theta \in \Theta} \left\{ J_{P_{J_{\theta}}}(x_1) + \mathcal{E}(J_{\theta}) \right\} + \epsilon \\ &\quad + \underbrace{\left\| P_{J_{\hat{\theta}}} - Q \right\|_1}_{\text{Off-policy error}} \max_{T \in \mathcal{T}} \sum_{x \in T} |J_{\hat{\theta}}(x) - LJ_{\hat{\theta}}(x)| \end{aligned}$$

Penalty function:

$$\mathcal{E}(J_{\theta}) = \sum_{T \in \mathcal{T}} \sum_{x \in T} \left( Q(T) e^{-\min(J_{\theta}, LJ_{\theta})} + P_{J_{\theta}}(T) \right) \underbrace{|J_{\theta}(x) - LJ_{\theta}(x)|}_{\text{Small if } J_{\theta} \text{ is close to the optimal value}}$$

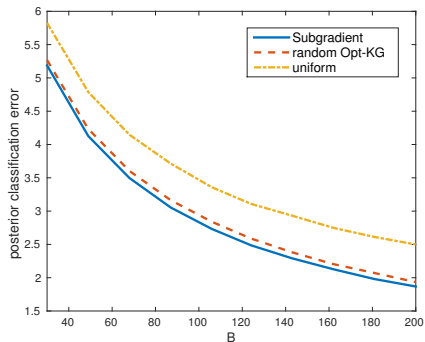
# Proof outline of main theorem

- $\left| J_{P_{J_{\theta^*}}}(x_1) - J_{\theta^*}(x_1) \right| = O(\|LJ_{\theta^*} - J_{\theta^*}\|)$
- Similarly bounding  $\left| J_{P_{J_{\hat{\theta}}}}(x_1) - J_{\hat{\theta}}(x_1) \right| = O(\|LJ_{\hat{\theta}} - J_{\hat{\theta}}\|)$
- $J_{\theta^*}$  and  $J_{\hat{\theta}}$  are close by the optimization

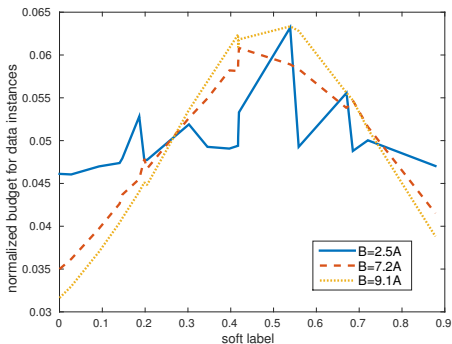


# Crowdsourcing

- 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  $a \in \{1, \dots, A\}$
  - ▶ 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



- 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
- Larger difference as  $B$  grows

Thanks!

# Crowdsourcing details

- Objective: posterior classification error
- Prior for label  $i$ :  $\text{Beta}(a_0^i, b_0^i)$
- State space: all possible integer increments for  $(a_0^i, b_0^i)$
- Define:  $l(a, b) = \Pr(\theta > .5 | \theta \sim \text{Beta}(a, b))$ ,  $h(x) = x \wedge (1 - x)$
- Opt-KG:  $p_i \propto [h(l(a_i + 1, b_i)) \wedge h(l(a_i, b_i + 1))] - h(l(a_i, b_i))$
- Base policy: Opt-KG
- Features: For each state  $\{a_i, b_i\}$ , all  $\mathbb{E}[X_i]$ ,  $1 - \mathbb{E}[X_i]$ , and  $\mathbb{E}[X_i^2]$  for  $X_i \sim \text{Beta}(a_i, b_i)$

# Proof part 1

## Lemma

Let  $u \in \mathbb{R}^{XA}$  be a vector with

$$\mathbf{1}^\top u = 1, \|u\|_1 \leq 1 + \epsilon_1, \|u^\top(P - B)\|_1 \leq \epsilon_2$$

For the stationary distribution  $\mu_u$  of policy  $u^+ = [u]_+ / \|[u]_+\|_1$ , we have

$$\|\mu_u - u\|_1 \leq \tau(\mu_u) \log(1/\epsilon_1) O(\epsilon_1 + \epsilon_2)$$

- Let  $\mu_t$  be  $u^+$  after  $t$  steps
- $\|\mu_t - u^+\|_1 = O(t(\epsilon_1 + \epsilon_2))$
- Mixing assumption:  $\|\mu_t - \mu_u\|_1 \leq e^{-t/\tau(u^+)}$
- Take  $t = \tau(u^+) \log(1/\epsilon')$  and use triangle inequality

# Applying SGD theorem

## Theorem (Lemma 3.1 of [?])

Assume we have

- Convex set  $\mathcal{Z} \subseteq B_2(Z, 0)$  and  $(f_t)_{t=1,2,\dots,T}$  convex functions on  $\mathcal{Z}$ .
- Gradient estimates  $f'_t$  with  $\mathbb{E}[f'_t|z_t] = \nabla f(z_t)$  and bound  $\|f'_t\|_2 \leq F$
- Sample Path  $z_1 = 0$  and  $z_{t+1} = \Pi_{\mathcal{Z}}(z_t - \eta f'_t)$  ( $\Pi_{\mathcal{Z}}$  Euclidean projection)

Then, for  $\eta = Z/(F\sqrt{T})$  and any  $\delta \in (0, 1)$ , the following holds with probability at least  $1 - \delta$ :

$$\sum_{t=1}^T f_t(z_t) - \min_{z \in \mathcal{Z}} \sum_{t=1}^T f_t(z) \leq O\left(Z\sqrt{T} \left(F + \sqrt{\log \frac{1}{\delta}}\right)\right)$$

## checking conditions of theorem

- Recall gradient: for  $(x, a) \sim q_1$  and  $y \sim q_2$ ,  $g_t(\theta) =$

$$\ell^\top \Phi - H \frac{\Phi_{(x,a),:}}{q_1(x, a)} \mathbb{I}\{\Phi_{(x,a),:} \theta < 0\} + H \frac{(P - B)^\top_{:,y} \Phi}{q_2(y)} \operatorname{sgn}((\Phi \theta)^\top (P - B)_{:,y})$$

- We can bound  $\|g_t(\theta)\|_2 \leq \sqrt{d} + H(C_1 + C_2) := F$
- $\mathbb{E}[g_t(\theta)] = \nabla c(\theta)$ .
- The SDG theorem gives

$$\ell^\top \Phi \hat{\theta}_T + H(V_1(\hat{\theta}) + V_2(\hat{\theta})) \leq \ell^\top \Phi \theta^* + H(V_1(\theta^*) + V_2(\theta^*)) + O\left(\frac{SH(C_1 + C_2)}{\sqrt{T}}\right)$$

## proof conclusion

- We take,  $i = 1, 2$

$$V_i(\hat{\theta}) \leq \frac{2 + 2S}{H} + V_1(\theta^*) + V_2(\theta^*) + O\left(\frac{C_1 + C_2}{\sqrt{T}}\right) := \epsilon'$$

- Apply the lemma to  $\Phi_{\hat{\theta}}$  and  $\Phi_{\theta^*}$ :

$$\begin{aligned} \ell^T \mu_{\hat{\theta}_T} - \ell^T \mu_{\theta^*} &\leq H(V_1(\theta^*) + V_2(\theta^*)) + O\left(\frac{H(C_1 + C_2)}{\sqrt{T}}\right) \\ &\quad + \tau(\mu_{\hat{\theta}_T}) \log(1/\epsilon') O(\epsilon') \\ &\quad + \tau(\mu_{\theta^*}) \log(1/(V_1(\theta^*) + V_2(\theta^*))) O(V_1(\theta^*) + V_2(\theta^*)) \\ &= HO(V_1(\theta^*) + V_2(\theta^*)) + O\left(\frac{H(C_1 + C_2)}{\sqrt{T}}\right) \end{aligned}$$

- Taking  $H = 1/\epsilon$  and  $T = 1/\epsilon^4$ :

$$\ell^T \mu_{\hat{\theta}_T} - \ell^T \mu_{\theta^*} \leq \frac{1}{\epsilon} (V_1(\theta^*) + V_1(\theta^*)) + O(\epsilon)$$