# I. DEFINITIONS AND PRINCIPLE OF OPTIMALITY

System characterized by:

- State $x \in \mathcal{X}$
- Control $u \in \mathcal{U}$
- "evolution operator": $x_{k+1} = \text{next}(x_k, u_k)$

GOAL: Minimize a __cost__ function $J(\{x\}, \{u\})$
by choosing $\{u\}$

$\vee$
· sequences of states/controls

Def. $\pi(x) = u$ is called · __control law__ · (or ~~policy~~)
and implements __closed loop control__

Under the __assumption__ that:

$$J(\{x\}, \{u\}) = \sum_{k=0}^{k_f - 1} \text{cost}(x_k, u_k) \qquad ①$$

$\hookrightarrow$ cost of ~~taking action~~
using control $u_k$ in state $x_k$

i.e. __cost__ is separable in time

Then
$\Longrightarrow$ it holds the __PRINCIPLE OF OPTIMALITY__:
(optimal substructure)

~~If~~ $\Gamma_{ac}^*$ achieves minimum cost ~~from~~ $a$ to $c$

$\Longrightarrow \Gamma_{bc}$ " " " from $b$ to $c$

$\mathcal{X}$

$\Gamma_{ac}$: short for
the seq. of states and controls that
induces the trajectory

Proof: (by contradiction)
If $\exists\, \Gamma_{bc}'$ s.t. $J(\Gamma_{bc}') < J(\Gamma_{bc})$

$\Longrightarrow J(\Gamma_{ab}) + J(\Gamma_{bc}') < J(\Gamma_{ab}) + J(\Gamma_{bc}) = J(\Gamma_{ac}^*)$

① $\longrightarrow$ "
$$J(\Gamma_{ac}') < J(\Gamma_{ac}^*)$$

①

CONTRADICTION!

<u>Def</u> INFINITE HORIZON PROBLEM :
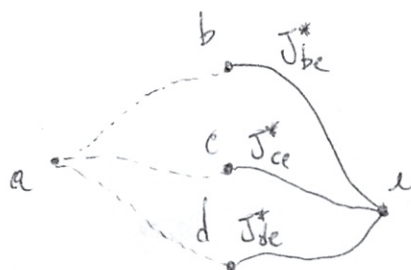
Optimal Control problem with no given endpoint

$$\Rightarrow J(\{x\}, \{u\}) = \sum_{k=0}^{\infty} \gamma^k \, cost(x_k, u_k) \quad , \quad \gamma < 1$$

---

# II. DYNAMIC PROGRAMMING (DP)

<u>DP</u>: Set of algorithms that exploit the optimality principle to solve optimization problems

<u>MOTIVATION</u>: Direct enumeration of all possible trajectories to find minimum is <u>exponential</u> <u>in the number of timesteps</u>

<u>INSTEAD</u>, consider:



Which of $\Gamma_{ab}$, $\Gamma_{ac}$, $\Gamma_{ad}$ is optimal?

If we know $J^*_{be}, J^*_{ce}, J^*_{de}$, it's sufficient to compute

$$J_{ae} = \min_{i=b,c,d} \left( J_{ai} + J^*_{ie} \right)$$

by starting from the endpoint, one can REUSE previously computed quantities

$\Rightarrow$ DP scales <u>linearly with the number of timesteps</u>

Def $V_{\tilde{\pi}}(x)$ : <u>value</u> of a state, i.e. total cost starting from $x$ to $x(T)$, following the control law $\tilde{\pi}(x)$

Def $V^*(x) := V_{\tilde{\pi}^*}(x)$, $\tilde{\pi}^*$: optimal control law

$\implies$ Iterative application of the optimality principle lead to:

$$V^*(x) = \min_u \left[ cost(x,u) + V^*(next(x,u)) \right] \qquad (2) \qquad \boxed{\text{Bellman equation}}$$

$$\tilde{\pi}^*(x) = \text{argmin}_u \left[ cost(x,u) + V^*(next(x,u)) \right]$$

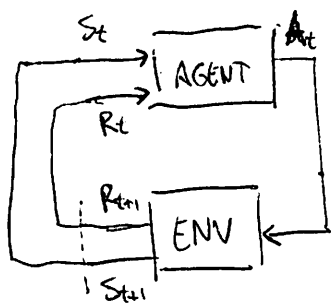Bellman eq. is a <u>RECURSIVE EQ</u> for the <u>OPTIMAL VALUE FUNCTION</u>

---

### III. a) WHAT is "next(x,u)"?

Optimal control : Dynamical system $\dot{\vec{x}} = f(\vec{x}, \vec{u}, t)$

OR

$$d\vec{x} = f(\vec{x}, \vec{u}, t) dt + F(\vec{x}, \vec{u}) d\vec{w}$$

Reinforcement learning (RL) : • Markov decision process (MDP)



$S_t, R_t, A_t$ are <u>random variables</u>, fully characterized by

$$P(s', r \mid s, a) : P(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$$

only previous state $\implies$ MARKOVIAN

b) WHAT is "cost(x,u)"?

- Optimal Control · Arbitrary function $\quad J = h(x(t_f)) + \int_0^{t_f} \ell(x(t), u(t), t)\, dt$

  EXAMPLES : Minimum effort ; Minimum time ; Endpoint control; tracking problems ...

- RL : ① Instead of min. cost, we want <u>max</u> Total <u>return</u> (EXPECTED)

  $\underline{Def}\quad G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad , \gamma < 1 \qquad \underline{Total\ return}$

  $\underline{Def}\quad \pi(a\,|\,s)$ POLICY , prob. of action $a$ given state $s$

  GOAL : Maximize $\quad E_\pi\left[ G_t \right]$

  $\underset{(VALUE)}{\underline{Def}}\quad V_\pi(s) = E_\pi\left[ G_t \mid S_t = s \right] = E_\pi\left[ R_{t+1} + \gamma\, V(S_{t+1}) \mid S_t = s \right] \qquad (3)$

- $V^*(s)$ : Optimal value function (achieves max return)

  $V^*(s)$ obeys

  $$\boxed{\begin{aligned} V^*(s) &= \max_a E\left[ R_{t+1} + \gamma \cdot V^*(S_{t+1}) \mid S_t = s,\ A_t = a \right] \\ &= \max_a \sum_{s',r} p(s',r \mid s,a)\left[ r + \gamma V^*(s') \right] \end{aligned}}$$

  Bellman optimality
  eq. in RL

  (4)

(4) is a set of $|S|$ nonlinear eqs. $\Rightarrow$ $\underline{\text{iterative}}$ methods to solve it

$\Rightarrow$ DP methods:

- $\underline{\text{Policy iteration}}$ : Alternate between:

  - Policy evaluation : $\pi(s) \longrightarrow V(s)$ by repeatedly use (3)

  - Greedy $\underline{\text{policy improvement}}$ :

    $$\tilde{\pi}(s) = \underset{a}{\text{argmax }} V(s)$$

- $\underline{\text{Value iteration}}$ : Turn (4) into an $\underline{\text{iterative update}}$ :

  $$V^{(k+1)}(s) = \underset{a}{\max} \left\{ \sum_{s',r} p(s',r|s,a) \left[ r + \gamma V^{(k)}(s') \right] \right.$$

  Can be seen as a policy iteration, with policy evaluation step truncated after one update

D.P. - Pros

- Reduces complexity thanks to principle of optimality
  (BOOTSTRAPPING)

D.P. - Cons

- Requires knowledge of $p(s',r|s,a)$
  (MODEL-BASED)

System: $\quad \vec{\delta x} = \vec{f}(\vec{x}, \vec{u}) \, \delta t + F(\vec{x}, \vec{u}) \, \vec{\delta W}$

$\qquad\qquad\qquad\qquad\qquad\qquad$ ↳ Brownian motion

$\qquad \vec{u}, \vec{x} \in \mathbb{R}^n, \quad F(\vec{x}, \vec{u}) \in \mathbb{R}^{n \times n}$

discretization
$\Longrightarrow \quad \vec{x}_{k+1} = \vec{x}_k + \Delta \cdot \vec{f}(\vec{x}_k, \vec{u}_k) + \sqrt{\Delta} \cdot F(\vec{x}_k, \vec{u}_k) \cdot \varepsilon_k \quad , \quad \varepsilon_k \sim N(0, I_n)$

$\qquad\qquad\qquad$ ↳ time step

Cost:

$$J(\vec{x}(\cdot), \vec{u}(\cdot)) = h(\vec{x}(t_f)) + \int_0^{t_f} \ell(\vec{x}(t), \vec{u}(t), t) \, dt \to h(x(t_f)) + \sum_{k=0}^{K-1} \ell(\vec{x}_k, \vec{u}_k, k \cdot \Delta) \Delta$$

$\Longrightarrow \underline{\vec{x}_{k+1} = \vec{x}_k + \Delta \cdot \vec{f}(\vec{x}_k, \vec{u}_k) + \xi} \quad , \quad \xi \sim N(0, \Delta S) \;, \; S = F(\vec{x}, \vec{u}) F^\top(\vec{x}, \vec{u})$

$\qquad\qquad$ Plug into Bellman eq (2)

$V(\vec{x}, k \cdot \Delta) = \min_{\vec{u}} \left\{ \Delta \cdot \ell(\vec{x}_k, \vec{u}_k, \Delta \cdot k) + E\left[ V\left( \vec{x}_k + \Delta \cdot \vec{f}(\vec{x}_k, \vec{u}_k) + \xi, (k+1)\Delta \right) \right] \right\}$

IDEA: Expand $V$ up to order $\Delta$.

$V(\vec{x} + \vec{s}) = V(\vec{x}) + \vec{s}^\top \cdot \vec{\partial}_x V(\vec{x}) + \frac{1}{2} \vec{s}^\top \cdot \partial_{xx} V(\vec{x}) \cdot \vec{s} \qquad$ (time index hidden)

$\vec{s} := \Delta \vec{f}(\vec{x}_k, \vec{u}_k) + \xi$

$\Rightarrow E[V(\vec{x} + \vec{s})] = V(\vec{x}) + \Delta \cdot \vec{f}^\top(\vec{x}_k, \vec{u}_k) \cdot \vec{\partial}_x V(\vec{x}_k) + E\left[ \frac{1}{2} \vec{\xi}^\top \cdot \partial_{xx} V(\vec{x}_k) \vec{\xi} \right]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underbrace{\frac{1}{2} E\left( \Delta \cdot \vec{\xi} \vec{\xi}^\top \partial_{xx} V(\vec{x}) \right)}_{} = \frac{\Delta}{2} Tr\left( S \cdot \partial_{xx} V(\vec{x}) \right)$

↳ $V(\vec{x}, k \cdot \Delta) - V(\vec{x}, (k+1) \cdot \Delta) = \min_{\vec{u}} \left\{ \Delta \cdot \ell(\vec{x}, \vec{u}, \Delta \cdot k) + \Delta \vec{f}^\top(\vec{x}, \vec{u}) \cdot \vec{\partial}_x V(\vec{x})(k+1)\Delta) + \frac{\Delta}{2} Tr\left( S \partial_{xx} V \right) \right\}$

$$\lim_{\Delta \to 0} \implies -\partial_t V(\vec{x},t) = \min_{\vec{u}}\left\{ \ell(\vec{x},\vec{u},t) + \vec{f}^{\,T}(\vec{x},\vec{u})\cdot\vec{\partial}_x V(\vec{x},t) + \frac{1}{2}\,\mathrm{Tr}\left(S\cdot\partial_{xx} V(\vec{x},t)\right)\right\} \quad \textcircled{7}$$

(5)

$$\boxed{\text{Hamilton - Jacobi - Bellman} \quad \text{Eq}} \quad (HJB)$$

· HJB can be applied to any system / cost $\varnothing$, including stochastic ones,

but PDE solvers are <u>exponential</u> in <u>M</u> $\implies$ <u>CURSE OF DIMENSIONALITY</u>

## VI. Detour: Maximum Principle from HJB (informal) —————

For a <u>deterministic</u> system

$$-\partial_t V(\vec{x},t) = \min_{\vec{u}}\left\{ \ell(\vec{x},\vec{u},t) + \vec{f}^{\,T}(\vec{x},\vec{u})\cdot\vec{\partial}_x V(\vec{x},t)\right\}$$

<u>Def</u>. $\vec{p} := \vec{\partial}_x V(\vec{x},t)$     COSTATE

<u>Def</u>. $H(\vec{x},\vec{p},\vec{u},t) := \ell(\vec{x},\vec{u},t) + \vec{f}^{\,T}(\vec{x},\vec{u})\cdot\vec{p}$

$$\implies -\partial_t V(\vec{x},t) = \min_{\vec{u}}\left\{ H(\vec{x},\vec{p},\vec{u},t)\right\}$$

$\implies$ Differs from H-J eq. of classical mechanics only because of "$\min_u$"

$\implies$ One can show that, as in Classical mechanics:

$$\begin{cases} \dot{x}_i = \dfrac{\partial H}{\partial p_i} \\[2mm] \dot{p}_i = -\dfrac{\partial H}{\partial x_i} \\[2mm] \vec{u} = \underset{\vec{u}}{\arg\min}\left[ H(\vec{x},\vec{p},\vec{u},t)\right] \end{cases}$$

with init. cond. $x_i(0)$
and $n$ boundary cond. $p_i(t_f) = \dfrac{\partial h(\vec{x}(t_f))}{\partial x_i}$

$\implies$ Systems of $n$ ODEs $\implies$ Avoids curse of dim.!

BUT : · No noise
· init. cond. fixed

# VII. Final Observations

- HJB can be solved ANALYTICALLY for Linear-Quadratic Gaussian regulator (LQG)

  $\Rightarrow$ Used as an approximation technique for NONLINEAR PROBLEMS : iLQG :

  Given $\vec{u}^{(k)}(t)$

  - Run nonlinear dyn. forward $\xrightarrow{\text{get}}$ $\bar{x}(t)$, J

  - Fit LQG approx. of $\bar{x}(t)$, J ~~xx~~ .

  - Solve analytically for $\hat{u}$

  - Set $\vec{u}^{(k+1)}(t) = \vec{u}^{(k)}(t) + \hat{u}(t)$

- RL — MODEL-FREE METHOD

  Combine Optimality principle + Montecarlo estimation of Value

  $\Rightarrow$ TD learning

  $$V(S_t) \leftarrow V(S_t) + \alpha \left[ R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right]$$

  $$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right]$$

  Requires $\{ S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1} \} \Rightarrow$ SARSA