Learning in Markov Decision Problems

Peter Bartlett

Computer Science and Statistics University of California at Berkeley

Mathematical Sciences Queensland University of Technology

> UCLA November 10, 2014

MDP: Managing Threatened Species

For t = 1, 2, ...:

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □









For t = 1, 2, ...:

- See state X_t of ecosystem
- 2 Play an action A_t intervention anti-poaching patrols
- 3 Incur loss $\ell(X_t, A_t)$





MDP: Managing Threatened Species

For t = 1, 2, ...:

- See state X_t of ecosystem
- Play an action A_t

anti-poaching patrols

intervention

Solution Incur loss $\ell(X_t, A_t)$ S, extinction

• State evolves to $X_{t+1} \sim P_{X_t,A_t}$

Transition matrix:

イロト 不得 とくき とくき とうき

2/27

 $P: \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})$

MDP: Managing Threatened Species

For t = 1, 2, ...:

- See state X_t of ecosystem
- ⁽²⁾ Play an action A_t intervention anti-poaching patrols
- Solution Incur loss $\ell(X_t, A_t)$ S, extinction

• State evolves to $X_{t+1} \sim P_{X_t,A_t}$

Transition matrix:

 $P: \mathcal{X} imes \mathcal{A} o \Delta(\mathcal{X})$ Policy: $\pi: \mathcal{X} o \Delta(\mathcal{A})$

Performance Measure: Regret

$$R_T = \mathbb{E}\sum_{t=1}^T \ell(X_t, A_t) - \min_{\pi} \mathbb{E}\sum_{t=1}^T \ell(X_t^{\pi}, \pi(X_t^{\pi})).$$

MDP: Managing Threatened Species

For t = 1, 2, ...:

- See state X_t of ecosystem
- 2 Play an action A_t
- 3 Incur loss $\ell(X_t, A_t)$ \$, extinction

anti-poaching patrols

• State evolves to $X_{t+1} \sim P_{X_t,A_t}$

Transition matrix:

 $P: \mathcal{X} \times \mathcal{A} \to \Delta(\mathcal{X})$ Policy: $\pi : \mathcal{X} \to \Delta(\mathcal{A})$ Stationary distribution: μ Average loss: $\mu^T \ell$.

Performance Measure: Regret

$$R_T = \mathbb{E}\sum_{t=1}^T \ell(X_t, A_t) - \min_{\pi} \mathbb{E}\sum_{t=1}^T \ell(X_t^{\pi}, \pi(X_t^{\pi})).$$

intervention

MDP: Managing Threatened Species

For t = 1, 2, ...:

- See state X_t of ecosystem
- Play an action A_t
 - anti-poaching patrols

intervention

- Solution Incur loss $\ell(X_t, A_t)$ \$, extinction
- State evolves to $X_{t+1} \sim P_{X_t,A_t}$

 $\begin{array}{l} \text{Transition matrix:} \\ P: \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X}) \\ \text{Policy:} \quad \pi: \mathcal{X} \rightarrow \Delta(\mathcal{A}) \\ \text{Stationary distribution: } \mu \\ \text{Average loss:} \quad \mu^{\mathsf{T}} \ell. \end{array}$

Performance Measure: Excess Average Loss

$$\mu_{\pi}^{\mathsf{T}}\ell - \min_{\pi} \mu_{\pi}^{\mathsf{T}}\ell$$

Large MDP Problems:

When the state space ${\boldsymbol{\mathcal{X}}}$ is large, we must scale back the ambition of optimal performance.

Large MDP Problems:

When the state space \mathcal{X} is large, we must scale back the ambition of optimal performance:

In comparison to a restricted family of policies Π.
 e.g., linear value function approximation.
 Want a strategy that competes with the best policy.

1. Large-Scale Policy Design

1. Large-Scale Policy Design

 Compete with a restricted family of policies Π: Linearly parameterized approximate stationary distributions.

1. Large-Scale Policy Design

 Compete with a restricted family of policies Π: Linearly parameterized exponentially transformed value function.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

- Changing MDP; complete information.
- Exponential weights strategy.
- Competitive with small comparison class Π .
- Computationally efficient if ∏ has polynomial size.
- Hard for shortest path problems.

Large-scale policy design

(with Yasin Abbasi-Yadkori and Alan Malek. ICML2014)

• Stationary distributions dual to value functions.

Large-scale policy design

(with Yasin Abbasi-Yadkori and Alan Malek. ICML2014)

- Stationary distributions dual to value functions.
- Consider a class of policies defined by feature matrix Φ, distribution μ₀, and parameters θ:

$$\pi_{\theta}(\mathbf{a}|\mathbf{x}) = \frac{[\mu_0(\mathbf{x}, \mathbf{a}) + \Phi_{(\mathbf{x}, \mathbf{a}), :\theta}]_+}{\sum_{\mathbf{a}'} [\mu_0(\mathbf{x}, \mathbf{a}') + \Phi_{(\mathbf{x}, \mathbf{a}'), :\theta}]_+}$$

Large-scale policy design (with Yasin Abbasi-Yadkori and Alan Malek. ICML2014)

- Stationary distributions dual to value functions.
- Consider a class of policies defined by **feature matrix** Φ , distribution μ_0 , and parameters θ :

$$\pi_{ heta}(a|x) = rac{[\mu_0(x,a) + \Phi_{(x,a),:} heta]_+}{\sum_{a'} [\mu_0(x,a') + \Phi_{(x,a'),:} heta]_+} \; .$$

• Let μ_{θ} denote the stationary distribution of policy π_{θ} .

Large-scale policy design (with Yasin Abbasi-Yadkori and Alan Malek. ICML2014)

- Stationary distributions dual to value functions.
- Consider a class of policies defined by **feature matrix** Φ , distribution μ_0 , and parameters θ :

$$\pi_{ heta}(a|x) = rac{[\mu_0(x,a) + \Phi_{(x,a),:} heta]_+}{\sum_{a'} [\mu_0(x,a') + \Phi_{(x,a'),:} heta]_+} \; .$$

- Let μ_{θ} denote the stationary distribution of policy π_{θ} .
- Find $\widehat{\theta}$ such that $\mu_{\widehat{\theta}}^{\top} \ell \leq \min_{\theta \in \Theta} \mu_{\theta}^{\top} \ell + \epsilon$.

Large-scale policy design (with Yasin Abbasi-Yadkori and Alan Malek. ICML2014)

- Stationary distributions dual to value functions.
- Consider a class of policies defined by **feature matrix** Φ , distribution μ_0 , and parameters θ :

$$\pi_{ heta}(a|x) = rac{[\mu_0(x,a) + \Phi_{(x,a),:} heta]_+}{\sum_{a'} [\mu_0(x,a') + \Phi_{(x,a'),:} heta]_+} \; .$$

- Let μ_{θ} denote the stationary distribution of policy π_{θ} .
- Find $\widehat{\theta}$ such that $\mu_{\widehat{\theta}}^{\top} \ell \leq \min_{\theta \in \Theta} \mu_{\theta}^{\top} \ell + \epsilon$.
- Large-scale policy design: Independent of size of X.

• Define a constraint violation function

$$V(\theta) = \underbrace{\|[\mu_0 + \Phi\theta]_-\|_1}_{\text{prob. dist.}} + \underbrace{\|(P - B)^\top (\mu_0 + \Phi\theta)\|_1}_{\text{stationary}}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで

• Define a constraint violation function

$$V(\theta) = \left\| \left[\mu_0 + \Phi \theta \right]_{-} \right\|_1 + \left\| (P - B)^\top (\mu_0 + \Phi \theta) \right\|_1$$

▲□▶ ▲□▶ ▲目▶ ▲目▶ - 目 - のへの

and consider the convex cost function $c(\theta) = \ell^{\top}(\mu_0 + \Phi\theta) + \alpha V(\theta).$

• Define a constraint violation function

$$V(\theta) = \left\| \left[\mu_0 + \Phi \theta \right]_{-} \right\|_1 + \left\| \left(P - B \right)^\top (\mu_0 + \Phi \theta) \right\|_1$$

and consider the convex cost function $c(\theta) = \ell^{\top}(\mu_0 + \Phi\theta) + \alpha V(\theta).$

• Stochastic gradient descent: $\theta_{t+1} = \theta_t - \eta g_t(\theta_t)$, $\hat{\theta}_T = \sum_{t=1}^T \theta_t / T$,

• Define a constraint violation function

$$V(\theta) = \left\| \left[\mu_0 + \Phi \theta \right]_{-} \right\|_1 + \left\| (P - B)^\top (\mu_0 + \Phi \theta) \right\|_1$$

and consider the convex cost function $c(\theta) = \ell^{\top}(\mu_0 + \Phi\theta) + \alpha V(\theta).$

- Stochastic gradient descent: $\theta_{t+1} = \theta_t \eta g_t(\theta_t)$, $\hat{\theta}_T = \sum_{t=1}^T \theta_t / T$,
- ... with cheap, unbiased stochastic subgradient estimates:

$$g_t(\theta) = \ell^\top \Phi - \alpha \frac{\Phi_{(x_t, a_t),:}}{q_1(x_t, a_t)} \mathbb{I}_{\{\mu_0(x_t, a_t) + \Phi_{(x_t, a_t),:}\theta < 0\}} + \alpha \frac{(P - B)_{:,x_t'}^\top \Phi}{q_2(x_t')} \operatorname{sign}((P - B)_{:,x_t'}^\top \Phi \theta).$$

<ロ> < 部> < 注> < 注> < 注) と 注 の Q (で 6/27

Main Result

For ${\cal T}=1/\epsilon^4$ gradient estimates, with high probability (under a mixing assumption),

$$\mu_{\widehat{\theta}_{\mathcal{T}}}^{\top} \ell \leq \min_{\theta \in \Theta} \left(\mu_{\theta}^{\top} \ell + \frac{V(\theta)}{\epsilon} + O(\epsilon) \right)$$

Main Result

For $\mathcal{T}=1/\epsilon^4$ gradient estimates, with high probability (under a mixing assumption),

$$\mu_{\widehat{ heta}_{ au}}^{ op}\ell\leq\min_{ heta\in\Theta}\left(\mu_{ heta}^{ op}\ell+rac{V(heta)}{\epsilon}+O(\epsilon)
ight) \;.$$

• Competitive with all policies (stationary distributions) in the linear subspace (i.e., $V(\theta) = 0$).

Main Result

For $\mathcal{T}=1/\epsilon^4$ gradient estimates, with high probability (under a mixing assumption),

$$\mu_{\widehat{\theta}_{\mathcal{T}}}^{\top} \ell \leq \min_{\theta \in \Theta} \left(\mu_{\theta}^{\top} \ell + \frac{V(\theta)}{\epsilon} + O(\epsilon) \right) \ .$$

- Competitive with all policies (stationary distributions) in the linear subspace (i.e., $V(\theta) = 0$).
- Competitive with other policies; comparison more favorable near some stationary distribution in the subspace.

Main Result

For $\mathcal{T}=1/\epsilon^4$ gradient estimates, with high probability (under a mixing assumption),

$$\mu_{\widehat{\theta}_{\mathcal{T}}}^{\top} \ell \leq \min_{\theta \in \Theta} \left(\mu_{\theta}^{\top} \ell + \frac{V(\theta)}{\epsilon} + O(\epsilon) \right) \; .$$

- Competitive with all policies (stationary distributions) in the linear subspace (i.e., $V(\theta) = 0$).
- Competitive with other policies; comparison more favorable near some stationary distribution in the subspace.
- Previous results of this kind:
 - require knowledge about optimal policy, or
 - require that the comparison class ∏ contains a near-optimal policy.

Simulation Results: Queueing



(Rybko and Stolyar, 1992; de Farias and Van Roy, 2003a)

Simulation Results: Queueing







8/27
1. Large-Scale Policy Design

 Compete with a restricted family of policies Π: Linearly parameterized approximate stationary distributions.

1. Large-Scale Policy Design

 Compete with a restricted family of policies Π: Linearly parameterized approximate stationary distributions. Linearly parameterized exponentially transformed value function.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized approximate stationary distributions. Linearly parameterized exponentially transformed value function.
- Stochastic gradient convex optimization.
- Competitive with policies in the approximating class.
- Simulation results: crowdsourcing.

Large-scale policy design

(with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

Large-scale policy design

(with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

Consider total cost:

(assume a.s. hit absorbing state with zero loss)



Large-scale policy design

(with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

Consider total cost:

(assume a.s. hit absorbing state with zero loss)



• Parameterized value functions, close to a reference policy

 $P_0: \mathcal{X} \to \Delta(\mathcal{X}).$

Large-scale policy design

(with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

Consider total cost:

(assume a.s. hit absorbing state with zero loss)

$$\mathbb{E}\sum_{t=1}^{\infty}\ell(X_t)$$

• Parameterized value functions, close to a reference policy

 $P_0: \mathcal{X} \to \Delta(\mathcal{X}).$

 Regularize with KL-divergence to P₀: (so optimization is linear; Todorov/Kappen/Fleming)

 $\ell(x, P) = \ell(x) + d_{\mathcal{KL}}(P(\cdot|x), P_0(\cdot|x)).$

P is transition matrix under policy.

Large-scale policy design (with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek) Consider a class of policies defined by feature matrix Φ and

 Consider a class of policies defined by feature matrix Φ, and parameters θ:

 $\Pi = \left\{ G\widehat{J}_{\theta} : \theta \in \Theta \right\}$ $G\widehat{J}(x) := \arg \min_{\pi} \left(\ell(x, P^{\pi}) + \mathbb{E}^{\pi} \left[\widehat{J}(x') | x \right] \right) \quad \text{greedy policies}$ $\widehat{J}_{\theta} = -\log(\Phi\theta). \quad \text{log linear}$

イロン イロン イヨン イヨン 三日

Large-scale policy design (with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

 Consider a class of policies defined by feature matrix Φ, and parameters θ:

$$\begin{split} &\Pi = \Big\{ G \widehat{J}_{\theta} : \theta \in \Theta \Big\} \\ &G \widehat{J}(x) := \arg\min_{\pi} \Big(\ell(x, P^{\pi}) + \mathbb{E}^{\pi} \left[\widehat{J}(x') | x \right] \Big) \quad \text{ greedy policies} \\ &\widehat{J}_{\theta} = -\log(\Phi \theta). \quad \text{ log linear} \end{split}$$

• Find parameters $\widehat{ heta}$ (hence policy $\hat{\pi} = G \widehat{J}_{\hat{ heta}}$) such that

$$egin{aligned} &J_{\hat{\pi}}(x_1) \leq \min_{\pi \in \Pi} J_{\pi}(x_1) + \epsilon \ &J_{\pi}(x) := \mathbb{E}^{\pi} \left[\left. \sum_{t=1}^{\infty} \ell(X_t) \right| X_1 = x
ight] \end{aligned}$$

Large-scale policy design (with Yasin Abbasi-Yadkori, Xi Chen and Alan Malek)

 Consider a class of policies defined by feature matrix Φ, and parameters θ:

$$\begin{split} &\Pi = \Big\{ G \widehat{J}_{\theta} : \theta \in \Theta \Big\} \\ &G \widehat{J}(x) := \arg\min_{\pi} \Big(\ell(x, P^{\pi}) + \mathbb{E}^{\pi} \left[\widehat{J}(x') | x \right] \Big) \quad \text{ greedy policies} \\ &\widehat{J}_{\theta} = -\log(\Phi \theta). \quad \text{ log linear} \end{split}$$

• Find parameters $\widehat{\theta}$ (hence policy $\hat{\pi}=G\widehat{J}_{\widehat{\theta}})$ such that

$$egin{aligned} &J_{\hat{\pi}}(x_1) \leq \min_{\pi \in \Pi} J_{\pi}(x_1) + \epsilon \ &J_{\pi}(x) := \mathbb{E}^{\pi} \left[\left. \sum_{t=1}^{\infty} \ell(X_t) \right| X_1 = x
ight] \end{aligned}$$

Approach: a Reduction to Convex Optimization

• Define a transformed Bellman error function

$$V(\theta) = \underbrace{\|\Phi\theta - \exp(-\ell(x))P_0\Phi\theta\|}_{\text{convex in }\theta}$$

Approach: a Reduction to Convex Optimization

• Define a transformed Bellman error function (|| · || is a 1-norm over trajectories)

$$V(\theta) = \|\Phi\theta - \exp(-\ell(x))P_0\Phi\theta\| = \left\|\exp\left(-\widehat{J}_{\theta}\right) - \exp\left(-T\widehat{J}_{\theta}\right)\right\|$$
$$T\widehat{J}(x) := \min_{\pi} \left(\ell(x, P^{\pi}) + \mathbb{E}^{\pi}\left[\widehat{J}(x')|x\right]\right) \qquad (\text{dynamic prog operator})$$

• Define a transformed Bellman error function (|| · || is a 1-norm over trajectories)

$$V(\theta) = \|\Phi\theta - \exp(-\ell(x))P_0\Phi\theta\| = \left\|\exp\left(-\widehat{J}_{\theta}\right) - \exp\left(-T\widehat{J}_{\theta}\right)\right\|$$
$$T\widehat{J}(x) := \min_{\pi} \left(\ell(x, P^{\pi}) + \mathbb{E}^{\pi}\left[\widehat{J}(x')|x\right]\right) \qquad (\text{dynamic prog operator})$$

and consider the convex cost function $c(\theta) = \widehat{J_{\theta}} + \alpha V(\theta).$ • Define a transformed Bellman error function (|| · || is a 1-norm over trajectories)

$$V(\theta) = \|\Phi\theta - \exp(-\ell(x))P_0\Phi\theta\| = \left\|\exp\left(-\widehat{J}_{\theta}\right) - \exp\left(-T\widehat{J}_{\theta}\right)\right\|$$
$$T\widehat{J}(x) := \min_{\pi} \left(\ell(x, P^{\pi}) + \mathbb{E}^{\pi}\left[\widehat{J}(x')|x\right]\right) \qquad (\text{dynamic prog operator})$$

イロン イロン イヨン イヨン 三日

and consider the convex cost function $c(\theta) = \widehat{J}_{\theta} + \alpha V(\theta).$

• Stochastic gradient descent

• Define a transformed Bellman error function (|| · || is a 1-norm over trajectories)

$$V(\theta) = \|\Phi\theta - \exp(-\ell(x))P_0\Phi\theta\| = \left\|\exp\left(-\widehat{J}_{\theta}\right) - \exp\left(-T\widehat{J}_{\theta}\right)\right\|$$
$$T\widehat{J}(x) := \min_{\pi} \left(\ell(x, P^{\pi}) + \mathbb{E}^{\pi}\left[\widehat{J}(x')|x\right]\right) \qquad (\text{dynamic prog operator})$$

and consider the convex cost function $c(\theta) = \widehat{J}_{\theta} + \alpha V(\theta).$

- Stochastic gradient descent
- ... with cheap, unbiased stochastic subgradient estimates.

Main Result

For $T = 1/\epsilon^4$ gradient estimates, with high probability,

$$J_{\hat{\pi}}(x_1) \leq \min_{\pi \in \Pi} \left(J_{\pi}(x_1) + rac{1}{\epsilon} \left\| \widehat{J}_{ heta} - \mathcal{T} \widehat{J}_{ heta}
ight\|
ight) + \left\| \widehat{J}_{\hat{ heta}} - \mathcal{T} \widehat{J}_{\hat{ heta}}
ight\|' + O(\epsilon).$$

Main Result

For $\mathcal{T}=1/\epsilon^4$ gradient estimates, with high probability,

$$J_{\hat{\pi}}(x_1) \leq \min_{\pi \in \Pi} \left(J_{\pi}(x_1) + \frac{1}{\epsilon} \left\| \widehat{J}_{\theta} - T \widehat{J}_{\theta} \right\| \right) + \left\| \widehat{J}_{\hat{\theta}} - T \widehat{J}_{\hat{\theta}} \right\|' + O(\epsilon).$$

• Competitive with all policies in the parameterized class, up to penalties involving the Bellman errors.

Main Result

For $\mathcal{T}=1/\epsilon^4$ gradient estimates, with high probability,

$$J_{\widehat{\pi}}(x_1) \leq \min_{\pi \in \Pi} \left(J_{\pi}(x_1) + rac{1}{\epsilon} \left\| \widehat{J}_{ heta} - \mathcal{T} \widehat{J}_{ heta}
ight\|
ight) + \left\| \widehat{J}_{\widehat{ heta}} - \mathcal{T} \widehat{J}_{\widehat{ heta}}
ight\|' + O(\epsilon).$$

- Competitive with all policies in the parameterized class, up to penalties involving the Bellman errors.
- Unfortunately:
 - require that the comparison class Π contains a near-optimal policy.

• Classification task.

64BCE25 X91	KION 2469/	RR9 PYAR	FUH DINE C
85% WREHS	TINEWHIG	ENW KX SAN HIE	(ČEŘJČĒBÚ,
(88H255R) -	T324 AEPS	KAWHENEY;	>VM16CTUFY
E3XNW/NE	121×1469月)	DXX SCHAND	SP?KF45NT

- Classification task.
- Crowdsource labels.

14BCE25 ¥91	KJON 2469/	RR9 PYAR	334TGFTU7
85% WR5HS	TIHMMATS	ENW KY Pr	(ČERJOBBU)
(88H255R) -	JT 3 Rt A EPQ	WAWH MEN;	>VM16CFUFT
E3XNW/NE	シアマヤを見り	DXXSCHAN	SP?KFUSN'ST

Imp://www.schriciafia.net/ http://www.schriciafia.net/ http://www.schriciafia.net/ 14 / 27

- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.

64BCE25 ¥91	KION 2489/	RR9 PYAR	334ITGFTU7
85% WREHS	TIHUNATS	ENW W FY T	(ČERJOBBÁ,
(88H255R) -	T3KAEPS	X AWHENEY;	>YM16CFUFJ
E3XNW/NE	シアマやき日	DITESCHARD	SP?KFUSNST



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i)

64BCE25 ¥91	KION 2489/	RR9 PYAR	334TGFTU7
85% WREHS	TIHMMATS	FINITS AND	(ČEŘJČBBÁ,
(88H255R) -	JT32742PD	WAWH MEN;	>VM16CFUFT
E3XNW/NE	这次王昭国	DXX SCHAND	SP ? KFUSRIST



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i)

64BCE25791	KIN 2469/	RR9 PYAR	334ITGFTU7
85% WR5HC	TIHMATT	ENWARS P	(ČEŘIČBBÁ,
(88H235R) -	JT324 AEPO	X AWHENEY;	>YM16CFUFJ
E3XNW/NE	シアマやの日	DXX SCHAN	SP ? KFUS NOT



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i)

14BEE25 X91	KJON 2469/	RR9 PYAR	334TGFTU7
85 W WREHS	TINNERG	ENW WX PY Pr	(ČERJOBBÁ,
(88H255R) -	JT324 AEPO	WAWHENEY;	>YM16CFUFJ
E3XNW/NE	シンマヤを見り	DXX SCHARD	SP?KFUSRIST



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i); p_i ∼ Beta.

64BCE25 ¥91	KION 2489/	RR9 PYAR	334ITGFTU7
85% WR5HC	TINERHTS	ENWARS P	(ČERJOBBA
(88H255R) -	T3KAEPS	X AWHENEY;	74M16CFUFJ
E3XNW YOUE	シアマやき日	DXX SCHAN	SP ? KFUSNST



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i); p_i ∼ Beta.
- State = posterior.

64BCE25 Y91	KION 2469/	RR9 PVARW	334TGFTU7
85% WREHS	TINNERG	ENW WX PY Pr	(Č ^e ŘJČĒ BÁ,
(88H255R) -	JT324 AEPO	WAWHENEY;	>VM16CFUFJ
E3XNW/NE	シンマヤを見り	DXX SCHARD	SP?KFUSRIST



- Classification task.
- Crowdsource: \$ for labels.
- Fixed budget; minimize errors.
- Bayesian model: binary labels, i.i.d. crowd; Y_i ∼ Bernoulli(p_i); p_i ∼ Beta.
- State = posterior.



1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

2. Learning Changing Dynamics

• Changing MDP; complete information.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

- Changing MDP; complete information.
- Exponential weights strategy.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

- Changing MDP; complete information.
- Exponential weights strategy.
- Competitive with small comparison class Π .

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

- Changing MDP; complete information.
- Exponential weights strategy.
- Competitive with small comparison class Π .
- Computationally efficient if **Π** has polynomial size.

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

- Changing MDP; complete information.
- Exponential weights strategy.
- Competitive with small comparison class Π .
- Computationally efficient if ∏ has polynomial size.
- Hard for shortest path problems.

(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

• Observe P_t , ℓ_t after round t.

(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

- Observe P_t , ℓ_t after round t.
- Consider a comparison class: $\Pi \subset \{\pi \mid \pi : \mathcal{X} \to \mathcal{A}\}$

(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

- Observe P_t , ℓ_t after round t.
- Consider a comparison class: $\Pi \subset \{\pi \mid \pi : \mathcal{X} \to \mathcal{A}\}$
- $\pi^* = \operatorname{argmin}_{\pi \in \Pi} \sum_{t=1}^T \ell_t(x_t^{\pi}, \pi(x_t^{\pi}))$
(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

- Observe P_t , ℓ_t after round t.
- Consider a comparison class: $\Pi \subset \{\pi \mid \pi : \mathcal{X} \to \mathcal{A}\}$
- $\pi^* = \operatorname{argmin}_{\pi \in \Pi} \sum_{t=1}^{T} \ell_t(x_t^{\pi}, \pi(x_t^{\pi}))$
- $R_T = \sum_{t=1}^T \ell_t(x_t, a_t) \sum_{t=1}^T \ell_t(x_t^{\pi^*}, \pi^*(x_t^{\pi^*}))$

(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

16/27

- Observe P_t , ℓ_t after round t.
- Consider a comparison class: $\Pi \subset \{\pi \mid \pi : \mathcal{X} \to \mathcal{A}\}$
- $\pi^* = \operatorname{argmin}_{\pi \in \Pi} \sum_{t=1}^{T} \ell_t(x_t^{\pi}, \pi(x_t^{\pi}))$
- $R_T = \sum_{t=1}^T \ell_t(x_t, a_t) \sum_{t=1}^T \ell_t(x_t^{\pi^*}, \pi^*(x_t^{\pi^*}))$
- Aim for low regret: $R_T/T \rightarrow 0$

(with Yasin Abbasi-Yadkori, Varun Kanade, Yevgeny Seldin, Csaba Szepesvari, NIPS2013)

- Observe P_t , ℓ_t after round t.
- Consider a comparison class: $\Pi \subset \{\pi \mid \pi : \mathcal{X} \to \mathcal{A}\}$
- $\pi^* = \operatorname{argmin}_{\pi \in \Pi} \sum_{t=1}^{T} \ell_t(x_t^{\pi}, \pi(x_t^{\pi}))$
- $R_T = \sum_{t=1}^T \ell_t(x_t, a_t) \sum_{t=1}^T \ell_t(x_t^{\pi^*}, \pi^*(x_t^{\pi^*}))$
- Aim for low regret: $R_T/T \rightarrow 0$
- Computationally efficient low regret strategies?

There is a strategy that (under a τ -mixing assumption) achieves

 $\mathbb{E}[R_T] \leq (4 + 2\tau^2)\sqrt{T \log |\Pi|} + \log |\Pi|.$

・ロ ・ ・ 一 ・ ・ 三 ・ ・ 三 ・ シ へ (* 17/27)

Strategy for a repeated game:

Choose action $a \in \mathcal{A}$ with probability proportional to

exp(total loss a has incurred so far).

Strategy for a repeated game:

Choose action $a \in \mathcal{A}$ with probability proportional to

exp(total loss *a* has incurred so far).

イロン イロン イヨン イヨン 三日

18 / 27

• Regret (total loss versus best in hindsight) for T rounds: $O\left(\sqrt{T \log |\mathcal{A}|}\right).$

Strategy for a repeated game:

Choose action $a \in \mathcal{A}$ with probability proportional to

exp(total loss *a* has incurred so far).

イロン イヨン イヨン イヨン 三日

18 / 27

- Regret (total loss versus best in hindsight) for T rounds: $O\left(\sqrt{T \log |\mathcal{A}|}\right).$
- Long history.

Strategy for a repeated game:

Choose action $a \in \mathcal{A}$ with probability proportional to

exp(total loss *a* has incurred so far).

- Regret (total loss versus best in hindsight) for T rounds: $O\left(\sqrt{T \log |\mathcal{A}|}\right).$
- Long history.
- Unreasonably broadly applicable:
 - Zero-sum games.
 - AdaBoost.
 - Bandit problems.
 - Linear programming.

- Shortest path problems.
- Fast max-flow.
- Fast graph sparsification.
- Model of evolution. ¹ ≥ ² ³ ³ ³ ³ ³

Strategy:

```
For all policies \pi \in \Pi, w_{\pi,0} = 1.

W_t = \sum_{\pi \in \Pi} w_{\pi,t}, p_{\pi,t} = w_{\pi,t-1}/W_{t-1}.

for t := 1, 2, ... do

w.p. \beta_t = \frac{W_{\pi_{t-1},t-1}}{W_{\pi_{t-1},t-2}}, \pi_t = \pi_{t-1}. Otherwise \pi_t \sim p_{.,t}.

Choose action a_t \sim \pi_t(.|x_t).

Observe dynamics P_t and loss \ell_t.

Suffer \ell_t(x_t, a_t).

For all policies \pi, w_{\pi,t} = w_{\pi,t-1} \exp(-\eta \mathbb{E}[\ell_t(x_t^{\pi}, \pi)]).

end for
```

イロト 不得下 イヨト イヨト 二日

19/27

• Exponential weights on Π .

Strategy:

```
For all policies \pi \in \Pi, w_{\pi,0} = 1.

W_t = \sum_{\pi \in \Pi} w_{\pi,t}, p_{\pi,t} = w_{\pi,t-1}/W_{t-1}.

for t := 1, 2, ... do

w.p. \beta_t = \frac{w_{\pi_{t-1},t-1}}{w_{\pi_{t-1},t-2}}, \pi_t = \pi_{t-1}. Otherwise \pi_t \sim p_{.,t}.

Choose action a_t \sim \pi_t(.|x_t).

Observe dynamics P_t and loss \ell_t.

Suffer \ell_t(x_t, a_t).

For all policies \pi, w_{\pi,t} = w_{\pi,t-1} \exp(-\eta \mathbb{E}[\ell_t(x_t^{\pi}, \pi)]).

end for
```

- Exponential weights on Π .
- Rare, random changes to π_t .



• Adversarial dynamics and loss functions.

Main Result There is a strategy that (under a τ -mixing assumption) achieves $\mathbb{E}[R_T] \leq (4 + 2\tau^2)\sqrt{T \log |\Pi|} + \log |\Pi|$.

20 / 27

- Adversarial dynamics and loss functions.
- Large state and action spaces.

There is a strategy that (under a τ -mixing assumption) achieves

```
\mathbb{E}\left[R_{T}\right] \leq (4 + 2\tau^{2})\sqrt{T \log |\Pi|} + \log |\Pi|.
```

- Adversarial dynamics and loss functions.
- Large state and action spaces.
- $\mathbb{E}[R_T]/T \to 0$ for $T = \omega(\log |\Pi|)$.

There is a strategy that (under a *T*-mixing assumption) achieves

```
\mathbb{E}\left[R_{T}\right] \leq (4 + 2\tau^{2})\sqrt{T \log |\Pi|} + \log |\Pi|.
```

- Adversarial dynamics and loss functions.
- Large state and action spaces.
- $\mathbb{E}[R_T]/T \to 0$ for $T = \omega(\log |\Pi|)$.
- Computationally efficient as long as $|\Pi|$ is polynomial.

There is a strategy that (under a *T*-mixing assumption) achieves

```
\mathbb{E}\left[R_{T}\right] \leq (4 + 2\tau^{2})\sqrt{T \log |\Pi|} + \log |\Pi|.
```

- Adversarial dynamics and loss functions.
- Large state and action spaces.
- $\mathbb{E}[R_T]/T \to 0$ for $T = \omega(\log |\Pi|)$.
- Computationally efficient as long as $|\Pi|$ is polynomial.
- No computationally efficient algorithm in general

Shortest Path Problem

Special case of MDP: node=state; action=edge; loss=weight.





http://www.google.com/

Shortest Path Problem

Special case of MDP: node=state; action=edge; loss=weight.





http://www.google.com/ http://www.meondirect.com/

Shortest Path Problem

Special case of MDP: node=state; action=edge; loss=weight.





http://www.google.com/ http://www.meondirect.com/

Hardness Result

Suppose there is a strategy for the online adversarial shortest path problem that:

- runs in time poly(n, T), and
- Solution Set in the set of t

Then there is an efficient algorithm for online agnostic parity learning with sublinear regret.

 Class of parity functions on {0,1}ⁿ: PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}

- Class of parity functions on {0,1}ⁿ: PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}
- Learning problem: given x_t ∈ {0,1}ⁿ, learner predicts ŷ_t ∈ {0,1}, observes the true label y_t and suffers loss I_{{ŷt≠yt}}

- Class of parity functions on {0,1}ⁿ: PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}
- Learning problem: given $x_t \in \{0, 1\}^n$, learner predicts $\hat{y}_t \in \{0, 1\}$, observes the true label y_t and suffers loss $\mathbb{I}_{\{\hat{y}_t \neq y_t\}}$
- $R_T = \sum_{t=1}^{T} \mathbb{I}_{\{\hat{y}_t \neq y_t\}} \min_{\mathsf{PAR}_S \in \mathsf{PARITIES}} \sum_{t=1}^{T} \mathbb{I}_{\{\mathsf{PAR}_S(x_t) \neq y_t\}}$

- Class of parity functions on {0,1}ⁿ:
 PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}
- Learning problem: given $x_t \in \{0,1\}^n$, learner predicts $\hat{y}_t \in \{0,1\}$, observes the true label y_t and suffers loss $\mathbb{I}_{\{\hat{y}_t \neq y_t\}}$
- $R_T = \sum_{t=1}^T \mathbb{I}_{\{\hat{y}_t \neq y_t\}} \min_{\mathsf{PAR}_S \in \mathsf{PARITIES}} \sum_{t=1}^T \mathbb{I}_{\{\mathsf{PAR}_S(x_t) \neq y_t\}}$
- Is there an efficient (time polynomial in n, T) learning algorithm with sublinear regret (R_T = O(poly(n)T^{1-δ}) for some δ > 0)?

- Class of parity functions on {0,1}ⁿ:
 PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}
- Learning problem: given $x_t \in \{0,1\}^n$, learner predicts $\hat{y}_t \in \{0,1\}$, observes the true label y_t and suffers loss $\mathbb{I}_{\{\hat{y}_t \neq y_t\}}$
- $R_T = \sum_{t=1}^T \mathbb{I}_{\{\hat{y}_t \neq y_t\}} \min_{\mathsf{PAR}_S \in \mathsf{PARITIES}} \sum_{t=1}^T \mathbb{I}_{\{\mathsf{PAR}_S(x_t) \neq y_t\}}$
- Is there an efficient (time polynomial in n, T) learning algorithm with sublinear regret (R_T = O(poly(n)T^{1-δ}) for some δ > 0)?
- Very well-studied.

- Class of parity functions on {0,1}ⁿ:
 PARITIES = {PAR_S | S ⊂ [n], PAR_S(x) = ⊕_{i∈S}x_i}
- Learning problem: given $x_t \in \{0,1\}^n$, learner predicts $\hat{y}_t \in \{0,1\}$, observes the true label y_t and suffers loss $\mathbb{I}_{\{\hat{y}_t \neq y_t\}}$
- $R_T = \sum_{t=1}^T \mathbb{I}_{\{\hat{y}_t \neq y_t\}} \min_{\mathsf{PAR}_S \in \mathsf{PARITIES}} \sum_{t=1}^T \mathbb{I}_{\{\mathsf{PAR}_S(x_t) \neq y_t\}}$
- Is there an efficient (time polynomial in n, T) learning algorithm with sublinear regret (R_T = O(poly(n)T^{1-δ}) for some δ > 0)?
- Very well-studied.
- Widely believed to be hard: used for cryptographic schemes.

Hardness Result

Suppose there is a strategy for the online adversarial shortest path problem that:

- runs in time poly(n, T), and
- Solution has regret $R_T = O(\text{poly}(n)T^{1-\delta})$ for some constant $\delta > 0$.

Then there is an efficient algorithm for online agnostic parity learning with sublinear regret.



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Edges (dynamics)	Weights (costs)	
Adversarial	Adversarial	As hard as noisy parity.
Stochastic	Adversarial	Efficient algorithm.
Adversarial	Stochastic	Efficient algorithm.

Outline

1. Large-Scale Policy Design

- Compete with a restricted family of policies Π: Linearly parameterized policies.
- Stochastic gradient convex optimization.
- Competitive with policies near the approximating class.
- Without knowledge of optimal policy.
- Simulation results: queueing, crowdsourcing.

2. Learning changing dynamics

- Changing MDP; complete information.
- Exponential weights strategy.
- Competitive with small comparison class Π .
- Computationally efficient if ∏ has polynomial size.
- Hard for shortest path problems.