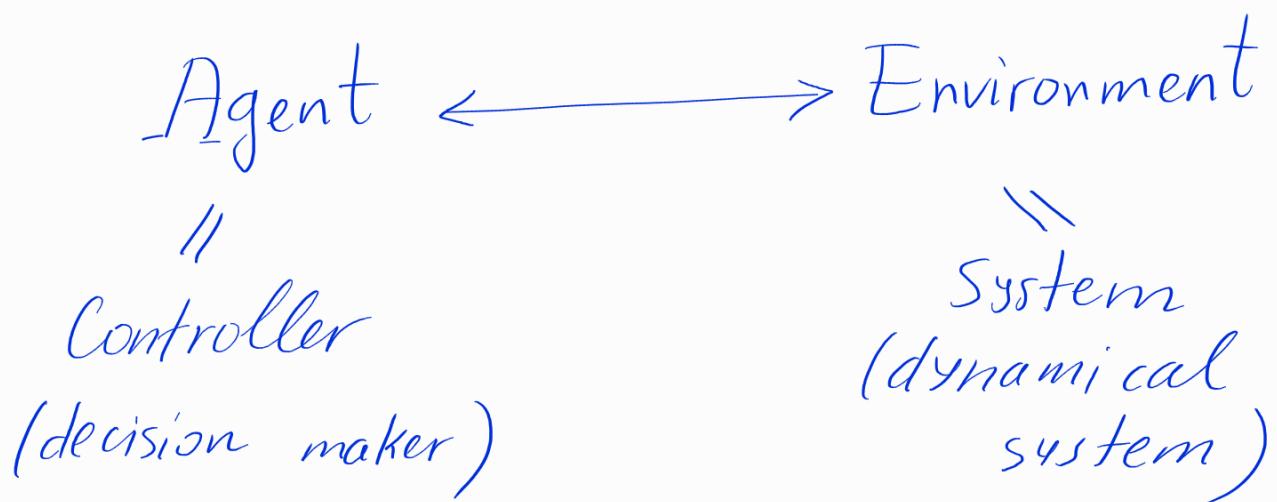


RL 2022 - lec 1

Introduction, context

RL ⊂ Optimal control



Terminology & notation:

x — state

Remark:

u — action

s — state

y — observation

a — action

o — observation

Environment evolution

$$x_{\text{next}} = f(x, u)$$

f - state transition law
(state dynamics function)

(Discrete-time, deterministic env.)

Now, stochastic case:

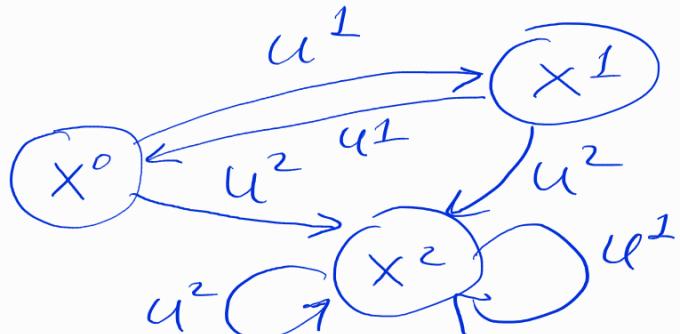
$$X_{\text{next}} \sim f(X_{\text{next}} | x, u)$$



like a probability
of the event $X_{\text{next}} = x_{\text{next}}$
given that the env. was in
the state x and action u
was taken

Example: there is 3 states

x^0, x^1, x^2 and two actions u^1, u^2



$$P[X = x^2 | X = x^1, U = u^2] = 50\%$$

$$x_{\text{next}} = f(x, u) + \eta$$

↑
noise term

$$\eta \sim N(0, \sigma_\eta^2)$$

↑
zero mean



When there is no designated time steps, time is treated continuous.

$$x_{\text{next}} = f(x, u)$$

↓

diff-n
w.r.t.
time t

$$\dot{x} = \frac{dx}{dt} = f(x, u)$$

Case 1

$$x_{\text{next}} = f(x, u)$$

Start at x_0 , x_5 - ?

$$x_1 = f(x_0, u_0)$$

$$x_2 = f(x_1, u_1) = f(f(x_0, u_0), \dots)$$

Case 2

$$\dot{x} = f(x, u)$$

What is the state after T seconds
starting at x_0 ?

The answer :

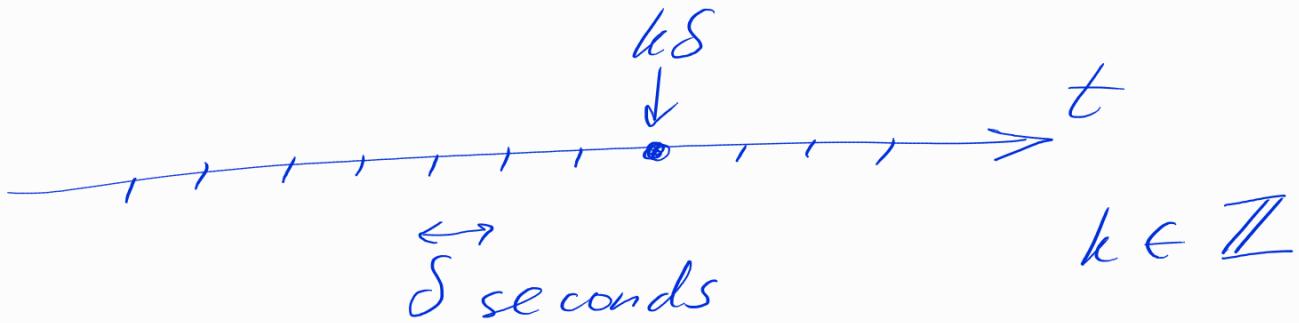
$$x(T) = x_0 + \int_0^T f(x(t), u(t)) dt$$

scipy.integrate

(ode_solve)

How to get from a cont time
to a discrete time?

$$\dot{x} = f(x, u)$$



Euler discretization

$$x_{k+1} = x_k + \delta f(x_k, u_k)$$

$$=: \hat{f}(x_k, u_k)$$

$$x_{k+1} = \hat{f}(x_k, u_k) + \eta_k$$

Why RL?

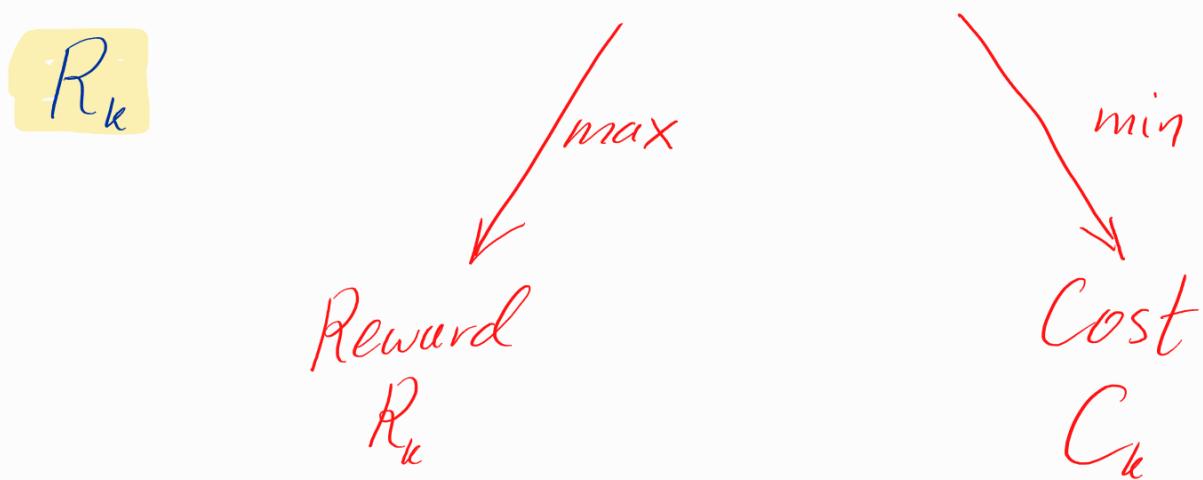
Optimal control (recall)

That is to say optimize

an objective

System: $X_{u+1} \sim f(x_{u+1} | x_u, u_u)$

In a moment in time
(say, at some step k), there
is a running objective



R is an instantaneous
performance mark

(it answers the question: how
good am I right now?)

Optimal control is usually
concerned with optimization
not just at a moment in

time, but in all moments
in time.

This can be formulated as:

sequence
of actions
 \downarrow
 $\{u\}_0^T$

$$\max_{\rho} \sum_{k=0}^T R_k$$

where T may be ∞ .

\downarrow
horizon

ρ -policy! (control law,
decision law,
action plan etc.)

Deterministic

$$\rho(x)$$

Stochastic

$$\rho(u/x)$$

Let's recap what we have so far:

$$\max_{\rho} \sum_{k=0}^{\infty} \mathbb{E}[\gamma^k R_k], \quad 0 < \gamma \leq 1$$

discount factor

$$\text{s.t. } X_{k+1} \sim f(x_{k+1} | x_k, u_k)$$

$$R_k \sim v(r_k | x_k, u_k)$$

$$U_k \sim p(u_k | x_k)$$

r - reward law

But, remember, the agent in general only sees an observation.

$$\max_{\rho} \sum_{k=0}^{\infty} \mathbb{E}[\gamma^k R_k]$$

$$\text{s.t. } X_{k+1} \sim f(x_{k+1} | x_k, u_k)$$

$$R_k \sim v(r_k | x_k, u_k)$$

$$U_k \sim p(u_k | y_k)$$

$$Y_k \sim h(y_k | x_k)$$

For simplicity, assume $\gamma=1$.

$$G_{k:k+N} = \sum_{i=k}^{k+N-1} R_i \quad \text{or}$$

$G_k = \sum_{i=k}^{\infty} R_i$ is called an outcome



Objective $\xrightarrow{\max}$ Value

$\xleftarrow{\min}$ Cost-to-go

$$\mathbb{E}[G_k]$$

$$\mathbb{V}^\rho = \mathbb{E}_{\text{the } S_u \sim \rho} [G]$$

\uparrow
action sequence
was distributed by ρ

$$V^* = \max_{\mathcal{P}} V^{\mathcal{P}}$$

optimal value

In max problems: $J^{\mathcal{P}}, J^*$

Optimal policy: \mathcal{P}^*

$$\mathcal{P}^* = \arg \max_{\mathcal{P}} V^{\mathcal{P}}$$