# What is optimality is time-average MDPs?

#### Nicolas Gast, Bruno Gaujal, Kimang Khun

Inria

Workshop Mama 2023, June 19, 2023

#### Motivation: Whittle index

"A problem is indexable if and only if the optimal policy  $\pi^*(\lambda)$  is non-decreasing".

# Motivation: Whittle index

"A problem is indexable if and only if the optimal policy  $\pi^*(\lambda)$  is non-decreasing".

Ambiguous definition:

- What is "optimal"?
- If the optimal: is it unique?

# Markov deicison processes (MDPs)

Markov decision processes = Markov chains + actions and rewards.



# Markov deicison processes (MDPs)

Markov decision processes = Markov chains + actions and rewards.



- Introduced in the 50s (Bellman)
- Very popular today because of reinforcement learning

Find a policy  $\pi: \mathcal{S} \to \mathcal{A}$  to maximize some *optimality criterion*.

#### What are optimality criteria?

- Finite horizon:  $\max \mathbb{E}[R_0 + R_1 + \cdots + R_T]$ .
- **2** Discounted:  $\max \mathbb{E} \left[ R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \right]$  for some  $\gamma < 1$ .
- Time-average (a.k.a. gain): max  $\lim_{T\to\infty} \frac{1}{T}\mathbb{E}[R_0 + R_1 + \cdots + R_T].$

#### What are optimality criteria?

- Finite horizon:  $\max \mathbb{E} [R_0 + R_1 + \cdots + R_T]$ .
- **2** Discounted:  $\max \mathbb{E} \left[ R_0 + \gamma R_1 + \gamma^2 R_2 + \dots \right]$  for some  $\gamma < 1$ .
- Time-average (a.k.a. gain):  $\max \lim_{T \to \infty} \frac{1}{T} \mathbb{E} [R_0 + R_1 + \cdots + R_T].$

What is the "good" notion of optimality for time-average MDPs?

- Time-average makes sense for queueing applications.
- Algorithms are designed for finite/discounted.
- Strong connections between the three notions

### Outline



2 Discounted MDP and *n*-sensitive optimality

3 Bellman optimality (a.k.a. canonical optimality)

#### 4 Conclusion

#### Definition of a MDP

At time *t*, you

- Observe  $S_t \in S$  and take an action  $A_t \in A$  (we assume S, A finite).
- Receive  $R_t = r(S_t, A_t)$ .
- $S_{t+1}$  jumps according to  $P(\cdot|S_t, A_t)$ .
- A (deterministic) policy is a function  $\pi: S \to A$ .

# Notion of gain-optimality

The gain of a policy is:

$$g^{\pi}(s) = \lim_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=0}^{T-1} r(S_t, A_t) \mid S_0 = s\right].$$

A gain-optimal policy<sup>1</sup>  $\pi^*$  maximizes the gain, *i.e.*, for all *s*:

 $g^{\pi^*}(s) = rg\max_{\pi} g^{\pi}(s).$ 

<sup>&</sup>lt;sup>1</sup>Theorem (e.g., Puterman 2005). When S and A are finite, gain and optimal policies are well defined.

# Is gain-optimal a good notion of optimality? Example of deterministic MDPs



# Is gain-optimal a good notion of optimality? Example of deterministic MDPs





# Is gain-optimal a good notion of optimality? Example of deterministic MDPs



Here: all policies are gain-optimal but some look better than others.

- What is the right notion of optimality?
- What do algorithms compute?

# Outline

#### 1 MDP and gain-optimality

#### 2 Discounted MDP and *n*-sensitive optimality

#### 3 Bellman optimality (a.k.a. canonical optimality)

#### 4 Conclusion

# Discounted optimality

For discounted problems, the quantity of importance is the value function:

$$V^{\pi}_{\gamma}(s) = \mathbb{E}\left[R_0 + \gamma R_1 + \gamma^2 R_2 + \cdots \mid S_0 = s
ight],$$

# Discounted optimality

For discounted problems, the quantity of importance is the value function:

$$V^{\pi}_{\gamma}(s) = \mathbb{E}\left[R_0 + \gamma R_1 + \gamma^2 R_2 + \cdots \mid S_0 = s
ight],$$

The optimal policy satisfies Bellman equation:

Bellman equation is very powerful:

- Unique solution for any  $\gamma < 1$ .
- Many algorithms to solve it (VI, PI, Q-learning).
- Can be used to define advantage  $Q^*(s, a) V^*(s)$  (e.g., actor-critic).

Discounted-optimality when  $\gamma \rightarrow 1$ 



Policy $\pi$	discounted gain $V^\pi_\gamma(A)$		
Stay left	$1 + \gamma + \gamma^2 + \dots$	$=rac{1}{1-\gamma}$	
Go right:	$2 + \gamma + \gamma^2 + \gamma^3 + \dots$	$=rac{1}{1-\gamma}+1.$	Best policy?
Alternate:	$2+2\gamma^2+2\gamma^4\dots$	$=rac{1}{1-\gamma}+rac{1}{1+\gamma}.$	

(Puterman 2015): For any policy  $\pi$ , there exists a bias vector  $h^{\pi}(s)$  s.t.:



(Puterman 2015): For any policy  $\pi$ , there exists a bias vector  $h^{\pi}(s)$  s.t.:

$$V^{\pi}_{\gamma}(s) = rac{1}{1-\gamma} \underbrace{g^{\pi}(s)}_{ ext{gain}} + \underbrace{h^{\pi}(s)}_{ ext{bias}} + o(1-\gamma).$$

A policy is:

- Gain-optimal if it maximizes  $g^{\pi}(s)$  (for all s).
- Bias-optimal if it maximizes  $g^{\pi}(s)$  (for all s) and then  $h^{\pi}(s)$ .

(Puterman 2015): For any policy  $\pi$ , there exists a bias vector  $h^{\pi}(s)$  s.t.:

$$V^{\pi}_{\gamma}(s) = rac{1}{1-\gamma} \underbrace{g^{\pi}(s)}_{ ext{gain}} + \underbrace{h^{\pi}(s)}_{ ext{bias}} + o(1-\gamma).$$

A policy is:

- Gain-optimal if it maximizes  $g^{\pi}(s)$  (for all s).
- Bias-optimal if it maximizes  $g^{\pi}(s)$  (for all s) and then  $h^{\pi}(s)$ .

#### What do algorithms compute? Gain-optimal policies? Bias-optimal policies?

(Puterman 2015): For any policy  $\pi$ , there exists a bias vector  $h^{\pi}(s)$  s.t.:

$$V_{\gamma}^{\pi}(s) = \frac{1}{1-\gamma} \underbrace{g^{\pi}(s)}_{\text{gain}} + \underbrace{h^{\pi}_{0}(s)}_{0\text{-bias}} + (1-\gamma) \underbrace{h^{\pi}_{1}(s)}_{1\text{-bias}} + (1-\gamma)^{2} \underbrace{h^{\pi}_{2}(s)}_{2\text{-bias}} + \dots$$

A policy is:

- Gain-optimal if it maximizes  $g^{\pi}(s)$  (for all s).
- Bias-optimal if it maximizes  $g^{\pi}(s)$  (for all s) and then  $h^{\pi}(s)$ .
- (Blackwell-optimal if it maximizes  $g^{\pi}(s)$  then  $h_0^{\pi}(s)$  then  $h_1^{\pi}(s),...)$

What do algorithms compute? Gain-optimal policies? Bias-optimal policies?

# Outline

1 MDP and gain-optimality

2 Discounted MDP and *n*-sensitive optimality

3 Bellman optimality (a.k.a. canonical optimality)

#### 4 Conclusion

# The (modified) Bellman equation

Let  $g^*$  be the optimal gain. There exists a bias vector h such that:

$$g^{*}(s) = \max_{a} g^{*}(s')p(s'|s,a)$$
  
$$h(s) + g^{*}(s) = \max_{a} r(s,a) + \sum_{s'} h(s')p(s'|s,a).$$
 (1)

Note: the solution is not unique (not just up to a constant).

# The (modified) Bellman equation

Let  $g^*$  be the optimal gain. There exists a bias vector h such that:

$$g^{*}(s) = \max_{a} g^{*}(s')p(s'|s,a)$$
  
$$h(s) + g^{*}(s) = \max_{a} r(s,a) + \sum_{s'} h(s')p(s'|s,a).$$
 (1)

Note: the solution is not unique (not just up to a constant).

We call a best response to (1) a Bellman-optimal policy.

What is a Bellman optimal policy?

Bellman-optimal policies are not gain/bias optimal policies

 $\mathsf{gain-optimal} \supsetneq \mathsf{Bellman-optimal} \supsetneq \mathsf{bias-optimal}$ 



Black is Bellman-opt All are Gain-optimal Bellman-optimal policies are not gain/bias optimal policies

 $\mathsf{gain}\mathsf{-}\mathsf{optimal} \supsetneq \mathsf{Bellman}\mathsf{-}\mathsf{optimal} \supsetneq \mathsf{bias}\mathsf{-}\mathsf{optimal}$ 



Black is Bellman-opt All are Gain-optimal



All policies are Bellman-optimal There exists a unique bias-optimal policy Bellman-optimal policies are not gain/bias optimal policies

 $\mathsf{gain-optimal} \supsetneq \mathsf{Bellman-optimal} \supsetneq \mathsf{bias-optimal}$ 



All are Gain-optimal

All policies are Bellman-optimal There exists a unique bias-optimal policy Bellman optimal policy is a natural notion

#### Theorem

The set of optimal policies that can be output of policy iteration of value iteration is the set of Bellman-optimal policy.

Bellman optimal policy is a natural notion

#### Theorem

The set of optimal policies that can be output of policy iteration of value iteration is the set of Bellman-optimal policy.

A policy  $\pi$  is *canonical optimal* if there exists a final reward F such that  $\pi$  is optimal for for all finite horizon with final reward F, *i.e.* it maximizes

 $\mathbb{E}\left[R_0+R_1+\cdots+R_T+F(S_T)\right].$ 

#### Theorem (Yushkevich 1974)

• A policy is Bellman-optimal if and only if it is canonical optimal.

#### Does it matter: Definition of computation of Whittle index

Consider a two-action MDP, with a penalty:

•  $P(\cdot|s_n, a_n)$  and  $r(s_n, a_n) - \lambda a_n$ .

# Does it matter: Definition of computation of Whittle index

Consider a two-action MDP, with a penalty:

• 
$$P(\cdot|s_n, a_n)$$
 and  $r(s_n, a_n) - \lambda a_n$ .

#### Classical definition of index

The Whittle index of s is a penalty  $\lambda_s$  such that that the optimal policy chooses  $\pi(s) = 1$  when  $\lambda < \lambda_s$  and  $\pi(s) = 0$  when  $\lambda > \lambda_s$ .

# Does it matter: Definition of computation of Whittle index

Consider a two-action MDP, with a penalty:

• 
$$P(\cdot|s_n, a_n)$$
 and  $r(s_n, a_n) - \lambda a_n$ .

#### Non-ambiguous definition of index

The Whittle index of s is the unique penalty  $\lambda_s$  such that that any (Bellman-)optimal policy chooses  $\pi(s) = 1$  when  $\lambda < \lambda_s$  and  $\pi(s) = 0$  when  $\lambda > \lambda_s$ .

### How to compute Whittle indices?

A Bellman-optimal policies satisfies Bellman equations:

$$g_{\lambda}^{*}(s) + h_{\lambda}^{*}(s) = \max_{a} r(s, a) + a\lambda + \sum_{s'} P(s'|s, a)h_{\lambda}^{*}(s')$$

We define the active advantage  $b_{\lambda}(s) := q_{\lambda}(s, 1) - q_{\lambda}(s, 0)$ .

# How to compute Whittle indices?

A Bellman-optimal policies satisfies Bellman equations:

$$g_{\lambda}^{*}(s) + h_{\lambda}^{*}(s) = \max_{a} r(s, a) + a\lambda + \sum_{s'} P(s'|s, a)h_{\lambda}^{*}(s')$$

We define the active advantage  $b_{\lambda}(s) := q_{\lambda}(s, 1) - q_{\lambda}(s, 0)$ .



Theorem (G,Gaujal,Khun, 22)

An arm is indexable if and only if for all s:  $b_{\lambda}(s, 1) = 0$  has a unique solution.

#### We can use to build an subcubic algorithm



Three ingredients:

**1** For MDP, the advantage function is piecewise linear:

$$b_{\lambda}^{\pi} = (A^{\pi})^{-1}(r + \lambda \pi).$$

#### We can use to build an subcubic algorithm



Three ingredients:

I For MDP, the advantage function is piecewise linear:

$$b_{\lambda}^{\pi} = (A^{\pi})^{-1}(r + \lambda \pi).$$

Sherman-Morisson formula: Let A be an invertible matrix, u and v vectors 1D such that  $1 + v^T A^{-1} u \neq 0$ . Then:

$$\left(A + uv^{T}\right)^{-1} = A^{-1} - \frac{A^{-1}uv^{T}A^{-1}}{1 + v^{T}A^{-1}u}$$

#### We can use to build an subcubic algorithm



Three ingredients:

• For MDP, the advantage function is piecewise linear:

$$b_{\lambda}^{\pi} = (A^{\pi})^{-1}(r + \lambda \pi).$$

Sherman-Morisson formula: Let A be an invertible matrix, u and v vectors 1D such that  $1 + v^T A^{-1} u \neq 0$ . Then:

$$\left(A + uv^{T}\right)^{-1} = A^{-1} - \frac{A^{-1}uv^{T}A^{-1}}{1 + v^{T}A^{-1}u}.$$

We can reorder operations to use Strassen's like operations.

# We obtain a theoretical complexity of $O(S^{2.53})$ and an efficient implemenation

https://pypi.org/project/markovianbandit-pkg/



# Outline

#### 1 MDP and gain-optimality

2 Discounted MDP and *n*-sensitive optimality

#### 3 Bellman optimality (a.k.a. canonical optimality)

#### 4 Conclusion

# Conclusion

Time-average MDPs are complicated:

- Gain-optimality is defined but stronger notions of optimality are used.
- Be careful about structure.

• Bellman-optimality allows to define the advantage (as for discounted problems).

```
http://polaris.imag.fr/nicolas.gast/
```

Computing Whittle (and Gittins) Index in Subcubic Time, G. Gaujal, Khun https://arxiv.org/abs/2203.05207