



Closed-loop control of complex nia systems using deep Reinforcement Learning **An off-policy approach**





1) LISN, CNRS, Université Paris-Saclay, Orsay (FR) 2) TAU, Inria, Université Paris-Saclay, CNRS, LISN, Orsay, (FR) 3) CNRS, Université de Poitiers, ISAE-ENSMA

Mail: michele-alessandro.bucci@inria.fr

"Flowcon", ANR-17-ASTR-0022

T. Guégan³, M.A. Bucci², O. Semeraro¹, L. Cordier³ & L. Mathelin¹







2 / 23





















2 / 23

Laminar to turbulent transition









- Laminar to turbulent transition
- Separation zone









- Laminar to turbulent transition
- Separation zone
- Buffet (λ)







- Laminar to turbulent transition
- Separation zone
- Buffet (λ)
- Secondary flows

2 / 23

- Laminar to turbulent transition
- Separation zone
- Buffet (λ)

- Secondary flows
- Fluid structure interaction

Interact with the system **locally** to:

Euromech Colloquium 614

2 / 23

- Laminar to turbulent transition
- Separation zone
- Buffet (λ)
- Secondary flows
- Fluid structure interaction

$\mathcal{J} = \min\{\text{drag}\}, \max\{\text{lift}\}, \min\{\text{noise}\}...$

- $\mathcal{J} = \min_{u} \int_{0}^{\infty} (x^{T}Qx + u^{T}Ru)dt \qquad x \in \mathbb{R}^{n} \text{ sensors signal (state)}$
 - s.t. $\dot{x} = Ax + Bu$ $u \in \mathbb{R}^m$ actuators signal (action)

Augmented Lagrangian

- $\mathscr{J} = \min_{u} \int_{0}^{\infty} (x^{T}Qx + u^{T}Ru)dt$ $x \in \mathbb{R}^{n}$ sensors signal (state)
 - s.t. $\dot{x} = Ax + Bu$ $u \in \mathbb{R}^m$ actuators signal (action)

Augmented Lagrangian

- $\mathscr{J} = \min_{u} \int_{0}^{\infty} (x^{T}Qx + u^{T}Ru)dt$ $x \in \mathbb{R}^{n}$ sensors signal (state)
 - $u \in \mathbb{R}^m$ actuators signal (action) s.t. $\dot{x} = Ax + Bu$

$$+ u^{T} R u) dt - \lambda (\dot{x} - Ax - Bu)$$

oction constraints

$$\lambda \in \mathbb{R}^n$$
 co-state

Augmented Lagrangian

$$\frac{\partial \mathcal{J}}{\partial \lambda} = 0 \Rightarrow \dot{x} = \left(\frac{\partial H}{\partial \lambda}\right)^T = Ax + Bu \quad x(0) = x_0$$
$$\frac{\partial \mathcal{J}}{\partial x} = 0 \Rightarrow -\dot{\lambda} = \left(\frac{\partial H}{\partial x}\right)^T = Qx + A^T\lambda \quad \lambda(T) = 0$$
$$\frac{\partial \mathcal{J}}{\partial u} = 0 \Rightarrow 0 = \left(\frac{\partial H}{\partial u}\right)^T = Ru + \lambda^T B = 0$$

Euromech Colloquium 614

3 / 23

- $\mathscr{J} = \min_{u} \int_{0}^{\infty} (x^{T}Qx + u^{T}Ru)dt$ $x \in \mathbb{R}^{n}$ sensors signal (state)
 - s.t. $\dot{x} = Ax + Bu$ $u \in \mathbb{R}^m$ actuators signal (action)

$$\frac{+ u^{T} R u}{dt} - \lambda (\dot{x} - A x - B u) \qquad \lambda \in \mathbb{R}^{n} \text{ co-state}$$

notion constraints

$$\frac{\partial \mathcal{J}}{\partial \lambda} = 0 \Rightarrow \dot{x} = \left(\frac{\partial H}{\partial \lambda}\right)^T = Ax + Bu \quad x(0) = x_0$$
$$\frac{\partial \mathcal{J}}{\partial x} = 0 \Rightarrow -\dot{\lambda} = \left(\frac{\partial H}{\partial x}\right)^T = Qx + A^T\lambda \quad \lambda(T) = 0$$
$$\frac{\partial \mathcal{J}}{\partial u} = 0 \Rightarrow 0 = \left(\frac{\partial H}{\partial x}\right)^T = Ru + \lambda^T B = 0$$

Euromech Colloquium 614

Linear methods

$\frac{\partial \mathcal{J}}{\partial \lambda} = 0 \Rightarrow .$
$\frac{\partial \mathcal{J}}{\partial x} = 0 \Rightarrow -\frac{\partial \mathcal{J}}{\partial x}$
$\frac{\partial \mathcal{J}}{\partial u} = 0 \Rightarrow 0$

Iterative methods: Direct-Adjoint

Linear iteration, Adjoint of direct adjoint (ADA), extension to nonlinear cases, Model Predictive Control (MPC)...

Direct methods: Control Algebraic **Riccati Equation**

Kim, J. and Bewley T.R. (2007). Annual Review of Fluid Mech. Luchini, P. and Bottaro, A. (2014). Annual Review of Fluid Mech.

$$\dot{k} = \left(\frac{\partial H}{\partial \lambda}\right)^T = Ax + Bu \quad x(0) = x_0$$
$$-\dot{\lambda} = \left(\frac{\partial H}{\partial x}\right)^T = Qx + A^T \lambda \quad \lambda(T) = 0$$
$$0 = \left(\frac{\partial H}{\partial x}\right)^T = Ru + \lambda^T B = 0$$

Ínia Optimal control: the *non*linear case

$$\mathcal{J}(x(t), t) = \max_{u} \int_{t}^{\infty} (x_{t}) dt$$

- The integrand of the cost function is the **reward** r(x, u)
- Non linear dynamical system

$x^TQx + u^TRu)dt$ st $\dot{x} = Ax + Bu$

the objective is to find the **policy** $u = \pi(x)$ that maximise the **value function**

Ínría Optimal control: the *non*linear case

$$\mathcal{J}(x(t), t) = \max_{u} \int_{t}^{\infty} r(t) dt$$

- The integrand of the cost function is the **reward** r(x, u)
- Non linear dynamical system

the objective is to find the **policy** $u = \pi(x)$ that maximise the **value function**

Ínría Hamilton-Jacobi-Bellman (HJB) equation

$$\mathcal{J}(x(t), t) = \max_{u} \int_{t}^{t+\Delta t} r(x, u) dt + \mathcal{J}(x(t+\Delta t), t+\Delta t)$$

maximised expanding in series and taking $\Delta t \rightarrow 0$

future value function

The principle of optimality requires that the **future value function** has to be

$$\mathcal{J}(x(t), t) = \max_{u} \int_{t}^{t+\Delta t} r(x, u) dt + \mathcal{J}(x(t+\Delta t), t+\Delta t)$$

maximised expanding in series and taking $\Delta t \rightarrow 0$

$$-\mathcal{J}(\dot{x(t)},t) = \max_{u} \{r(x,u) + \nabla \mathcal{J}(x(t),t)f(x,u)\}$$

The Hamilton-Jacobi-Bellman eq. (1954) is a functional equation. The solution is the **optimal polic**

ellman (HJB) equation

future value function

The principle of optimality requires that the **future value function** has to be

$$\mathbf{y} \ u = \pi(x)$$

Ínia–HJB in discrete form

Discrete time observations

Bellman optimality equation

7 / 23

We introduce in the HJB the term $\gamma = e^{-\rho t}$ where $\rho > 0$ is the **discount factor**

$$t) = \max_{u} \int_{t}^{T} e^{-\rho t} r(x, u) dt$$

$$u, t) = \max_{u} \{ r(x, u) + \nabla \mathcal{J}(x(t), t) f(x, u) \}$$

 $\mathcal{J}(x_n) = \max\{r(x_n, u) + \gamma \mathcal{J}(x_{n+1})\}$ \mathcal{U}

Ínia Bellman equation

8 / 23

Euromech Colloquium 614

Bellman optimality equation $\mathscr{J}(x_n) = \max\{r(x_n, u) + \gamma \mathscr{J}(x_{n+1})\}$ \mathcal{U}

Ínia Bellman equation

Bellman optimality equation

Evolution of the discounted cost function following an optimal policy

 $V_{\pi}(x_{n}) = r(x_{n}, u_{n}) + \gamma V_{\pi}(x_{n+1})$ $Q_{\pi}(x_n, u_n) = r(x_n, u_n) + \gamma Q_{\pi}(x_n, u_n) + \gamma Q_{\pi}(x_n,$

8 / 23

$$\mathcal{J}(x_n) = \max_{u} \{ r(x_n, u) + \gamma \mathcal{J}(x_{n+1}) \}$$

Value function

$$(x_{n+1}, u_{n+1})$$
 Quality function

Ínia Bellman equation

Bellman optimality equation
$$\mathcal{J}(x_n) = \max_{u} \{r(x_n, u) + \gamma \mathcal{J}(x_{n+1})\}$$

Evolution of the discounted cost function following an optimal policy
 $V_{\pi}(x_n) = r(x_n, u_n) + \gamma V_{\pi}(x_{n+1})$ Value function
 $Q_{\pi}(x_n, u_n) = r(x_n, u_n) + \gamma Q_{\pi}(x_{n+1}, u_{n+1})$ Quality function
Bellman's principle of optimality

Bellman optimality equation
$$\mathcal{J}(x_n) = \max_u \{r(x_n, u) + \gamma \mathcal{J}(x_{n+1})\}$$
Evolution of the discounted cost function following an optimal policy $V_{\pi}(x_n) = r(x_n, u_n) + \gamma V_{\pi}(x_{n+1})$ Value function $Q_{\pi}(x_n, u_n) = r(x_n, u_n) + \gamma Q_{\pi}(x_{n+1}, u_{n+1})$ Quality functionBellman's principle of optimality

The discounted infinite-horizon optimal problem can be decomposed in a series of local optimal problems.

Ínría Bellman equation

Bellman optimality equation
$$\mathcal{J}(x_n) = \max_{u} \{r(x_n, u) + \gamma \mathcal{J}(x_{n+1})\}$$

Evolution of the discounted cost function following an optimal policy
 $V_{\pi}(x_n) = r(x_n, u_n) + \gamma V_{\pi}(x_{n+1})$ Value function
 $Q_{\pi}(x_n, u_n) = r(x_n, u_n) + \gamma Q_{\pi}(x_{n+1}, u_{n+1})$ Quality function

Bellman optimality equation
$$\mathcal{J}(x_n) = \max_{u} \{r(x_n, u) + \gamma \mathcal{J}(x_{n+1})\}$$

Evolution of the discounted cost function following an optimal policy
 $V_{\pi}(x_n) = r(x_n, u_n) + \gamma V_{\pi}(x_{n+1})$ Value function
 $Q_{\pi}(x_n, u_n) = r(x_n, u_n) + \gamma Q_{\pi}(x_{n+1}, u_{n+1})$ Quality function

Bellman's principle of optimality

The discounted infinite-horizon optimal problem can be decomposed in a series of local optimal problems.

Reinforcement Learning

It is possible to learn from the observation of the interaction between the system and the plant

Case study Control of Kuramoto-Sivashinsky system

Bucci M. A., Semeraro O., Allauzen A., Wisniewski G., Cordier L. & Mathelin L.. "Control of chaotic systems by deep reinforcement learning." *Proceedings of the Royal Society A* (2019)

Ínría KS system as a plant

$$\frac{\partial u}{\partial t} = -\nabla^4 u - \nabla^2 u - \frac{1}{2} |\nabla u|$$

If numerical domain length L = 22 chaotic

Plant setup

8 sensors (local measurements)

4 actuators (Gaussian body forcing)

Euromech Colloquium 614

Critic $Q(x_t, u_t | \theta)$

PODMD: partial observable Markov Decision Proces

Euromech Colloquium 614

PODMD: partial observable Markov Decision Proces

Euromech Colloquium 614

Euromech Colloquium 614

Euromech Colloquium 614

 ω_{ij} are the weights of the Actor Neural Network

 ω_{ij} are the θ_{ij} are the

Euromech Colloquium 614

12/23

- ω_{ij} are the weights of the Actor Neural Network
- θ_{ij} are the weights of the Critic Neural Network

 \sum

Bucc

Critic update

256 Neurons 3 Layers

Actor update

128 Neurons 3 Layers

Learn from the observation

 $\nabla_{\omega}Q = \frac{\partial Q_t}{\partial u_t} \frac{\partial \pi_t}{\partial \omega_i}$

$$\omega_i' = \omega_i + \alpha \frac{\partial Q}{\partial \omega_i}$$

Explore the action-state space

Perturbation of the parameters of the policy, **Ornstein-Uhlenbeck process**

Ínia Deep Deterministic Policy Gradient $\nabla_{\theta} TD = \frac{\partial \|Q_t - (r_t + Q_{t+1})\|}{\partial \theta_i}$ $\theta_i' = \theta_i - \alpha \frac{\partial TD}{\partial \theta_i}$ **Critic update**

256 Neurons 3 Layers

Actor update

128 Neurons 3 Layers

Learn from the observation

Remark

13/23

 $\nabla_{\omega}Q = \frac{\partial Q_t}{\partial u_t} \frac{\partial \pi_t}{\partial \omega_i}$ $\omega_i' = \omega_i + \alpha \frac{\partial Q}{\partial \omega_i}$

Explore the action-state space

Perturbation of the parameters of the policy, **Ornstein-Uhlenbeck process**

For the Critic update we don't need to know the policy used to generate a_t : off-policy approach

Silver, D., et al. (2014), ICML Bucci, M.A. et al. (2019), PRSA

14/23

Actor-Critic algorithm: deep deterministic policy gradient (DDPG)

Three policies are identified

Reward defined with respect of the target states

$$r_t = -\|x_t - x_{Target}\|$$

Discount factor

On average: **1 hour training** (Intel I3 CPU, 2015)

Silver, D., et al. (2014), ICML Bucci, M.A. et al. (2019), PRSA

14/23

Actor-Critic algorithm: deep deterministic policy gradient (DDPG)

Three policies are identified

Reward defined with respect of the target states

$$r_t = -\|x_t - x_{Target}\|$$

Discount factor

On average: **1 hour training** (Intel I3 CPU, 2015)

Euromech Colloquium 614

15/23

modal solution ansatz around $x = \hat{x}e^{(\lambda + i\omega)t}$

Uncontrolled KS:

$$(\lambda + i\omega)\hat{x} = J_{KS}\hat{x}$$

two unstable eigenvalues

15/23

modal solution ansatz around $x = \hat{x}e^{(\lambda + i\omega)t}$

Uncontrolled KS:

$$(\lambda + i\omega)\hat{x} = J_{KS}\hat{x}$$

two unstable eigenvalues

Controlled KS:

 $(\lambda + i\omega)\hat{x} = (J_{KS} + \pi_x)\hat{x}$

no unstable eigenvalues

The identified policies are robust with respect of the initial conditions due to the Markiovanity of the Q-function, and the exploration

Case study The fluidic Pinball

<u>Guegan, T.</u>, Bucci, M. A., Semeraro, O., Cordier, L., & Mathelin, L. (2020). Control by Deep Reinforcement Learning of a separated flow. *Bulletin of the American Physical Society*.

Ínría The fluidic pinball case (CFD)

Deng Nan, et al. (2020), JFM

Euromech Colloquium 614

Fluidic Pinball case

Re = **[100**, 150, 200]

Inflow: uniform velocity field

Ínría The fluidic pinball case (CFD)

Deng Nan, et al. (2020), JFM

Euromech Colloquium 614

Fluidic Pinball case Re = [100, 150, 200]Inflow: uniform velocity field

State: **9 pressure** probes Action: **angular velocity**

Ínría The fluidic pinball case (CFD)

Deng Nan, et al. (2020), JFM

Fluidic Pinball case Re = [100, 150, 200]Inflow: uniform velocity field State: 9 pressure probes

Action: angular velocity

$$r(x, u) = -F_D^2 - \alpha ||a_t||_2$$

 $\alpha \in [0, 0.6]$ action **penalty**

 $a_t \in [-3, 3]$ action **bounds**

The employed agent (i.e. learning algorithm) is the **Twin Delayed Deep Deterministic** policy gradient: TD3

$$TD = Q_{\theta}(x_n, u_n) - [r(x_n, u_n) + \max_{u} \gamma Q_{\theta'}(x_{n+1}, u)]$$

With a noisy Q function the max operation can overestimate its evaluation

- Clipped double Q-learning
- **Delayed Policy Update**

4 vortex shedding per epoch and 20 epochs are enough to converge the policy $\pi(x)$.

Fujimoto, Scott, Herke Van Hoof, and David Meger. "Addressing function approximation error in actor-critic methods." *arXiv preprint arXiv:1802.09477* (2018).

19/23

lnia Results $\alpha = 0, 0.6$

Euromech Colloquium 614

20/23

Without action penalty $\alpha = 0,\,65\,\%$ drag reduction

Similar results obtained with genetic algorithm

Guy Y. Cornejo Maceda et al. (2021), JFM

Ínia Results $\alpha = 0, 0.6$

Without action penalty $\alpha = 0,\,65\,\%$ drag reduction

Similar results obtained with genetic algorithm

With action penalty $\alpha = 0.6$ 50% drag reduction

Counter rotation of A_B and A_T

Guy Y. Cornejo Maceda et al. (2021), JFM

lnría Low observability

With 9 probes:

- $\alpha = 0$: drag reduction $\approx 65\%$
- $\alpha = 0.6$: drag reduction $\approx 50\%$

With 3 probes:

• $\alpha = 0.6$: drag reduction $\approx 35\%$

With 1 probe:

• $\alpha = 0.6$: drag reduction $\approx 25\%$

21/23

linia Low observability

With 9 probes:

- $\alpha = 0$: drag reduction $\approx 65\%$
- $\alpha = 0.6$: drag reduction $\approx 50\%$

With 3 probes:

• $\alpha = 0.6$: drag reduction $\approx 35\%$

With 1 probe:

• $\alpha = 0.6$: drag reduction $\approx 25\%$

Reducing the number of probes the **full-observability hypothesis**, required by off-policy approach, is no longer satisfied \rightarrow Introduce **temporal embedding**

Ínría Results $\alpha = 0.6$ + temporal embedding

With 3 probes + temporal embedding (Taken's theorem):

• $\alpha = 0.6$: drag reduction $\approx 50\%$

as for the 9 probes penalised case

Ínría Results $\alpha = 0.6$ + temporal embedding

With 3 probes + temporal embedding (Taken's theorem):

 $\sim \alpha = 0.6$: drag reduction $\approx 50\%$

as for the 9 probes penalised case

The learning is **robust** w.r.t. the initialisation of the Neural Network **parameters** θ and the initial condition x_0

- **Deep reinforcement learning** is a powerful method for **non-linear** control Full knowledge of the system is not required
- of the Bellman equation
- Successfully tested on the KS chaotic system (PRSA, 2019)
- And on fluidic pinball case (APS 2020):
 - Different *Re* numbers tested
 - Fast training \approx 4 vortex shedding per epochs and 20 epochs

- decoupling of the action and reward evaluation location in POMDP
- experimental open cavity

In principle, the policy is a global optimum if the cost function is solution

Big dream: train an off-policy control agent with time delay (credit assignment)

Optimal control theory Lev Pontryagin x = stateu = action $c(u, x) = \cot(-r(s, a))$

Pontryagin maximum principle

To find the optimal policy it is necessary to solve the Hamiltonian system (i.e. variational problem on $[t_0, t_1]$) and a maximum condition of the Hamiltonian control.

Dynamic programming Richard Bellman

- s = state
- a = action
- r(s, a) = reward

Bellman's principle of optimality:

An optimal policy has the property that whatever the initial state and initial decision are, the remaining decision must constitute an optimal policy with regard to the state resulting from the first decision.

We have to solve the variational problem $J = \int_{-\pi}^{\pi} r(x(t), u(t)) dt \quad \text{s.t. } \dot{x} = f(x, u)$

Introducing the Lagrangian multiplier $\lambda(t)$ $H(x, u, \lambda, t) = \lambda^{\mathrm{T}} f(x, u) + r(x, u)$ we can solve:

1. $H_{\lambda} = \dot{x} = f(x, u)$ **2.** $H_x = -\dot{\lambda} = \lambda^T f_x(x, u) + r_x(x, u)$ **3.** $H_{\mu} = 0 \Rightarrow u^* = \pi^*(x^*)$

It is a **necessary** condition for optimality

Euromech Colloquium 614

Given V(x, t) with $x \in \mathbb{R}^n$, $V : \mathbb{R}^n \times \mathbb{R} \to \mathbb{R}$ solution of the Hamilton-Jacobi-Bellman equation,

 $-V(\dot{x},t) = \max(\nabla V(x,t) \cdot f(x,t))$ the optimal action u^* at x_t is the one that maximise $V(x_{t+1})$

It is a sufficient and necessary condition for optimality

$$u) + r(x, u))$$

More efficient (fewer samples)

> model-based shallow RL

model-based deep RL

Picture from Levine's DRL class

