

# Near-optimal Reinforcement Learning in Factored MDPs

Ian Osband   Benjamin Van Roy

Management Science and Engineering  
Stanford University  
iosband@stanford.edu

November 7, 2014

# Table of contents

Reinforcement Learning

Efficient RL

Factored MDPs

Main results

# Reinforcement Learning

- We imagine an agent taking actions within an environment.
- Actions serve two purposes:
  - Instantaneous loss/reward to the agent.
  - Influence the state of the environment.
- The agent wants to maximize cumulative reward through time.
- How can an agent learn to take “good” actions?

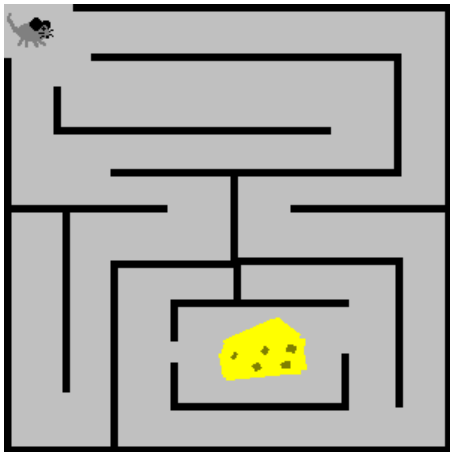
# Reinforcement Learning

- We imagine an agent taking actions within an environment.
  - Actions serve two purposes:
    - Instantaneous loss/reward to the agent.
    - Influence the state of the environment.
  - The agent wants to maximize cumulative reward through time.
  - How can an agent learn to take “good” actions?
- 
- Multi-armed bandit with added state transitions.
  - Statistical estimation + optimal control.

# Reinforcement Learning

- We imagine an agent taking actions within an environment.
  - Actions serve two purposes:
    - Instantaneous loss/reward to the agent.
    - Influence the state of the environment.
  - The agent wants to maximize cumulative reward through time.
  - How can an agent learn to take “good” actions?
- 
- Multi-armed bandit with added state transitions.
  - Statistical estimation + optimal control.
  - ... this could get hard!

# Mouse in a maze



What's the best way to the cheese?

# Self-driving car



Drive me from A to B

# Self-driving car



Drive me from A to B... also don't kill anyone.



# Markov decision process (MDP)

- Model decision making in a stochastic environment.
- At time  $t$  the agent in state  $s_t \in \mathcal{S}$  chooses action  $a_t \in \mathcal{A}$ .
- Reward  $r_t \sim R(s_t, a_t)$  and transition to  $s_{t+1} \sim P(\cdot | s_t, a_t)$ .
- Given  $R, P$  there is some optimal policy  $\pi^* : s_t \rightarrow a_t^*$ .
- Generally this is computed via Dynamic Programming.

# Markov decision process (MDP)

- Model decision making in a stochastic environment.
  - At time  $t$  the agent in state  $s_t \in \mathcal{S}$  chooses action  $a_t \in \mathcal{A}$ .
  - Reward  $r_t \sim R(s_t, a_t)$  and transition to  $s_{t+1} \sim P(\cdot | s_t, a_t)$ .
  - Given  $R, P$  there is some optimal policy  $\pi^* : s_t \rightarrow a_t^*$ .
  - Generally this is computed via Dynamic Programming.
- 
- In reinforcement learning, the agent is unsure of  $P, R$ .
  - The learning algorithm will pick  $\pi_t^L : s_t \rightarrow a_t^L$ .
  - We would like to get average reward  $\rho(\pi^L)$  close to  $\rho(\pi^*)$ .

# Efficient reinforcement learning

- Will we learn the best policy?

$$\rho(\pi_t^L) \rightarrow \rho(\pi^*)$$

# Efficient reinforcement learning

- Will we learn the best policy?

$$\rho(\pi_t^L) \rightarrow \rho(\pi^*)$$

- How long do we have to wait to do well? (Sample complexity)

$$\forall t > T(\text{MDP}, \pi^L) \quad \rho(\pi_t^L) \geq \rho(\pi^*) - \epsilon$$

# Efficient reinforcement learning

- Will we learn the best policy?

$$\rho(\pi_t^L) \rightarrow \rho(\pi^*)$$

- How long do we have to wait to do well? (Sample complexity)

$$\forall t > T(\text{MDP}, \pi^L) \quad \rho(\pi_t^L) \geq \rho(\pi^*) - \epsilon$$

- How badly do we do while we're learning? (Regret)

$$\text{Regret}(T, \pi^L) = \sum_{t=1}^T r_t^* - r_t \leq f(\text{MDP}, T, \pi^L)$$

## Reward and transition functions

We will specialize our analysis to two important function classes:

**Definition (Reward functions  $\in \mathcal{P}_{\mathcal{S} \times \mathcal{A}, \mathbb{R}}^{C, \sigma}$ )**

$\mathcal{P}_{\mathcal{X}, \mathbb{R}}^{C, \sigma}$  is the set of functions from  $\mathcal{X}$  to  $\sigma$ -sub gaussian measures over  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with mean in  $[0, C]$ .

**Definition (Transition functions  $\in \mathcal{P}_{\mathcal{S} \times \mathcal{A}, \mathcal{S}}$ )**

$\mathcal{P}_{\mathcal{X}, \mathcal{Y}}$  is the set of functions mapping elements of a finite set  $\mathcal{X}$  to probability mass functions over a finite set  $\mathcal{Y}$ .

# Regret bounds

- Greedy algorithms may never learn  $\rho(\pi_t^L) \rightarrow \rho(\phi^*)$ .
- Naïve exploration leads to regret exponential in  $|\mathcal{S}|, |\mathcal{A}|$ .
- Efficient algorithms must guide their exploration:

# Regret bounds

- Greedy algorithms may never learn  $\rho(\pi_t^L) \rightarrow \rho(\phi^*)$ .
- Naïve exploration leads to regret exponential in  $|\mathcal{S}|, |\mathcal{A}|$ .
- Efficient algorithms must guide their exploration:
  - “Optimism in the face of uncertainty” (OFU).
  - “Posterior sampling” (PS)



# Regret bounds

- Greedy algorithms may never learn  $\rho(\pi_t^L) \rightarrow \rho(\phi^*)$ .
- Naïve exploration leads to regret exponential in  $|\mathcal{S}|, |\mathcal{A}|$ .
- Efficient algorithms must guide their exploration:
  - “Optimism in the face of uncertainty” (OFU).
  - “Posterior sampling” (PS)
- Efficient algorithms with Regret  $\tilde{O}(|\mathcal{S}|\sqrt{|\mathcal{A}|T})$ .
- Lower bound on Regret  $\Omega(\sqrt{|\mathcal{S}||\mathcal{A}|T})$ .
- But in many cases of interest  $\mathcal{S}, \mathcal{A}$  are huge...

# Regret bounds

- Greedy algorithms may never learn  $\rho(\pi_t^L) \rightarrow \rho(\phi^*)$ .
- Naïve exploration leads to regret exponential in  $|\mathcal{S}|, |\mathcal{A}|$ .
- Efficient algorithms must guide their exploration:
  - “Optimism in the face of uncertainty” (OFU).
  - “Posterior sampling” (PS)
- Efficient algorithms with Regret  $\tilde{O}(|\mathcal{S}|\sqrt{|\mathcal{A}|T})$ .
- Lower bound on Regret  $\Omega(\sqrt{|\mathcal{S}||\mathcal{A}|T})$ .
- But in many cases of interest  $\mathcal{S}, \mathcal{A}$  are huge... ☹

# Factored MDPs



Each transition only depends directly on a subset of the MDP.

## Factored MDPs

Let  $\mathcal{S} \times \mathcal{A} = \mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_n$ ,  $Z \subseteq [n]$  and  $\mathcal{X}[Z] := \bigotimes_{i \in Z} \mathcal{X}_i$

**Definition ( Factored reward functions  $R \in \mathcal{R} \subseteq \mathcal{P}_{\mathcal{X}, \mathbb{R}}^{C, \sigma}$ )**

$\mathcal{R}$  is factored over  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_n$  with scopes  $Z_1, \dots, Z_l \iff$ , for all  $R \in \mathcal{R}, x \in \mathcal{X}$  there exist,

$$\mathbb{E}[R(x)] = \sum_{i=1}^l \mathbb{E}[R_i(x[Z_i])]$$

with each  $r_i \sim R_i(x[Z_i])$  and individually observed.

**Definition ( Factored transition functions  $P \in \mathcal{P} \subseteq \mathcal{P}_{\mathcal{X}, \mathcal{S}}$ )**

$\mathcal{P}$  is factored over  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_n$  and  $\mathcal{S} = \mathcal{S}_1 \times \dots \times \mathcal{S}_m$  with scopes  $Z_1, \dots, Z_m \iff$ , for all  $P \in \mathcal{P}, x \in \mathcal{X}, s \in \mathcal{S}$  there exist,

$$P(s|x) = \prod_{i=1}^m P_i \left( s[i] \mid x[Z_i] \right)$$

## Factored MDPs

- Agent knows  $\mathcal{G} = (\{\mathcal{S}_i\}_{i=1}^m; \{\mathcal{X}_i\}_{i=1}^n; \{Z_i^R\}_{i=1}^l; \{Z_i^P\}_{i=1}^m)$ .
- Must learn  $(\{R_i\}_{i=1}^l; \{P_i\}_{i=1}^m)$  from experience.
- Algorithms ignoring  $\mathcal{G}$  lead to exponential regret bounds.
- $|\mathcal{S}| = \prod_{i=1}^m |\mathcal{S}_i| = |\mathcal{S}_1|^m, \quad |\mathcal{S}||\mathcal{A}| = \prod_{i=1}^n |\mathcal{X}_i| = |\mathcal{X}_1|^n$

## Factored MDPs

- Agent knows  $\mathcal{G} = (\{\mathcal{S}_i\}_{i=1}^m; \{\mathcal{X}_i\}_{i=1}^n; \{Z_i^R\}_{i=1}^l; \{Z_i^P\}_{i=1}^m)$ .
- Must learn  $(\{R_i\}_{i=1}^l; \{P_i\}_{i=1}^m)$  from experience.
- Algorithms ignoring  $\mathcal{G}$  lead to exponential regret bounds.
- $|\mathcal{S}| = \prod_{i=1}^m |\mathcal{S}_i| = |\mathcal{S}_1|^m, \quad |\mathcal{S}||\mathcal{A}| = \prod_{i=1}^n |\mathcal{X}_i| = |\mathcal{X}_1|^n$
- Efficient complexity bounds exist polynomial in  $|\mathcal{X}_i|, |\mathcal{S}_i|$
- No such results for regret. . .

## Factored MDPs

- Agent knows  $\mathcal{G} = (\{\mathcal{S}_i\}_{i=1}^m; \{\mathcal{X}_i\}_{i=1}^n; \{Z_i^R\}_{i=1}^l; \{Z_i^P\}_{i=1}^m)$ .
- Must learn  $(\{R_i\}_{i=1}^l; \{P_i\}_{i=1}^m)$  from experience.
- Algorithms ignoring  $\mathcal{G}$  lead to exponential regret bounds.
- $|\mathcal{S}| = \prod_{i=1}^m |\mathcal{S}_i| = |\mathcal{S}_1|^m, \quad |\mathcal{S}||\mathcal{A}| = \prod_{i=1}^n |\mathcal{X}_i| = |\mathcal{X}_1|^n$
- Efficient complexity bounds exist polynomial in  $|\mathcal{X}_i|, |\mathcal{S}_i|$
- No such results for regret. . . until now.

# Posterior Sampling for Reinforcement Learning

- 
- 
- 1: **Input:** Prior  $\phi$  encoding  $\mathcal{G}$ ,  $t = 1$
  - 2: **for** episodes  $k = 1, 2, \dots$  **do**
  - 3:   sample  $M_k \sim \phi(\cdot | H_t)$
  - 4:   compute  $\mu_k = \Gamma(M_k, \sqrt{\tau/k}) \leftarrow (\text{ADP planner})$
  - 5:   **for** timesteps  $j = 1, \dots, \tau$  **do**
  - 6:     sample and apply  $a_t = \mu_k(s_t, j)$
  - 7:     observe  $r_t^1, \dots, r_t^l$  and  $s_{t+1}^1, \dots, s_{t+1}^m$
  - 8:      $t = t + 1$
  - 9:   **end for**
  - 10: **end for**
-



# UCRL-Factored

- 
- 1: **Input:** Graph structure  $\mathcal{G}$ , confidence  $\delta$ ,  $t = 1$
  - 2: **for** episodes  $k = 1, 2, \dots$  **do**
  - 3:    $d_t^{R_i} = 4\sigma^2 \log(4l|\mathcal{X}[Z_i^R]|k/\delta)$  for  $i = 1, \dots, l$
  - 4:    $d_t^{P_j} = 4|\mathcal{S}_j| \log(4m|\mathcal{X}[Z_j^P]|k/\delta)$  for  $j = 1, \dots, m$
  - 5:    $\mathcal{M}_k = \{M \mid \mathcal{G}, \bar{R}_i \in \mathcal{R}_t^i(d_t^{R_i}), P_j \in \mathcal{P}_t^j(d_t^{P_j}) \forall i, j\} \leftarrow (\text{Confidence sets})$
  - 6:   compute  $\mu_k = \tilde{\Gamma}(\mathcal{M}_k, \sqrt{\tau/k}) \leftarrow (\text{ADP planner})$
  - 7:   **for** timesteps  $u = 1, \dots, \tau$  **do**
  - 8:     sample and apply  $a_t = \mu_k(s_t, u)$
  - 9:     observe  $r_t^1, \dots, r_t^l$  and  $s_{t+1}^1, \dots, s_{t+1}^m$
  - 10:     $t = t + 1$
  - 11:   **end for**
  - 12: **end for**
-

## Main results

### Theorem (Expected regret for PSRL in factored MDPs)

*If the prior  $\phi$  is the distribution of  $M^*$  and  $\Psi$  is the span of the optimal value function:*

$$\mathbb{E} [\text{Regret}(T, \pi_{\tau}^{\text{PS}}, M^*)] = \tilde{O} \left( \sigma \sum_{i=1}^{d_1} \sqrt{|\mathcal{X}[Z_i^R]| T} + \mathbb{E}[\Psi] \sum_{j=1}^{d_2} \sqrt{|\mathcal{X}[Z_j^P]| |\mathcal{S}_j| T} \right) \quad (1)$$

### Theorem (High probability regret for UCRL-Factored)

*If  $D$  is the diameter of  $M^*$ , then for any  $M^*$  can bound the regret of UCRL-Factored:*

$$\text{Regret}(T, \pi_{\tau}^{\text{UC}}, M^*) = \tilde{O} \left( \sigma \sum_{i=1}^{d_1} \sqrt{|\mathcal{X}[Z_i^R]| T} + CD \sum_{j=1}^{d_2} \sqrt{|\mathcal{X}[Z_j^P]| |\mathcal{S}_j| T} \right) \quad (2)$$

*with probability at least  $1 - \delta$*

## Clean bounds in the symmetric case

Let  $\mathcal{Q}$  be shorthand for the structure  $\mathcal{G}$  such that  $l + 1 = m$ ,  $C = \sigma = 1$ ,  $|\mathcal{S}_i| = |\mathcal{X}_i| = K$  and  $|Z_i^R| = |Z_i^P| = \zeta$  for all suitable  $i$  and write  $J = K^\zeta$ . In this case  $\Psi, D \leq \tau$  trivially.

### Corollary (Clean bounds for PSRL)

$$\mathbb{E} \left[ \text{Regret}(T, \pi_\tau^{\text{PS}}, M^*) \right] \leq 15m\tau \sqrt{JKT \log(2mJT)} \quad (3)$$

### Corollary (Clean bounds for UCRL-Factored)

$$\text{Regret}(T, \pi_\tau^{\text{UC}}, M^*) \leq 15m\tau \sqrt{JKT \log(12mJT/\delta)} \quad (4)$$

*with probability at least  $1 - \delta$ .*

## Clean bounds in the symmetric case

Let  $\mathcal{Q}$  be shorthand for the structure  $\mathcal{G}$  such that  $l + 1 = m$ ,  $C = \sigma = 1$ ,  $|\mathcal{S}_i| = |\mathcal{X}_i| = K$  and  $|\mathcal{Z}_i^R| = |\mathcal{Z}_i^P| = \zeta$  for all suitable  $i$  and write  $J = K^\zeta$ . In this case  $\Psi, D \leq \tau$  trivially.

The key point is that we go from  $\mathcal{G}$ -agnostic

$$\tilde{O}(|\mathcal{S}| \sqrt{|\mathcal{A}| T}) = \tilde{O}(\sqrt{J^{m/\zeta} K^m T})$$

to the new bounds

$$\tilde{O}\left(\sum_{j=1}^m \sqrt{|\mathcal{X}[\mathcal{Z}_j^P]| |\mathcal{S}_j| T}\right) = \tilde{O}(m \sqrt{J K T})$$

which can be exponentially tighter.

## Analysis outline

We consider  $\text{Regret}(T) = \sum_{k=1}^m \Delta_k$ , regret within each episode.

$$\Delta_k = V_{*,1}^*(s) - V_{k,1}^*(s) = \left( V_{k,1}^k(s) - V_{k,1}^*(s) \right) + \left( V_{*,1}^*(s) - V_{k,1}^k(s) \right) \quad (5)$$

Where  $V_{\mu,1}^M(s)$  is the value of employing policy  $\mu$  on MDP  $M$  for one episode. We write  $*, k$  as shorthand for the optimal and algorithmic MDPs at each stage.

$$(V_{k,1}^k - V_{k,1}^*)(s_{t_k+1}) = \sum_{i=1}^{\tau} (\mathcal{T}_{k,i}^k - \mathcal{T}_{k,i}^*) V_{k,i+1}^k(s_{t_k+i}) + \sum_{i=1}^{\tau} d_{t_k+i}. \quad (6)$$

Where  $d_t$  is a bounded martingale difference and the first term  $A$ :

$$A \leq \sum_{i=1}^{\tau} |\bar{R}^k(x_{k,i}) - \bar{R}^*(x_{k,i})| + \frac{1}{2} \Psi_k \|P^k(\cdot|x_{k,i}) - P^*(\cdot|x_{k,i})\|_1 \quad (7)$$

# Key lemma

## Lemma (Bounding factored deviations)

*Let the transition function class  $\mathcal{P} \subseteq \mathcal{P}_{\mathcal{X}, \mathcal{S}}$  be factored over  $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_n$  and  $\mathcal{S} = \mathcal{S}_1 \times \dots \times \mathcal{S}_m$  with scopes  $Z_1, \dots, Z_m$ . Then, for any  $P, \tilde{P} \in \mathcal{P}$  we may bound their L1 distance by the sum of the differences of their factorizations:*

$$\|P(x) - \tilde{P}(x)\|_1 \leq \sum_{i=1}^m \|P_i(x[Z_i]) - \tilde{P}_i(x[Z_i])\|_1$$

Using Azuma-Hoeffding with a union bound we can bound the Bellman error in terms of a concentration which depends on  $|\mathcal{X}[Z_i^P]|$  as opposed to  $|\mathcal{X}|$ . From the previous slide this gives us bounds on regret.

# References

Please see arXiv for a full version of the paper.

# References

Please see arXiv for a full version of the paper.

Thanks!