

Blackwell Optimality in Robust MDPs

Niels van Duijl

Vineet Goyal and Julien Grand-Clément, Robust Markov Decision Process: Beyond Rectangularity, 2023

Julien Grand-Clement, Marek Petrik, Nicolas Vieille, Beyond discounted returns: Robust Markov decision processes with average and Blackwell optimality, 2024

Contents

- **Introduction**
- **Main Paper: Robust Markov Decision Process: Beyond Rectangularity**
- **Second Paper: Beyond discounted returns: Robust Markov decision processes with average and Blackwell optimality**

Introduction

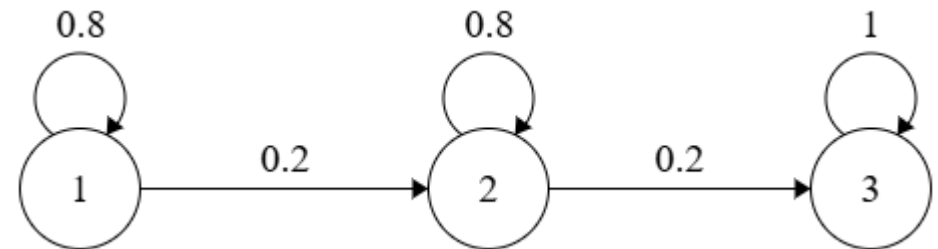
Markov Chains

- Discrete-time stochastic process
- State space \mathcal{S} , time horizon T

$$\{X_t, t = 0, 1, 2, \dots, T\} \quad X_t \in \mathcal{S}, \forall t$$

- Transition probabilities P
- **Markov Chain** when P only depends on the current state.

$$P\{X_{t+1} = j \mid X_t = i, X_{t-1} = i_{t-1}, \dots, X_1 = i_1, X_0 = i_0\} = P_{ij}$$



Introduction

Markov Decision Process (MDP)

- Markov Chain, but with an action space \mathbf{A}
- P_{ij} now depends on the chosen action as well
- Each state-action pair has a reward $r_{s,a}$
- Choose actions based on a policy π
- Stationary policy:

$$\pi = \{\pi_s(a), a \in A, s \in S\}$$

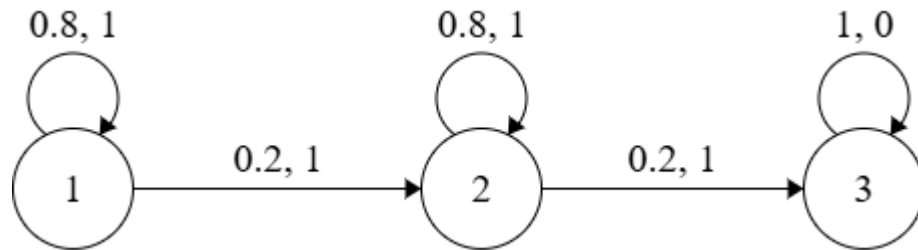
- Deterministic when $\pi_s(a) \in \{0, 1\}$

Introduction

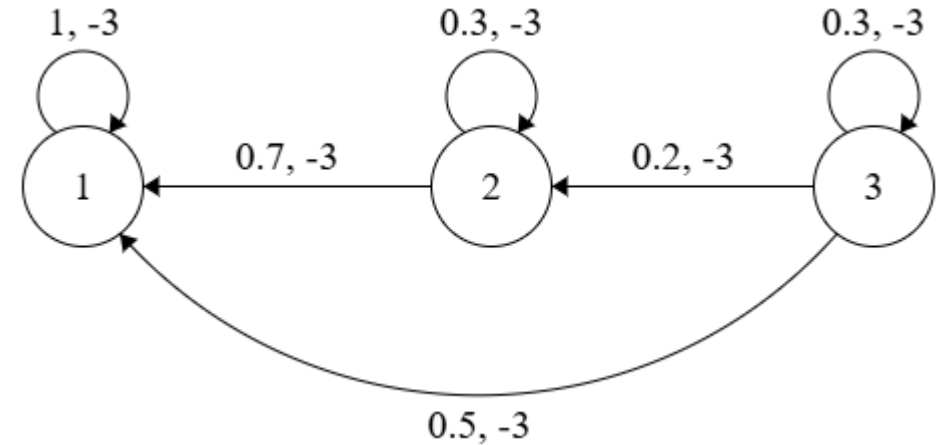
MDP Example

- We now have rewards (second label on edge)
- We have *wait* and *repair* as actions

- P_{ij} when choosing *wait*:



- P_{ij} when choosing *repair*:



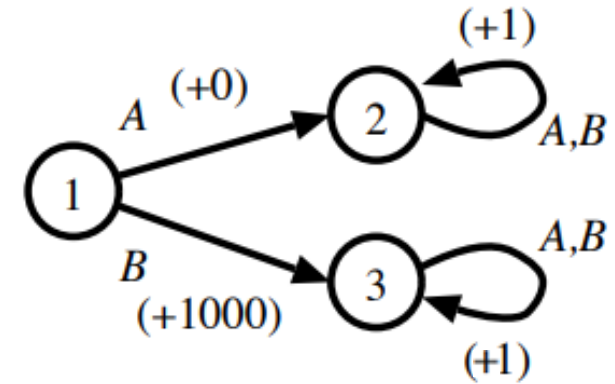
Introduction

Bellman Equations

- Computing the optimal policy
- How to deal with infinite time horizon?
- Discounted vs average reward:

$$v(\pi^*, s) = \max_{a \in A} \{ r_{s,a} + \lambda \sum_{s' \in S} P(s' | s, a) v(\pi^*, s') \} \quad \lambda \in (0, 1).$$

$$v(\pi^*, s) = \max_{a \in A} \{ r_{s,a} + \sum_{s' \in S} P(s' | s, a) v(\pi^*, s') \} - g(\pi^*).$$



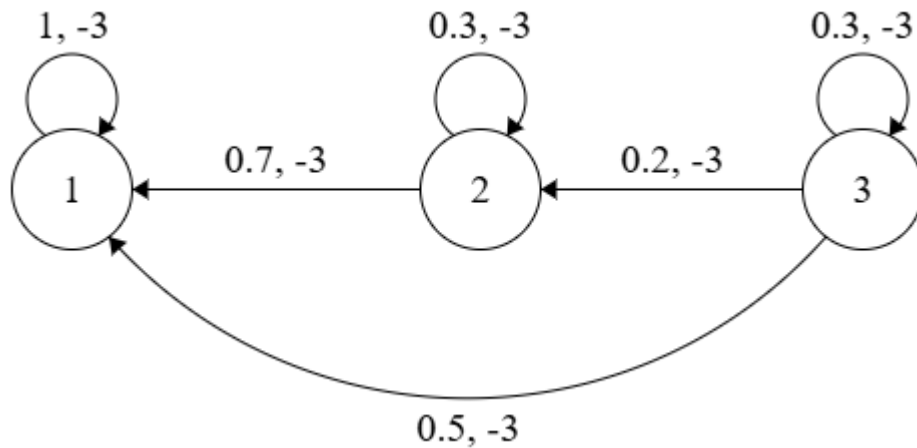
- A policy π is **Blackwell optimal** if it is optimal for all λ close enough to 1
- These are also optimal policies for the average reward counterpart.
- We will use discounted reward for the rest of the talk

Graph Source: Schwartz, A. (1993). *A Reinforcement Learning Method for Maximizing Undiscounted Rewards*. 298–305. <https://doi.org/10.1016/B978-1-55860-307-3.50045-9>

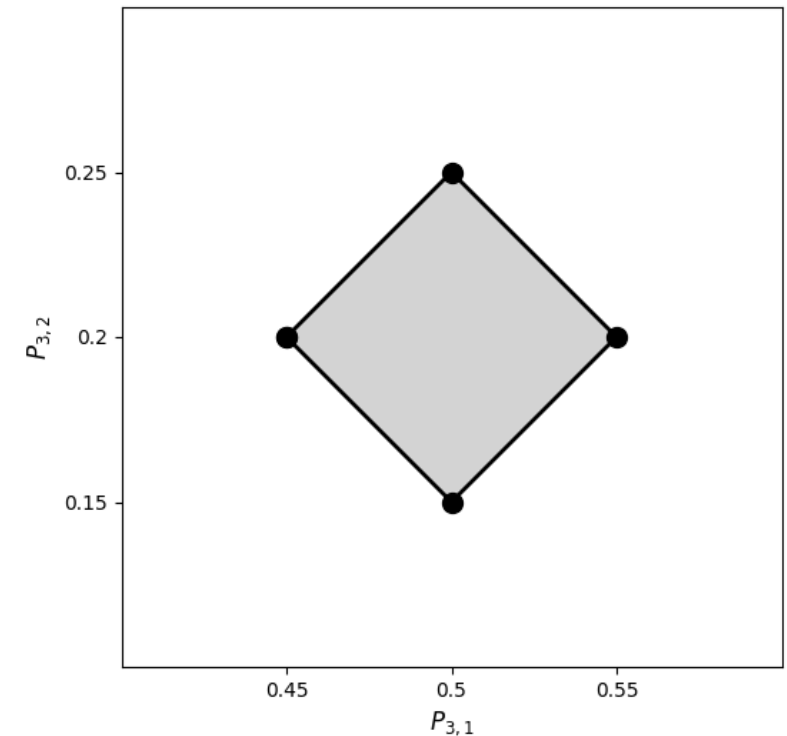
Introduction

Robust MDPs

- The probabilities are now uncertain, picked from an uncertainty set \mathbb{P}
- Knowing exact probabilities is hard when deriving from data
- $(P_{3,3} = 1 - P_{3,2} - P_{3,1})$
- P_{ij} when choosing *repair*:



- P_3 when choosing *repair*:



Robust MDPs

- The probabilities are now uncertain, picked from an uncertainty set \mathbb{P}
- Knowing exact probabilities is hard when deriving from data
- We consider the worst-case probabilities

$$\max_{\pi \in \Pi^G} \min_{\mathbf{P} \in \mathbb{P}} R(\pi, \mathbf{P}).$$

- Adversary MDP: Second player that controls the factor matrix that tries to minimize our reward
- Two-player game of (normal) MDPs.

Rectangularity

- Generally, it is NP-hard to compute optimal policy
- Independence assumptions are needed: this is rectangularity
- Neglecting constraints only makes the worst-case worse.

Introduction

(s,a)-rectangularity, s-rectangularity

- (s,a)-rect: Each \mathbf{P} can be chosen from the uncertainty set independently of others
- Optimal policy that is stationary and deterministic, and can be computed efficiently

$$\mathbb{P} = \prod_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathbb{P}_{sa}, \quad \mathbb{P}_{sa} \subseteq \mathbb{R}_+^{|\mathcal{S}|}.$$

- s-rect: Probs may be dependent on probs for different actions in the same state. Still independent between states.
- Here the optimal policy is stationary but may not be deterministic

$$\mathbb{P} = \prod_{s \in \mathcal{S}} \mathbb{P}_s, \quad \mathbb{P}_s \subseteq \mathbb{R}_+^{|\mathcal{S}| \times |\mathcal{A}|}.$$

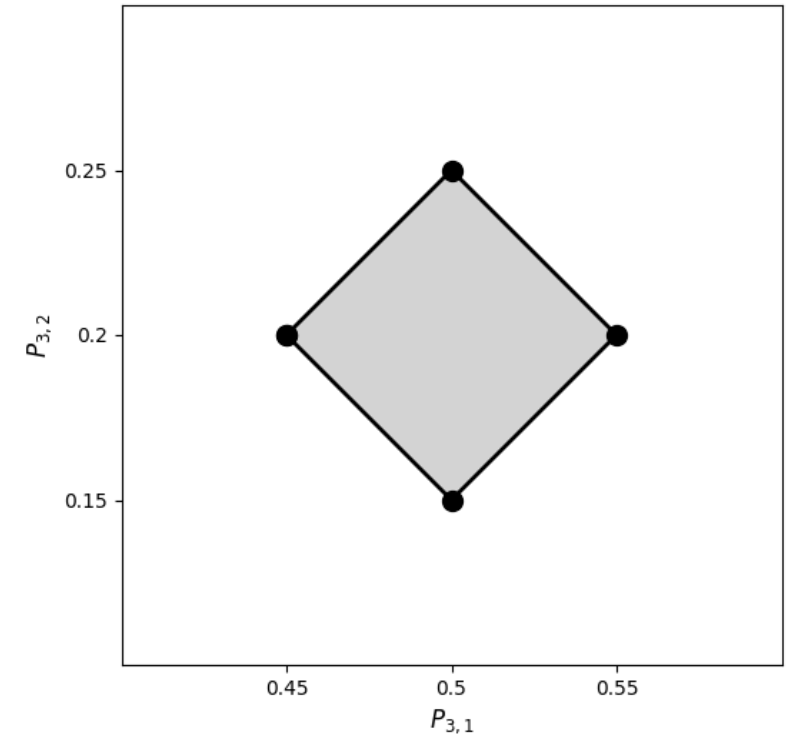
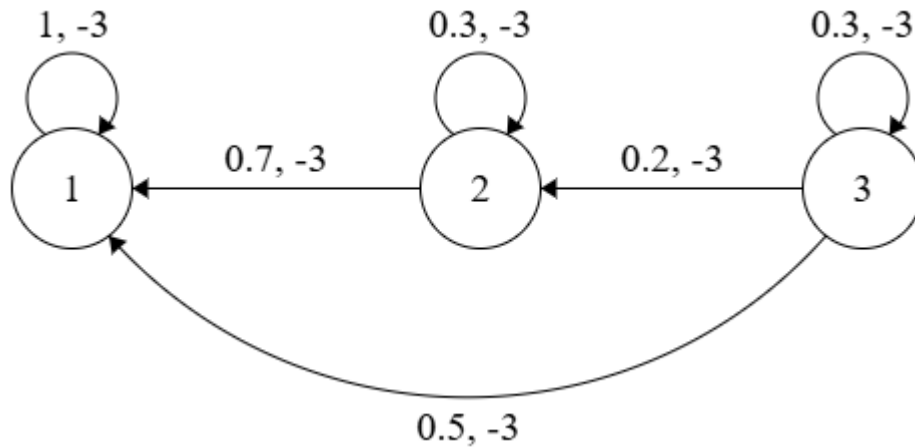
Introduction

(s,a)-rectangularity, s-rectangularity

- (s,a)-rect: Each \mathbf{P} can be chosen from the uncertainty set independently of others
- s-rect: Probs may be dependent on probs for different actions in the same state.

- \mathbf{P}_3 when choosing *repair*:

- \mathbf{P}_{ij} when choosing *repair*:



Paper 1: Main Contributions

- **Introduce a new type of rectangularity**
- **Min-max duality:**

$$\max_{\pi \in \Pi} \min_{\mathbf{P} \in \mathcal{P}} R(\pi, \mathbf{P}) = \min_{\mathbf{P} \in \mathcal{P}} \max_{\pi \in \Pi} R(\pi, \mathbf{P}).$$

- **Algorithm to compute the optimal policy**
- **Blackwell optimality**

A new type of uncertainty set

- Idea: common underlying factors (healthcare)
- Fixed Coefficients \mathbf{u} , factors \mathbf{w} themselves are uncertain
- Each factor is a probability distribution over the next state
- r is not the reward!

$$\mathbb{P} = \left\{ \left(\sum_{i=1}^r u_{sa}^i w_{i,s'} \right)_{sas'} \mid \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_r) \in \mathcal{W} \subseteq \mathbb{R}^{S \times r} \right\}$$

$$\sum_{i=1}^r u_{sa}^i = 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A}, \sum_{s'=1}^S w_{i,s'} = 1, \forall i \in [r],$$

A new type of rectangularity

- r-rectangularity: when the factors are independent

$$\mathcal{W} = \mathcal{W}^1 \times \dots \times \mathcal{W}^r, \text{ where } \mathcal{W}^1, \dots, \mathcal{W}^r \subset \mathbb{R}_+^S.$$

- (s,a)-rect \rightarrow r-rect
- s-rect and r-rect not related

$$\mathbb{P} = \left\{ \left(\sum_{i=1}^r u_{sa}^i w_{i,s'} \right)_{sas'} \mid \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_r) \in \mathcal{W} \subseteq \mathbb{R}^{S \times r} \right\}$$

Assumption 2.4

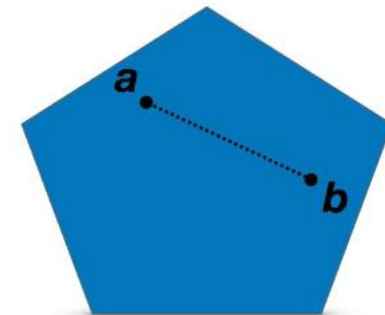
- Used in most results of the rest of the paper
- The uncertainty sets need to be convex compact.

Assumption 2.4 *The sets $\mathcal{W}^1, \dots, \mathcal{W}^r$ are convex compact. Moreover, for any $i \in [r]$, for any cost vector $\mathbf{c} \in \mathbb{R}^S$, we can find an ϵ -optimal solution of*

$$\min_{\mathbf{w}_i \in \mathcal{W}^i} \mathbf{c}^\top \mathbf{w}_i$$

in $O(\text{comp}(\mathcal{W}^i) \log(\epsilon^{-1}))$ arithmetic operations, where $\text{comp}(\mathcal{W}^i) \in \mathbb{R}$ depends on the structure of the uncertainty set \mathcal{W}^i .

- Compact: closed and bounded
- Convex: Every line segment is contained in the set



Paper 1: Evaluating a policy
Adversary MDP

- Adversary MDP: r States, S actions, W is the policy

$$Prob(i \xrightarrow{\text{action } s} j) = \sum_{a \in \mathbb{A}} \pi_{sa} u_{sa}^j, \quad Reward(i, \text{action } s) = \sum_{a \in \mathbb{A}} \pi_{sa} r_{sa}.$$

$$\mathbf{T}_\pi = \left(\sum_{a \in \mathbb{A}} \pi_{sa} u_{sa}^i \right)_{(s,i) \in \mathbb{S} \times [r]} \in \mathbb{R}_+^{S \times r}.$$

- Bellman equation of this adversary gives us value function β for the adversary:

$$v_s^\pi = \sum_{a \in \mathbb{A}} \pi_{sa} (r_{sa} + \lambda \mathbf{P}_{sa}^\top \mathbf{v}^\pi), \quad \forall s \in \mathbb{S}, \quad \longrightarrow \quad \beta_i = \mathbf{w}_i^\top (\mathbf{r}_\pi + \lambda \cdot \mathbf{T}_\pi \beta), \quad \forall i \in [r].$$

- R-rect is necessary to allow the adversary to independently optimize each vector.

$$\mathbb{P} = \left\{ \left(\sum_{i=1}^r u_{sa}^i w_{i,s'} \right)_{sas'} \mid \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_r) \in \mathcal{W} \subseteq \mathbb{R}^{S \times r} \right\}$$

Paper 1: Min-Max Duality

Theorem of Duality

- Lemma 4.1: If \mathbb{P} is r -rectangular and the sets $\mathcal{W}^1, \dots, \mathcal{W}^r$ are convex compact, there exists a stationary, deterministic optimal policy.
- Duality Theorem:

Theorem 4.2 *Under Assumption 2.4, let (π^*, \mathbf{W}^*) be a solution to the robust MDP problem (1.2) with r -rectangular uncertainty set, with π^* deterministic. Then*

$$\mathbf{W}^* \in \arg \min_{\mathbf{W} \in \mathcal{W}} R(\pi^*, \mathbf{W}) \text{ and } \pi^* \in \arg \max_{\pi \in \Pi} R(\pi, \mathbf{W}^*). \quad (4.6)$$

Moreover, the following strong min-max duality holds.

$$\max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}) = \min_{\mathbf{W} \in \mathcal{W}} \max_{\pi \in \Pi} R(\pi, \mathbf{W}). \quad (4.7)$$

- Result: (π^*, \mathbf{W}^*) is an equilibrium in the two-player game.
- Does not hold for s -rectangular uncertainty sets.

Robust Maximum Principle

- Maximum principle: the optimal policy attains the highest value regardless of starting state.
- Robust equivalent:

Proposition 6.1 *Under Assumption 2.4, let \mathbb{P} be an r -rectangular uncertainty set.*

1. *Let π be a policy and $\mathbf{W}^1 \in \arg \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W})$. Then*

$$v_s^{\pi, \mathbf{W}^1} \leq v_s^{\pi, \mathbf{W}^0}, \quad \forall \mathbf{W}^0 \in \mathcal{W}, \forall s \in \mathcal{S}.$$

2. *Let $(\pi^*, \mathbf{W}^*) \in \arg \max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W})$. Then*

$$\forall \pi \in \Pi, \forall \mathbf{W}^1 \in \arg \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}), v_s^{\pi, \mathbf{W}^1} \leq v_s^{\pi^*, \mathbf{W}^*}, \quad \forall s \in \mathcal{S}.$$

- Follows from the Duality Theorem

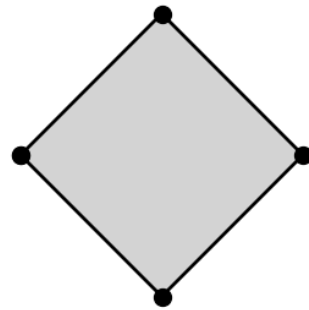
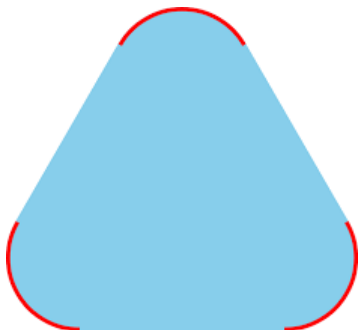
Paper 1: Blackwell Optimality

Blackwell Optimality

- Recall: A policy π is **Blackwell optimal** if it is optimal for all λ close enough to 1
- Proven for r -rect sets with finitely many extreme points.

Proposition 6.2 *Let \mathbb{P} be an r -rectangular uncertainty set. Assume that the sets $\mathcal{W}^1, \dots, \mathcal{W}^r$ have finitely many extreme points. Then there exists a stationary deterministic policy π^* , a factor matrix \mathbf{W}^* , and a discount factor $\lambda_0 \in (0, 1)$, such that for all $\lambda \in (\lambda_0, 1)$, the pair (π^*, \mathbf{W}^*) remains an optimal solution to the robust MDP problem (1.2).*

- Proof idea: both \mathbf{W} and Π have finitely many extreme points.
- Extreme point: A point which does not lie in any open line segment joining two points in the set.



Proof of Blackwell Optimality

Lemma G.1 *Let $(u_n)_{n \geq 0}$ be a sequence with values in a finite set \mathcal{E} . Then there exists an element $e \in \mathcal{E}$ that is attained infinitely often by $(u_n)_{n > 0}$.*

- So, we can construct:

$$\lambda_n \rightarrow 1, \text{ and } (\pi^*, \mathbf{W}^*) \in \arg \max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}, \lambda_n), \forall n \geq 0.$$

- By contradiction, suppose we can also construct:

$$\gamma_n \rightarrow 1, \text{ and } (\pi^*, \mathbf{W}^*) \notin \arg \max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}, \gamma_n), \forall n \geq 0.$$

- Take $(\tilde{\pi}, \tilde{\mathbf{W}})$ to be optimal for all γ_n

- Since (π^*, \mathbf{W}^*) not optimal, we can find: $v_{\gamma_n, x_1}^{\pi^*, \mathbf{W}^{**}} < v_{\gamma_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}}$

$$v_{\gamma_n, x_1}^{\pi^*, \mathbf{W}^{**}} < v_{\gamma_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}}, v_{\lambda_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}} \leq v_{\lambda_n, x_1}^{\pi^*, \mathbf{W}^{**}}.$$

- Continuous, rational function that takes on the value 0 an infinite number of times:

$$f : (0, 1) \rightarrow \mathbb{R}, f(t) = v_{t, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}} - v_{t, x_1}^{\pi^*, \mathbf{W}^{**}}.$$

Proposition 6.3

- This proof works for any λ !
- We conclude the following proposition:

Proposition 6.3 *Let \mathbb{P} be an r -rectangular uncertainty set. Assume that the sets $\mathcal{W}^1, \dots, \mathcal{W}^r$ have finitely many extreme points. Then there exists an integer $p \in \mathbb{N}$, there exists some scalars $\lambda_0 = 0 < \lambda_1 < \dots < \lambda_p = 1$ such that for all $j \in \{0, \dots, p-1\}$, the same pair of stationary deterministic policy and factor matrix (π_j, \mathbf{W}_j) is an optimal solution to the robust MDP problem (1.2) for all $\lambda \in (\lambda_j, \lambda_{j+1})$.*

The rest of the paper

- Algorithm to compute the optimal policy:

$$\mathbf{v}^0 = \mathbf{0}, v_s^{k+1} = \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \min_{\mathbf{w}_i \in \mathcal{W}^i} \mathbf{w}_i^\top \mathbf{v}^k \right\}, \forall s \in \mathbb{S}, \forall k \geq 0.$$

- Numerical Experiments
- Conclusion: R-rectangularity outperforms s-rectangularity

Paper 2: Summary

- Average optimality for RMDPs
- **Blackwell optimality for RMDPs**
- Algorithms to compute the optimal average reward
- Numerical experiments
- **Does not talk about r-rect.**

Uncertainty set \mathcal{U}	Discount optimality	Average optimality	Blackwell optimality
Singleton (MDPs)	stationary, deterministic	stationary, deterministic	stationary, deterministic
sa-rectangular, compact, convex	stationary, deterministic	stationary, deterministic	<ul style="list-style-type: none"> • may not exist • $\exists \pi$ stationary deterministic, π ϵ-Blackwell optimal, $\forall \epsilon > 0$ • π also average optimal
sa-rectangular, compact, convex, definable	stationary, deterministic	stationary, deterministic	<ul style="list-style-type: none"> • stationary, deterministic • π also average optimal
s-rectangular, compact convex	stationary, randomized	history-dependent, randomized	may not exist

Motivation

- Less assumptions made than in previous work
- Again, uncertainty set is assumed to be convex, compact.
- Adversary MDP makes the problem similar to stochastic games
- However, the adversary is restricted to stationary policies

Introduction

- Reminder: A policy π is **Blackwell optimal** if it is optimal for all discount factors close enough to 1

$$\inf_{P \in \mathcal{U}} R_\gamma(\pi, P) \geq \sup_{\pi' \in \Pi_H} \inf_{P \in \mathcal{U}} R_\gamma(\pi', P), \quad \forall \gamma \in (\gamma_0, 1).$$

- ϵ -Blackwell optimality:

$$\min_{P \in \mathcal{U}} (1 - \gamma) R_\gamma(\pi, P) \geq \sup_{\pi' \in \Pi_S} \min_{P \in \mathcal{U}} (1 - \gamma) R_\gamma(\pi', P) - \epsilon, \quad \forall \gamma \in (\gamma_\epsilon, 1).$$

- Normalised and a difference of ϵ is allowed

General (s,a)-rect uncertainty sets

Theorem 4.4 *There exists an sa-rectangular robust MDP instance, with a compact convex uncertainty set \mathcal{U} , and with no Blackwell optimal policy:*

$$\forall \pi \in \Pi_{\text{H}}, \forall \gamma \in (0, 1), \exists \gamma' \in (\gamma, 1), \min_{\mathbf{P} \in \mathcal{U}} R_{\gamma'}(\pi, \mathbf{P}) < \sup_{\pi' \in \Pi_{\text{S}}} \min_{\mathbf{P} \in \mathcal{U}} R_{\gamma'}(\pi', \mathbf{P}).$$

- Based on two distinct sets $\mathcal{U}_{s_0 a_1}$ $\mathcal{U}_{s_0 a_2}$, of which the boundaries intersect infinitely often.
- However: we can find a policy that is ϵ -Blackwell optimal for every $\epsilon > 0$:

Theorem 4.5 *Let \mathcal{U} be an sa-rectangular compact uncertainty set. Then there exists a stationary deterministic policy that is ϵ -Blackwell optimal for all $\epsilon > 0$, i.e., $\exists \pi \in \Pi_{\text{SD}}, \forall \epsilon > 0, \exists \gamma_{\epsilon} \in (0, 1)$ such that*

$$\min_{\mathbf{P} \in \mathcal{U}} (1 - \gamma) R_{\gamma}(\pi, \mathbf{P}) \geq \sup_{\pi' \in \Pi_{\text{S}}} \min_{\mathbf{P} \in \mathcal{U}} (1 - \gamma) R_{\gamma}(\pi', \mathbf{P}) - \epsilon, \forall \gamma \in (\gamma_{\epsilon}, 1).$$

Paper 2: Definability

Definability

- A subset of \mathbb{R}^n is definable if it is of the form:

$$\{\mathbf{x} \in \mathbb{R}^n \mid \exists k \in \mathbb{N}, \exists \mathbf{y} \in \mathbb{R}^k, P(x_1, \dots, x_n, y_1, \dots, y_k, \exp(x_1), \dots, \exp(x_n), \dots, \exp(y_1), \dots, \exp(y_k)) = 0\}$$

- A function $f : \Omega \rightarrow \mathbb{R}^m$, $\Omega \subset \mathbb{R}^n$ is definable if its graph is definable:

$$\{(\mathbf{x}, \mathbf{y}) \in \Omega \times \mathbb{R}^m \mid \mathbf{y} = f(\mathbf{x})\}$$

- Intuitively, a set is definable if it is constructed based on polynomials, the exponential function, and canonical projections (elimination of variables).

- Simple example:

$$x \geq 0 \longrightarrow P(x, y) = x - y^2$$

- Lemma 4.15**
1. *The only definable subsets of \mathbb{R} are the finite union of open intervals and singletons.*
 2. *If $A, B \subset \mathbb{R}^n$ are definable sets, then $A \cup B$, $A \cap B$ and $\mathbb{R}^n \setminus A$ are definable sets.*
 3. *Let f, g be definable functions. Then $f \circ g$, $-f$, $f + g$, $f \times g$ are definable.*
 4. *For A, B two definable sets and $g : A \times B \rightarrow \mathbb{R}$ a definable function, then the functions $\mathbf{x} \mapsto \inf_{\mathbf{y} \in B} g(\mathbf{x}, \mathbf{y})$ and $\mathbf{x} \mapsto \sup_{\mathbf{y} \in B} g(\mathbf{x}, \mathbf{y})$ (defined over A) are definable functions.*
 5. *If $A, B \subseteq \mathbb{R}$ and $g : A \rightarrow \mathbb{R}$ are definable then $g^{-1}(B)$ is a definable set.*

An example

$$\{\mathbf{x} \in \mathbb{R}^n \mid \exists k \in \mathbb{N}, \exists \mathbf{y} \in \mathbb{R}^k, P(x_1, \dots, x_n, y_1, \dots, y_k, \exp(x_1), \dots, \exp(x_n), \dots, \exp(y_1), \dots, \exp(y_k)) = 0\}$$

Example 4.16 (ℓ_p -norms are definable.) Consider an ℓ_p -norm for $p \in \mathbb{N}$. Then its graph is $\{(\mathbf{x}, y) \in \mathbb{R}^S \times \mathbb{R} \mid \sum_{s \in S} |x_s|^p - y^p = 0, y \geq 0\}$, which is a definable set. Therefore, ℓ_p -norms are definable.

- They can have infinitely many extreme points but are definable!

Definable (s,a)-rect uncertainty sets

Theorem 4.17 (Theorem 2.1, Coste [2000]) *Let $f : (a, b) \rightarrow \mathbb{R}$ be a definable function. Then there exists a finite subdivision of the interval (a, b) as $a = a_1 < a_2 < \dots < a_k = b$ such that on each (a_i, a_{i+1}) for $i = 1, \dots, k - 1$, f is continuous and either constant or strictly monotone.*

Proposition 4.19 *Assume that \mathcal{U} is sa-rectangular and definable. Then for any policy $\pi \in \Pi_S$, the function $\gamma \mapsto v_{\gamma}^{\pi, \mathcal{U}}$ is a definable function.*

Theorem 4.25 *Consider an sa-rectangular robust MDP with a definable compact uncertainty set \mathcal{U} . Then there exists a stationary deterministic Blackwell optimal policy:*

$$\exists \pi \in \Pi_{SD}, \exists \gamma_0 \in (0, 1), \forall \gamma \in (\gamma_0, 1), \min_{P \in \mathcal{U}} R_{\gamma}(\pi, P) \geq \sup_{\pi' \in \Pi_S} \min_{P \in \mathcal{U}} R_{\gamma}(\pi', P).$$

- Proof based on $v_{\gamma, s}^{\pi, \mathcal{U}} - v_{\gamma, s}^{\pi', \mathcal{U}}$

Proof of Blackwell Optimality

Lemma G.1 *Let $(u_n)_{n \geq 0}$ be a sequence with values in a finite set \mathcal{E} . Then there exists an element $e \in \mathcal{E}$ that is attained infinitely often by $(u_n)_{n > 0}$.*

- So, we can construct:

$$\lambda_n \rightarrow 1, \text{ and } (\pi^*, \mathbf{W}^*) \in \arg \max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}, \lambda_n), \forall n \geq 0.$$

- By contradiction, suppose we can also construct:

$$\gamma_n \rightarrow 1, \text{ and } (\pi^*, \mathbf{W}^*) \notin \arg \max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}, \gamma_n), \forall n \geq 0.$$

- Take $(\tilde{\pi}, \tilde{\mathbf{W}})$ to be optimal for all γ_n

- Since (π^*, \mathbf{W}^*) not optimal, we can find: $v_{\gamma_n, x_1}^{\pi^*, \mathbf{W}^{**}} < v_{\gamma_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}}$

$$v_{\gamma_n, x_1}^{\pi^*, \mathbf{W}^{**}} < v_{\gamma_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}}, v_{\lambda_n, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}} \leq v_{\lambda_n, x_1}^{\pi^*, \mathbf{W}^{**}}.$$

- Continuous, rational function that takes on the value 0 an infinite number of times:

$$f : (0, 1) \rightarrow \mathbb{R}, f(t) = v_{t, x_1}^{\tilde{\pi}, \tilde{\mathbf{W}}} - v_{t, x_1}^{\pi^*, \mathbf{W}^{**}}.$$

The rest of the paper: average optimality results

- Blackwell optimal policies also average optimal
- Algorithms to compute the optimal gain, but not the actual policies
- No polynomial-time algorithm known for (s,a)-rect RMDPs.
- Experiments on the newly defined algorithms.

Uncertainty set \mathcal{U}	Discount optimality	Average optimality	Blackwell optimality
Singleton (MDPs)	stationary, deterministic	stationary, deterministic	stationary, deterministic
sa-rectangular, compact, convex	stationary, deterministic	stationary, deterministic	<ul style="list-style-type: none"> • may not exist • $\exists \pi$ stationary deterministic, π ϵ-Blackwell optimal, $\forall \epsilon > 0$ • π also average optimal
sa-rectangular, compact, convex, definable	stationary, deterministic	stationary, deterministic	<ul style="list-style-type: none"> • stationary, deterministic • π also average optimal
s-rectangular, compact convex	stationary, randomized	history-dependent, randomized	may not exist

Questions?

Introduction

$$R_{\text{avg}}(\pi, \mathbf{P}) = \mathbb{E}_{\pi, \mathbf{P}} \left[\limsup_{T \rightarrow +\infty} \frac{1}{T+1} \sum_{t=0}^T r_{s_t a_t s_{t+1}} \mid s_0 \sim \mathbf{p}_0 \right].$$

- To solve: $\sup_{\pi \in \Pi_{\text{H}}} \inf_{\mathbf{P} \in \mathcal{U}} R_{\text{avg}}(\pi, \mathbf{P})$
- Inf is used instead of min because it might not exist

Proposition 3.2 *There exists a robust MDP instance with an sa-rectangular compact convex uncertainty set, for which $\inf_{\mathbf{P} \in \mathcal{U}} R_{\text{avg}}(\pi, \mathbf{P})$ is not attained for any $\pi \in \Pi_{\text{S}}$.*

(s,a)-rectangular uncertainty sets

- The optimal policy is stationary and deterministic:

$$\sup_{\pi \in \Pi_H} \inf_{P \in \mathcal{U}} R_{\text{avg}}(\pi, P) = \max_{\pi \in \Pi_{SD}} \inf_{P \in \mathcal{U}} R_{\text{avg}}(\pi, P).$$

- Strong duality also holds.

$$\sup_{\pi \in \Pi_H} \inf_{P \in \mathcal{U}} R_{\text{avg}}(\pi, P) = \inf_{P \in \mathcal{U}} \sup_{\pi \in \Pi_H} R_{\text{avg}}(\pi, P)$$

$$\max_{\pi \in \Pi_{SD}} \inf_{P \in \mathcal{U}} R_{\text{avg}}(\pi, P) = \inf_{P \in \mathcal{U}} \max_{\pi \in \Pi_{SD}} R_{\text{avg}}(\pi, P)$$

- From these results, it also follows that: $\sup_{\pi \in \Pi_H} \inf_{P \in \mathcal{U}_H} R_{\text{avg}}(\pi, P) = \sup_{\pi \in \Pi_H} \inf_{P \in \mathcal{U}} R_{\text{avg}}(\pi, P).$
- So we can justifiably restrict the Adversary to stationary policies

Paper 2: Average optimality

S-rectangular uncertainty sets

- History-dependent policies may be optimal
- Not true for discounted s-rect, or avg (s,a)-rect!

Proposition 2.5

- Under Assumption 2.4, if r-rect, there exists an optimal stationary policy
- Proven using the duality result of later in the paper. $\max_{\pi \in \Pi_S} \min_{\mathbf{P} \in \mathbb{P}} R(\pi, \mathbf{P}) = \min_{\mathbf{P} \in \mathbb{P}} \max_{\pi \in \Pi_S} R(\pi, \mathbf{P})$.

$$\max_{\pi \in \Pi_S} \min_{\mathbf{P} \in \mathbb{P}} R(\pi, \mathbf{P}) \leq \max_{\pi \in \Pi} \min_{\mathbf{P} \in \mathbb{P}} R(\pi, \mathbf{P})$$

$$\leq \min_{\mathbf{P} \in \mathbb{P}} \max_{\pi \in \Pi} R(\pi, \mathbf{P})$$

$$= \min_{\mathbf{P} \in \mathbb{P}} \max_{\pi \in \Pi_S} R(\pi, \mathbf{P})$$

$$= \max_{\pi \in \Pi_S} \min_{\mathbf{P} \in \mathbb{P}} R(\pi, \mathbf{P})$$

Lemma's 4.1 and 4.3

Lemma 4.1 *Let \mathbb{P} be an r -rectangular uncertainty set. Under Assumption 2.4, there exists a stationary and deterministic policy solution of the policy improvement problem.*

Lemma 4.3 *Let $\pi \in \Pi$ and $\mathbf{W} \in \mathcal{W}$. Let \mathbf{v} be the value function of the decision maker and β be the value function of the adversary. Then $\mathbf{W}^\top \mathbf{v} = \beta$.*

- This is the (unique) solution to the Bellman equation of the adversary

$$\beta_i = \mathbf{w}_i^\top (\mathbf{r}_\pi + \lambda \cdot \mathbf{T}_\pi \beta), \forall i \in [r].$$

- So it suffices to show that

$$(\mathbf{W}^\top \mathbf{v})_i = \mathbf{w}_i^\top (\mathbf{r}_\pi + \lambda \mathbf{T}_\pi \mathbf{W}^\top \mathbf{v}), \forall i \in [r].$$

Proof of 4.3

- To proof: $(\mathbf{W}^\top \mathbf{v})_i = \mathbf{w}_i^\top (\mathbf{r}_\pi + \lambda \mathbf{T}_\pi \mathbf{W}^\top \mathbf{v}), \forall i \in [r]$.
- Bellman equation for the decision maker:

$$v_s = \sum_{a \in \mathbb{A}} \pi_{sa} (r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \mathbf{w}_i^\top \mathbf{v}), \forall s \in \mathbb{S}.$$

$$\mathbf{T}_\pi = (\sum_{a \in \mathbb{A}} \pi_{sa} u_{sa}^i)_{(s,i) \in \mathbb{S} \times [r]} \in \mathbb{R}_+^{S \times r}.$$

$$\mathbf{v} = \mathbf{r}_\pi + \lambda \mathbf{T}_\pi \mathbf{W}^\top \mathbf{v}.$$

$$\mathbf{w}_i^\top \mathbf{v} = \mathbf{w}_i^\top (\mathbf{r}_\pi + \lambda \cdot \mathbf{T}_\pi \mathbf{W}^\top \mathbf{v}), \forall i \in [r].$$

$$\mathbf{W}^\top \mathbf{v} = (\mathbf{w}_i^\top \mathbf{v})_{i \in [r]}.$$

Paper 1: Min-Max Duality
Theorem 4.2: Proof

$$\max_{\pi \in \Pi} \min_{\mathbf{W} \in \mathcal{W}} R(\pi, \mathbf{W}) \leq \min_{\mathbf{W} \in \mathcal{W}} \max_{\pi \in \Pi} R(\pi, \mathbf{W}).$$

Our goal is to show $\pi^* \in \arg \max_{\pi \in \Pi} R(\pi, \mathbf{W}^*)$.

- Equivalent to:

$$v_s^* = \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \mathbf{w}_i^{*\top} \mathbf{v}^* \right\}, \forall s \in \mathbb{S}.$$

- Proof:

$$\begin{aligned} v_s^* &= r_{sa^*(s)} + \lambda \cdot \sum_{i=1}^r u_{sa^*(s)}^i \mathbf{w}_i^{*\top} \mathbf{v}^* \\ &= r_{sa^*(s)} + \lambda \cdot \sum_{i=1}^r u_{sa^*(s)}^i \beta_i^* \\ &= \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \beta_i^* \right\} \\ &= \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \mathbf{w}_i^{*\top} \mathbf{v}^* \right\}, \end{aligned}$$

Theorem 5.1

$$v_s^* = \max_{\pi_s \in \Delta} \{r_{\pi_s, s} + \lambda \cdot (\mathbf{T}_\pi \mathbf{W}^{* \top} \mathbf{v}^*)_s\}, \forall s \in \mathbb{S},$$

$$\mathbf{v} = \mathbf{r}_\pi + \lambda \mathbf{T}_\pi \mathbf{W}^\top \mathbf{v}.$$

$$\beta_i^* = \min_{\mathbf{w}_i \in \mathcal{W}^i} \{\mathbf{w}_i^\top (\mathbf{r}_{\pi^*} + \lambda \cdot \mathbf{T}_{\pi^*} \beta^*)\}, \forall i \in [r].$$

- By substitution of: $\mathbf{W}^{* \top} \mathbf{v}^* = \beta^*$,

$$v_s^* = \max_{\pi_s \in \Delta} \{r_{\pi_s, s} + \lambda \cdot (\mathbf{T}_\pi (\min_{\mathbf{w}_i \in \mathcal{W}^i} \{\mathbf{w}_i^\top \mathbf{v}^*\})_{i \in [r]})_s\}, \forall s \in \mathbb{S}.$$

$$v_s^* = \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \min_{\mathbf{w}_i \in \mathcal{W}^i} \mathbf{w}_i^\top \mathbf{v}^* \right\}, \forall s \in \mathbb{S}.$$

- So, these two are equivalent:

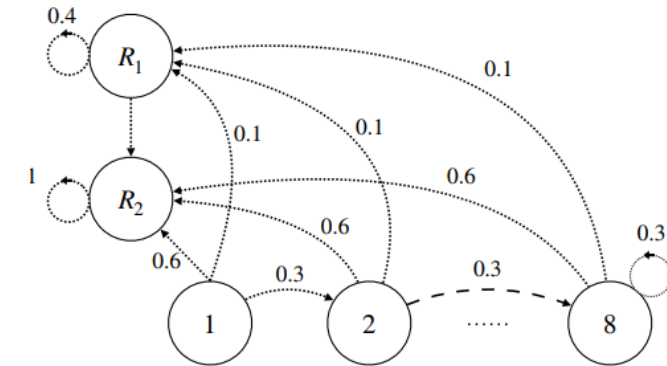
$$F(\mathbf{v})_s = \max_{\pi_s \in \Delta} \{r_{\pi_s, s} + \lambda \cdot (\mathbf{T}_\pi (\min_{\mathbf{w}_i \in \mathcal{W}^i} \{\mathbf{w}_i^\top \mathbf{v}\})_{i \in [r]})_s\}, \forall s \in \mathbb{S}, \forall \mathbf{v} \in \mathbb{R}_+^{\mathbb{S}}.$$

$$\mathbf{v}^0 = \mathbf{0}, v_s^{k+1} = \max_{a \in \mathbb{A}} \left\{ r_{sa} + \lambda \cdot \sum_{i=1}^r u_{sa}^i \min_{\mathbf{w}_i \in \mathcal{W}^i} \mathbf{w}_i^\top \mathbf{v}^k \right\}, \forall s \in \mathbb{S}, \forall k \geq 0.$$

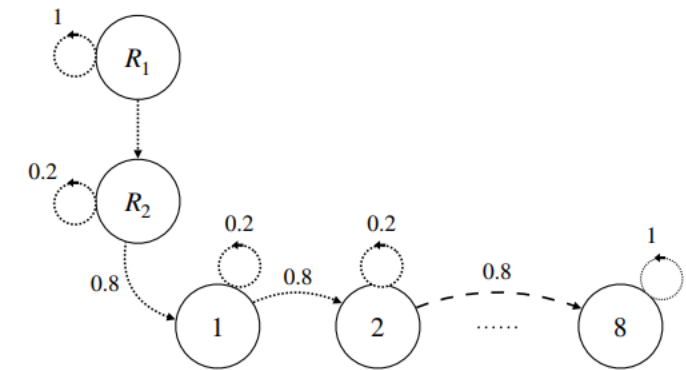
Paper 1: Numerical Experiments
Testing

- Test 1: Machine replacement problem
- Conclusion: R-rect performs better than s-rect.

Budget of deviation τ	0.05	0.07	0.09
Worst-case of π^{nom} for $\mathbb{P}^{(r)}$	94.40	92.21	90.04
Worst-case of π^{nom} for $\mathbb{P}^{(s)}$	91.74	88.56	85.46
Budget of deviation τ	0.05	0.07	0.09
Nominal reward of $\pi^{\text{rob},r}$	100.00	100.00	100.00
Worst-case of $\pi^{\text{rob},r}$ for $\mathbb{P}^{(r)}$	94.40	92.21	90.04
Nominal reward of $\pi^{\text{rob},s}$	99.28	98.53	97.81
Worst-case of $\pi^{\text{rob},s}$ for $\mathbb{P}^{(s)}$	91.90	89.09	86.62



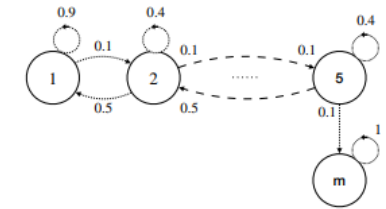
(a) Transition probabilities for action *repair*.



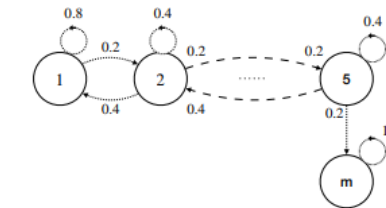
(b) Transition probabilities for action *wait*.

Paper 1: Numerical Experiments
Testing

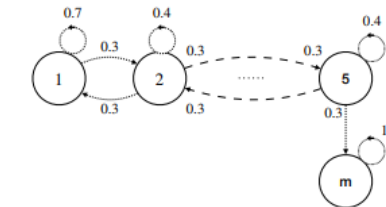
- Test 2: Inspired by healthcare
- Conclusion: Worst-case performance is actually worse, but average isn't!



(a) Transition probabilities for action *high dosage*.



(b) Transition probabilities for action *medium dosage*.



(c) Transition probabilities for action *low dosage*.

Budget of deviation τ	$\tau = 0.05$	$\tau = 0.07$	$\tau = 0.09$
Nominal reward of π^{nom}	100.00	100.00	100.00
Worst-case of π^{nom} for $\mathbb{P}(r)$	50.26	41.74	35.63
Worst-case of π^{nom} for $\mathbb{P}(s)$	45.75	37.37	31.51
Nominal reward of $\pi^{\text{rob},r}$	100.00	92.92	92.92
Worst-case of $\pi^{\text{rob},r}$ for $\mathbb{P}(r)$	50.26	42.29	36.56
Nominal reward of $\pi^{\text{rob},s}$	91.48	91.35	89.56
Worst-case of $\pi^{\text{rob},s}$ for $\mathbb{P}(s)$	52.09	44.39	38.69