

# TRADING-OFF REWARD MAXIMIZATION AND STABILITY IN SEQUENTIAL DECISION MAKING





FEDERICO CORSO, MARCO MUSSI, AND ALBERTO MARIA METELLI

{federico.corso, marco.mussi, albertomaria.metelli}@polimi.it

#### MOTIVATION

RL objective: maximize expected reward.

Limitation: Reward maximization alone is often *insufficient* in real-world, safety-critical domains. Properties beyond pure *performance* are difficult to encode in the reward.

**Our idea:** *Enforce stability within RL*, inspired by control theory, where stability is fundamental for systems' robustness, safety, and reliability.

## OBJECTIVE

Find a control policy that *trades* off between maximizing reward and ensuring a stable behavior.

#### SETTING: UNDISCOUNTED MDPs

We consider infinite-horizon undiscounted MDPs:

$$\mathcal{M}\coloneqq \langle \mathcal{S}, \mathcal{A}, \mathbf{P}, \mathbf{r}, oldsymbol{\eta}_0 
angle$$

#### **Assumptions:**

- *Finite* state and action spaces.
- *Ergodicity:* for any  $\pi \in \Pi$ , the induced chain admits a unique stationary distribution  $\eta^{\pi}$ :

$$\lim_{t o\infty}oldsymbol{\eta}_t^{oldsymbol{\pi}}=oldsymbol{\eta}^{oldsymbol{\pi}}$$

#### A MEASURE OF STABILITY

$$\mathcal{M} \, + \, oldsymbol{\pi} \ \Rightarrow \ \mathcal{M}^{oldsymbol{\pi}} \coloneqq \langle \mathcal{S}, \mathbf{P}^{oldsymbol{\pi}}, oldsymbol{\eta}_0 
angle$$

A policy  $\pi$  is *more stable* if the Markov Chain (MC)  $\mathcal{M}^{\pi}$  induced converges faster to its stationary distribution  $\eta^{\pi}$ .

The asymptotic convergence rate is governed by the Second Largest Eigenvalue Modulus (SLEM):

$$\lim_{t \to \infty} \left( \max_{s \in \mathcal{S}} \left\| \mathbf{p}_t^{\boldsymbol{\pi}}(\cdot|s) - \boldsymbol{\eta}^{\boldsymbol{\pi}} \right\|_{\text{TV}} \right)^{1/t} = \mu(\mathbf{P}^{\boldsymbol{\pi}})$$

with  $\mu(\mathbf{P}^{\boldsymbol{\pi}}) \coloneqq \max_{i \geq 2} |\lambda_i|$ .

# SINGLE OBJECTIVE: reward AND stability

$$\max_{\pi \in \Pi} \ (\boldsymbol{\eta}^{\boldsymbol{\pi}})^{\top} \mathbf{r} \ - \ \mu(\mathbf{P}^{\boldsymbol{\pi}})$$

## AN EXPLICIT FORM OF THE OBJECTIVE: BILINEARITY ISSUES

maximize 
$$\boldsymbol{\eta}^{\top}(\boldsymbol{\Pi}\mathbf{r}) - \left\| \mathbf{D}_{\boldsymbol{\eta}}^{1/2}(\boldsymbol{\Pi}\mathbf{P}) \mathbf{D}_{\boldsymbol{\eta}}^{-1/2} - \sqrt{\boldsymbol{\eta}}\sqrt{\boldsymbol{\eta}}^{\top} \right\|_{2}$$
subject to  $\boldsymbol{\eta} = (\boldsymbol{\Pi}\mathbf{P})^{\top}\boldsymbol{\eta} \to \text{Stationary constraint}$ 
Bilinear in  $(\boldsymbol{\eta}, \boldsymbol{\Pi})$ 

The problem is **not jointly-convex** in  $(\eta, \Pi)$  and the stationary constraints implies fixing  $\eta$  once  $\Pi$  is fixed

## A SURROGATE OBJECTIVE IN $\mathcal{S} imes \mathcal{A}$ : Bypassing Bilinearity

### (1) Introduce a New MC over $S \times A$

$$\mathcal{T}\coloneqq \langle \mathcal{S} imes \mathcal{A},\, \mathbf{T},\, \mathbf{x}_0
angle$$

where:

- $T = P\Pi$  (state-action transition matrix)
- $\mathbf{x} = \boldsymbol{\eta}^{\top} \boldsymbol{\Pi}$  (state–action *stationary* distribution)

## Eigenvalue relationship between $P\Pi$ and $\Pi P$

$$\Lambda(\mathbf{P\Pi}) = \Lambda(\mathbf{\Pi}\mathbf{P}) \cup \{0\}^{|\mathcal{S}||\mathcal{A}|-|\mathcal{S}|}$$

(2) Rewrite the Objective in  $S \times A$ 

$$\frac{\mu(\Pi \mathbf{P})}{\|\mathbf{\Pi} \mathbf{r}\|} - \left\| \mathbf{D}_{\boldsymbol{\eta}}^{1/2}(\Pi \mathbf{P}) \mathbf{D}_{\boldsymbol{\eta}}^{-1/2} - \sqrt{\boldsymbol{\eta}} \sqrt{\boldsymbol{\eta}}^{\top} \right\|_{2}$$

$$\| \mathbf{x}^{\top} \mathbf{r}\| - \left\| \mathbf{D}_{\mathbf{x}}^{1/2} \mathbf{P} \Pi \mathbf{D}_{\mathbf{x}}^{-1/2} - \sqrt{\mathbf{x}} \sqrt{\mathbf{x}}^{\top} \right\|_{2}$$

$$\| \mathbf{\sigma}_{2}(\mathbf{P} \mathbf{\Pi}) \|_{2}$$

Lifting removes the  $oldsymbol{\eta}^{ op} oldsymbol{\Pi}$  bilinearity in the reward term

## (3) Bilinearity Resolution

- $\bullet \;$  Introduce the auxiliary variable  $\mathbf{X} = \mathbf{D_x} \mathbf{T}$
- Rewrite the problem in (x, X)

maximize 
$$\mathbf{X}, \mathbf{X}$$
  $\mathbf{X}, \mathbf{X}$   $\mathbf{X}, \mathbf{X}$  subject to  $\mathbf{X}, \mathbf{X}, \mathbf{X}$   $\mathbf{X}, \mathbf{X}$   $\mathbf{X}$   $\mathbf{X}, \mathbf{X}$   $\mathbf{X}$   $\mathbf{X}$ 

#### (4) The Final Problem

$$\begin{aligned} & \underset{\mathbf{X} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}{\operatorname{maximize}} & (\mathbf{X}\mathbf{1})^{\top}\mathbf{r} - \left\| \mathbf{D}_{\mathbf{x}^{\star}}^{-1/2} \mathbf{X} \mathbf{D}_{\mathbf{x}^{\star}}^{-1/2} - \sqrt{\mathbf{x}^{\star}} \sqrt{\mathbf{x}^{\star}}^{\top} \right\|_{2} \\ & \text{subject to} & \mathbf{1}^{\top} \mathbf{X} \mathbf{1} = 1, \\ & \| \mathbf{X} \mathbf{1} - \mathbf{x}^{\star} \|_{2}^{2} \leq \delta^{2} \rightarrow \textit{Relaxed} \; \text{stationary constraints} \end{aligned}$$

It is **convex** in **X** 

x\*: solution of the standard undiscounted MDP

#### RECOVERING THE POLICY

#### Given:

- $X^{\dagger}$ : Solution to the final problem.
- $\mathbf{x}^{\dagger} = \mathbf{X}^{\dagger}\mathbf{1}$ : Stationary distribution associated to the MC with transition  $\mathbf{T}^{\dagger} = \mathbf{D}_{\mathbf{x}^{\dagger}}^{-1}\mathbf{X}^{\dagger}$ .

We can compute the optimal policy:

$$\pi^{\dagger}(a \mid s) \coloneqq \frac{x^{\dagger}(s, a)}{\sum_{a \in \mathcal{A}} x^{\dagger}(s, a)}$$

Let  $\pi^{\dagger} \in \Pi$  be a solution, then it holds:

$$(\mathbf{x}^{\star} - \mathbf{x}^{\dagger})^{\top} \mathbf{r} \leq \delta R_{\max} \sqrt{|\mathcal{S}||\mathcal{A}|},$$

anc

$$\sigma_2(\mathbf{D}_{\mathbf{x}^{\star}}^{1/2}\mathbf{P}\mathbf{\Pi}^{\dagger}\mathbf{D}_{\mathbf{x}^{\star}}^{-1/2}) \leq \sigma_2(\mathbf{D}_{\mathbf{x}^{\star}}^{1/2}\mathbf{P}\mathbf{\Pi}^{\star}\mathbf{D}_{\mathbf{x}^{\star}}^{-1/2}).$$

#### MAIN LIMITATION

Improvement on the upper-bound of the SLEM is **not guaranteed**:

$$\sigma_2(\mathbf{D}_{\mathbf{x}^{\dagger}}^{1/2}\mathbf{T}^{\dagger}\mathbf{D}_{\mathbf{x}^{\dagger}}^{-1/2}) \leq \sigma_2(\mathbf{D}_{\mathbf{x}^{\star}}^{1/2}\mathbf{T}^{\star}\mathbf{D}_{\mathbf{x}^{\star}}^{-1/2})$$

#### FUTURE WORKS

- Guarantee monotonic improvement of the SLEM upper-bound
- Experimental evaluation
- Consider unknown model setting
- Extend to continuous state and action spaces

#### REFERENCES

Stephen P. Boyd, Persi Diaconis, and Lin Xiao. Fastest mixing markov chain on a graph. *SIAM Review*, 2004.

Ilse C. F. Ipsen and Teresa M. Selee. Ergodicity coefficients defined by vector norms. *SIAM Journal on Matrix Analysis and Applications*, 2011.

Mirco Mutti and Marcello Restelli. An intrinsically-motivated approach for learning highly exploring and fast mixing policies. In *AAAI Conference on Artificial Intelligence*, 2020.

Jean Tarbouriech and Alessandro Lazaric. Active exploration in markov decision processes. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019