

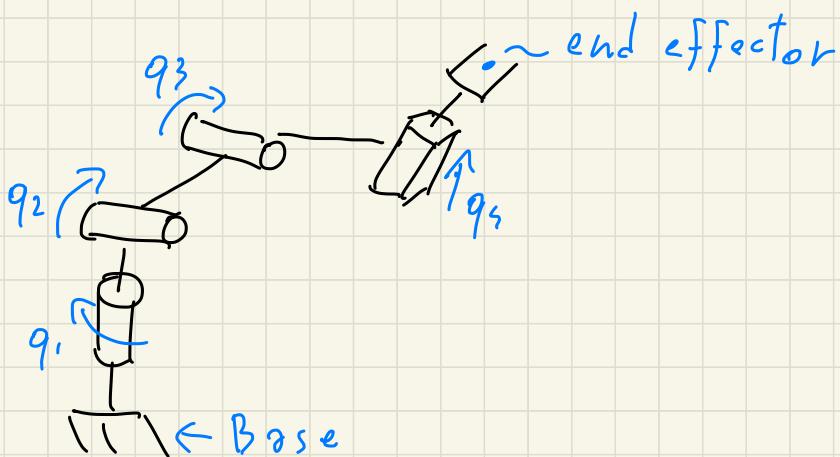
# ROBOT DYNAMIC MODELING

Definitions:

- Kinematics (forward)
- Statics
- Dynamics
- Homogeneous matrix
- Jacobian
- . . .
- $\ddot{M}\ddot{q} + C\dot{q} + g(q) = \tau + J^T u$

# ROBOT MANIPULATOR

- Set of rigid bodies connected by joints
  - revolute
  - prismatic
- Typically joints are actuated by system of motor + transmission (e.g., gear box)

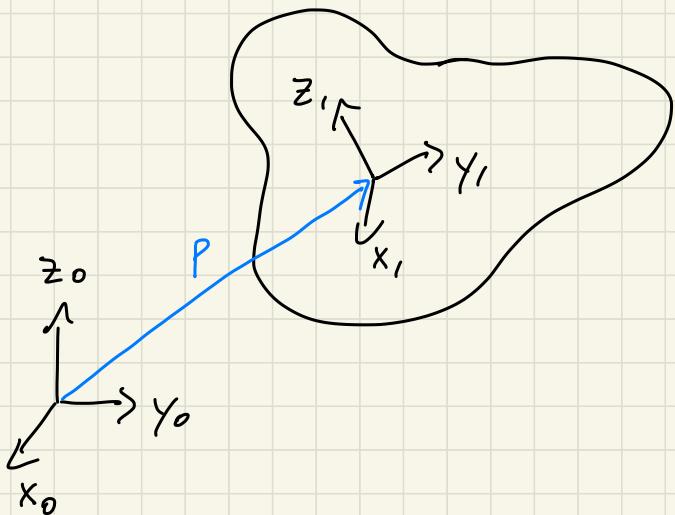


$q_i$ : i-th joint coordinate

$$q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_n \end{bmatrix}$$

# HOMOGENEOUS TRANSFORMATION

Def: **Pose**  $\triangleq$  position and orientation

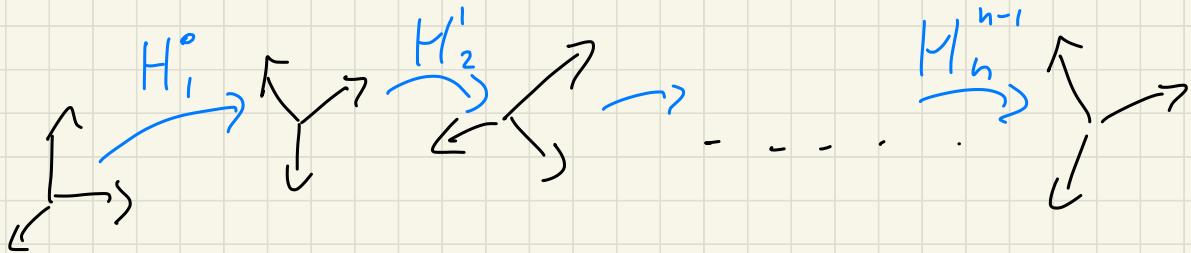


$$H_i^o = \begin{bmatrix} R_i^o & p \\ 0_{1 \times 3} & 1 \end{bmatrix}$$

$$R_i^o = [x_i^o \ y_i^o \ z_i^o]$$

$$R_i^o = R_o^{i^{-1}} = R_i^o{}^T$$

# COMPOSITION OF TRANSFORMATION



$$H_n^o = H_1^o H_2^1 H_3^2 \dots H_n^{n-1}$$

$H_i^{i-1}(q_i) \Rightarrow$  depends on i-th joint angle only

$$H_n^o(q_1, q_2, \dots, q_n) = H_n^o(q)$$

# REPRESENTATION OF ORIENTATION

- Rotation matrix  $[3 \times 3]$
  - Euler angles  $[3]$
  - Roll-Pitch-Yaw  $[3]$
  - Angle-axis  $[5]$
  - Quaternions  $[4]$
- $\left. \begin{matrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{matrix} \right\}$  Minimal representations  
 $\left. \begin{matrix} \cdot \\ \cdot \\ \cdot \end{matrix} \right\}$  Have singularities

# DIFFERENTIAL KINEMATICS

## GEOMETRICAL APPROACH

We want to get  ${}^0V_e = \begin{bmatrix} \dot{P}_e \\ \omega_e \end{bmatrix}$  linear vel  
angular vel

$${}^0\dot{P}_n = {}^0\dot{P}_n + {}^0\dot{P}_{n-1} = {}^0\dot{P}_n + {}^{n-2}\dot{P}_{n-1} + {}^0\dot{P}_{n-2} = \sum_{i=0}^{n-1} {}^i\dot{P}_{i+1}$$

Rotation axis of joint i

$${}^{i-1}\dot{P}_i = \begin{cases} \text{Revolute: } [{}^i\hat{z}_i \times (\dot{P}_i - \dot{P}_{i-1})] \dot{q}_i = J_{p_i} \dot{q}_i \\ \text{Prismatic: } {}^i\hat{z}_i \dot{q}_i = J_{o_i} \dot{q}_i \end{cases}$$

$$J_p \triangleq [J_{p_1} \ J_{p_2} \ \dots \ J_{p_n}]$$

$${}^0\dot{P}_n = \sum_{i=0}^{n-1} J_{p_{i+1}} \dot{q}_{i+1} = \overline{J_p} \dot{q}$$

$$J_o \triangleq [J_{o_1} \ J_{o_2} \ \dots \ J_{o_n}]$$

$${}^0\dot{\omega}_n = \sum_{i=0}^{n-1} J_{o_i} \dot{q}_i = \overline{J_o} \dot{q}$$

Dimension:  $6 \times N$ , where N is the number of joints

$${}^0V_e = \begin{bmatrix} \overline{J_p} \\ \overline{J_o} \end{bmatrix} \dot{q} = J(\dot{q}) \dot{q}$$

Jacobian matrix

# DIFFERENTIAL KINEMATICS

## ANALYTICAL APPROACH

$$\phi_n = \begin{bmatrix} x_n \\ y_n \\ z_n \\ \varrho_n \\ \theta_n \\ \psi_n \end{bmatrix} = \phi_n(q)$$

$$\dot{\phi} = \frac{\partial \phi}{\partial q} \cdot \frac{dq}{dt} = \frac{\partial \phi}{\partial q} \dot{q}$$

↑  
Analytical Jacobian  $J_A$

$$\begin{bmatrix} J_P \\ J_{A_0} \end{bmatrix} = J_A \neq J = \begin{bmatrix} J_P \\ J_O \end{bmatrix}$$

# STATICS

Problem: given wrench (Force + moment) applied at the end-effector  $\Rightarrow$  compute joint torques generated by this wrench

$$W = \begin{bmatrix} F \\ m \end{bmatrix}$$

$$\text{Power at the joints: } P_j = \tau^T \dot{q}$$

$$\text{Power at end-effector: } P_e = f^T \dot{p} + m^T \dot{w} = [f^T \ m^T] \begin{bmatrix} \dot{p} \\ \dot{w} \end{bmatrix} =$$

We can simplify  $\dot{q}$  because we know that this is valid for each joint  $\tau^T v = W^T v = W^T J \dot{q}$

$$\tau^T \dot{q} = W^T J \dot{q} \Rightarrow \tau^T = W^T J \Rightarrow \tau = J^T W$$

Should hold  $\forall \dot{q}$

# DYNAMICS

Problems:

- Direct dynamics:  $q, \dot{q}, \tau \Rightarrow \ddot{q}$  SIMULATION
- Inverse dynamics:  $q, \dot{q}, \ddot{q} \Rightarrow \tau$  CONTROL

$$M(q) \ddot{q} + \underbrace{((q, \dot{q}) \dot{q} + g(q))}_{h(q, \dot{q})} = \tau + J^T(q) w$$

- $M(q) \in \mathbb{R}^{n \times n}$ : mass matrix, positive-definite, symmetric
- $((q, \dot{q})) \in \mathbb{R}^{n \times n}$ : centrifugal + coriolis effects
- $g(q) \in \mathbb{R}^n$ : gravity torques

-Nonlinear equations

-Linear in  $\ddot{q}, \tau, w$  if  $(q, \dot{q})$  is fixed

# JOINT MOTION CONTROL

Given a reference joint trajectory  $q^{\text{ref}}(t)$ , compute  $\tau(t)$  such that  $q(t) \approx q^{\text{ref}}(t)$ .

**Assumption:** We know system dynamics  $M(q)\ddot{q} + h(q, \dot{q}) = \tau$

## IDEA #1: OPEN-LOOP CONTROL

- Compute  $\dot{q}^{\text{ref}}(t)$ ,  $\ddot{q}^{\text{ref}}(t)$
- Compute  $\tau(t) = M(q^{\text{ref}}(t))\ddot{q}^{\text{ref}}(t) + h(q^{\text{ref}}(t), \dot{q}^{\text{ref}}(t))$
- Assuming  $q(0) = q^{\text{ref}}(0)$ ,  $\dot{q}(0) = \dot{q}^{\text{ref}}(0)$  then  
$$\cancel{M(q)\ddot{q} + h(q, \dot{q}) = \tau} = \cancel{M(q^{\text{ref}})\ddot{q}^{\text{ref}} + h(q^{\text{ref}}, \dot{q}^{\text{ref}})} \Rightarrow \ddot{q} = \ddot{q}^{\text{ref}}$$

## #2: PID CONTROL

$$\tau(t) = K_p(q^{\text{ref}}(t) - q(t)) + K_d(\dot{q}^{\text{ref}}(t) - \dot{q}(t)) + K_i \int_0^t (q^{\text{ref}}(s) - q(s)) ds$$

$K_p, K_d, K_i$  are diagonal positive-definite matrices

Guarantee that  $q(t) \rightarrow q^{\text{ref}}(t)$  if  $\dot{q}^{\text{ref}} = \emptyset$

$$M\ddot{q} + C\dot{q} + g = \tau = K_p e + K_d \dot{e} + K_i \int e ds$$

At steady state ( $\dot{q} = \ddot{q} = 0$ ) I have

$$g = K_p e + K_i \int e ds \Rightarrow K_p \dot{e} + K_i e = 0$$

Stable first-order linear dynamical system

$$\dot{e} = -K_i^{-1} K_p e$$

$$\begin{aligned} e &\rightarrow 0 \\ q &\rightarrow q^{\text{ref}} \end{aligned}$$

$$\text{If } k_r = 0 \Rightarrow k_p e = g \Rightarrow e = k_p^{-1} g$$

### #3 : PD + GRAVITY COMPENSATION

$$\tau(t) = K_p (q^{\text{ref}}(t) - q(t)) + K_d (\dot{q}^{\text{ref}}(t) - \dot{q}(t)) + g(q(t))$$

$e$                            $\dot{e}$

Better performance than PID alone.

Steady-state analysis:

$$g = K_p e + g \Rightarrow K_p e = 0 \Rightarrow e = 0$$

# # 5 : INVERSE DYNAMICS CONTROL

$$\ddot{\gamma} = M(\ddot{q}^{\text{ref}} + k_p e + k_d \dot{e}) + h$$

A.V.A. FEEDBACK LINEARIZATION  
COMPUTED TORQUE

System closed-loop dynamics:

~~$$M\ddot{q} + h = M(\ddot{q}^{\text{ref}} + k_p e + k_d \dot{e}) + h$$~~

$$\ddot{q} = \ddot{q}^{\text{ref}} + k_p e + k_d \dot{e}$$

$$\ddot{e} + k_d \dot{e} + k_p e = 0$$

$$y \triangleq \begin{bmatrix} e \\ \dot{e} \end{bmatrix}$$

$$\dot{y} = \begin{bmatrix} \dot{e} \\ \ddot{e} \end{bmatrix}$$

$$\dot{y} = \begin{bmatrix} 0 & I \\ -k_p & -k_d \end{bmatrix} y$$

## # VARIANT OF INVERSE DYNAMICS

$$\tau = M \ddot{q}^{\text{ref}} + h + K_p e + K_d \dot{e}$$

Closed-loop behavior:

$$M \ddot{q} + h = M \ddot{q}^{\text{ref}} + h + K_p e + K_d \dot{e}$$

$$\ddot{q} = \ddot{q}^{\text{ref}} + \tilde{M}' K_p e + \tilde{M}' K_d \dot{e}$$

$$\ddot{q}^{\text{ref}} - \ddot{q} + \bar{K}_p e + \bar{K}_d \dot{e} = 0$$

$$\begin{aligned}\bar{K}_p &\triangleq \tilde{M}' K_p | > 0 \\ \bar{K}_d &\triangleq \tilde{M}' K_d | > 0\end{aligned}$$

Positive definite

Holds because  $K_p$  and  $K_d$  are diagonal

## CARTESIAN-SPACE MOTION CONTROL

Given a reference trajectory for end-effector  $x^{\text{ref}}(t)$ ,  
compute  $\tau(t)$  such that  $x(t) \approx x^{\text{ref}}(t)$ .

[Assume  $x \in \mathbb{R}^3$  represents the e-e position]

Inverse Kinematics

IDEA #1: Compute  $q^{\text{ref}}(t)$  such that  $\overbrace{x^{\text{ref}}(t)}^{\text{You cannot do this computation offline}} = \underbrace{FG(q^{\text{ref}}(t))}_{\text{Forward Geometry}}$

where  $FG(\cdot)$  is forward geometry function.

Then apply joint motion control.

# OPERATIONAL SPACE CONTROL

Method to do cartesian space control

Operational space dynamics

$$J\tilde{M}' \left[ M\ddot{q} + h = \tau \right]$$

Dynamics of the robot  
in joint space

$$\dot{J}\ddot{q} + J\tilde{M}'h = J\tilde{M}'\tau$$

$$\ddot{x} - \dot{J}\dot{q} + J\tilde{M}'h = J\tilde{M}'\tau$$

$$\Lambda \ddot{x} + \underbrace{\Lambda (J\tilde{M}'h - \dot{J}\dot{q})}_{\mu} = \Lambda J\tilde{M}'\tau$$

$$\Lambda \ddot{x} + \mu = \Lambda J\tilde{M}'\tau$$

$$\Lambda(q)\ddot{x} + \nu(q,\dot{q}) = f$$

Pseudoinverse

$$\bar{\Lambda} = (J\tilde{M}'J^T + \varepsilon I)^{-1}$$

$$\varepsilon > 0$$

$$\varepsilon = 10^{-5}$$

$$\ddot{x} = J\ddot{q} + \dot{J}\dot{q}$$

$$\Lambda \triangleq (J\tilde{M}'J^T)^{-1}$$

Assume  $\tau = J^T f$

$$[\Lambda J\tilde{M}'J^T = I]$$

J is full  
row rank

Choose  $f = \underbrace{A(\ddot{x}^{\text{ref}} + K_p(x^{\text{ref}} - x) + K_d(\dot{x}^{\text{ref}} - \dot{x}))}_{\ddot{x}^d} + \nu$

$$\begin{aligned}\tau &= J^T f = J^T(A \ddot{x}^d + \nu) = \\ &= J^T A \ddot{x}^d + J^T A (J M^{-1} h - J \ddot{q})\end{aligned}$$

If I apply this control law I get  $\ddot{x} = \ddot{x}^d$

A simplified version of this control law is:

$$\tau = J^T A \ddot{x}^d + h$$

Let's check it works.

$$M \ddot{q} + h = J^T A \ddot{x}^d + h$$

$$\ddot{\mathbf{J}}\ddot{\mathbf{q}} + \ddot{\mathbf{J}}\ddot{\mathbf{M}}^T \ddot{\mathbf{h}} = \ddot{\mathbf{J}}\ddot{\mathbf{M}}^T \ddot{\mathbf{J}}^T \ddot{\mathbf{A}} \ddot{\mathbf{x}}^d + \ddot{\mathbf{J}}\ddot{\mathbf{M}}^T \ddot{\mathbf{h}}$$

$$\ddot{\mathbf{J}}\ddot{\mathbf{q}} = \ddot{\mathbf{x}}^d$$

$$\ddot{\mathbf{x}} - \underbrace{\ddot{\mathbf{J}}\ddot{\mathbf{q}}}_{\text{Typically small}} = \ddot{\mathbf{x}}^d \Rightarrow \ddot{\mathbf{x}} \simeq \ddot{\mathbf{x}}^d$$

## REDUNDANCY

$$A\ddot{x} + \nu = (J\tilde{M}'J^T)^{-1}J\tilde{M}'\tau = J^{T\#}\tau$$

$J^{T\#}$  is the pseudo-inverse of  $J^T$  because  $J^{T\#}J^T = I$   
 $J^T J^{T\#} \neq I$

Any  $\tau_0$  such that  $J^{T\#}\tau_0 = 0$  does not affect  $\ddot{x}$ .

$$A\ddot{x} + \nu = J^{T\#}(J + \tau_0) = J^{T\#}\tau$$

I can compute  $\tau_0$  as  $\tau_0 = (I - J^T J^{T\#})\tau_1$

Verify  $J^{T\#}\tau_0 = J^{T\#}(I - J^T J^{T\#})\tau_1 = (J^{T\#} - J^{T\#}J^T J^{T\#})\tau_1$

$$= 0$$

So I can extend my control law as:

$$\tau = \bar{J}^T \bar{J} \ddot{x}^d + h + (I - \bar{J}^T \bar{J}^{T\#}) \tau_i$$

How do I compute  $\tau_i$ ?

If I have a reference posture  $q^{ref}$  I can compute  $\tau_i$  as

$$\tau_i = M(K_p(q^{ref} - q) - K_d \dot{q})$$

# INVERSE KINEMATICS + INVERSE DYNAMICS

$$\ddot{x} = J \ddot{q} + \dot{J} \dot{q} \Rightarrow \ddot{q}^d = J^\# (\ddot{x}^d - \dot{J} \dot{q})$$

$$[J J^\# = I]$$

$$\tau = M \ddot{q}^d + h$$

$$J^\# = J^T (J J^T)^{-1}$$

Verify closed-loop behavior:

~~$$M \ddot{q} + h = M \ddot{q}^d + h$$~~

$$\ddot{q} = J^\# (\ddot{x}^d - \dot{J} \dot{q})$$

~~$$J \ddot{q} = J J^\# (\ddot{x}^d - \dot{J} \dot{q})$$~~

$$J \ddot{q} + \dot{J} \dot{q} = \ddot{x}^d \Rightarrow \ddot{x} = \ddot{x}^d$$

ASSUMPTIONS:

- [J is full row rank]
- [All rows are linearly independent]

# Q.P.-BASED REACTIVE CONTROL

- Joint space control
- Review of convex optimization
- Task-Space control
- Underactuated robots
- Contact constraints
- Multi-Task control

# JOINT-SPACE INVERSE DYNAMICS

$$M(q) \ddot{q} + h(q, \dot{q}) = \tau$$

Given  $q^r(t)$ , find  $\tau(t)$  such that  $q(t) \approx q^r(t)$

$$\tau^d = M \ddot{q}^d + h \quad \ddot{q}^d = \ddot{q}^r + k_d (\dot{q}^r - \dot{q}) + k_p (q^r - q)$$


---

$$(\tau^*, \ddot{q}^*) = \arg \min_{\tau, \ddot{q}} \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2 = \frac{1}{2} \ddot{q}^T \ddot{q} - \ddot{q}^d T \ddot{q} + \frac{1}{2} \cancel{\ddot{q}^T \cancel{\ddot{q}}}$$

constant

subject to  $M \ddot{q} + h = \tau$

$$\left. \begin{array}{l} \tau^* = \tau^d = M \ddot{q}^d + h \\ \ddot{q}^* = \ddot{q}^d \end{array} \right\} \text{SOLUTION}$$

# TORQUE, CURRENT, VELOCITY LIMITS

$$\underset{\tau, \dot{q}}{\text{minimize}} \quad \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2$$

$$\text{subject to} \quad M\dot{q} + h = \tau$$

$$\tau^{\min} \leq \tau \leq \tau^{\max}$$

$$i^{\min} \leq k_T \tau \leq i^{\max}$$

$$\dot{q}^{\min} \leq \dot{q} + \Delta t \ddot{q} \leq \dot{q}^{\max}$$

We know current is proportional to torque:

$$i = k_T \tau$$

---

Assuming  $\dot{q}$  remains constant during time step  $\Delta t$ :

$$\dot{q}(t + \Delta t) = \dot{q}(t) + \Delta t \ddot{q}$$

## JOINT POSITION LIMITS

Could use same approach as for joint vel.:

$$q(t + \Delta t) = q(t) + \Delta t \dot{q}(t) + \frac{\Delta t^2}{2} \ddot{q}$$

$$q^{\min} \leq q + \Delta t \dot{q} + \frac{\Delta t^2}{2} \ddot{q} \leq q^{\max}$$

0 → 0 → 0 → 0 → 0



This could result in large accelerations, incompatible with torque/current limits  $\Rightarrow$  QP unfeasible

Heuristics: use larger values of  $\Delta t$  (e.g. [10, 100] ms)  
improve behavior, but still no guarantees.

# TAXONOMY OF CONVEX OPTIMIZATION PROBLEMS

- Least-Squares Program (LSP)
  - Linear eq./ineq. constraints ( $Ax \leq b$ ,  $Ax = b$ )
  - 2-norm of affine function  $\|Ax - b\|^2$
- LSP are subclass of Quadratic Programs (QP):
  - Linear eq./ineq. constraints
  - Convex quadratic cost:  $\frac{1}{2}x^T H x + g^T x$ ,  $H \overset{\text{Positive Semi Definite}}{\underset{\text{eig}(H) \geq 0}{\circlearrowleft}}$

LSP and QP are convex optim. problems and can be solved efficiently with off-the-shelf software.

$\nwarrow$   
 $\text{<1ms}$

# LEAST-SQUARES PROGRAMS

$$\min_x \frac{1}{2} \|Ax - b\|^2 = \frac{1}{2} \underbrace{x^T A^T A x}_{H \succ 0} - b^T A x + \frac{1}{2} b^T b \triangleq f(x)$$

$$\nabla f = A^T A x - A^T b = \emptyset \Rightarrow \text{if } A^T A \text{ is invertible}$$

$$x^* = \underbrace{(A^T A)^{-1}}_{\text{(left) pseudo-inverse of } A} A^T b$$

$$(A^T A)^{-1} A^T A = I \quad \text{(left) pseudo-inverse of } A$$

If  $A^T A$  not invertible, but  $AA^T$  is invertible:

$$x^* = \underbrace{A^T (AA^T)^{-1}}_{\text{(right) pseudo-inverse}} b \Leftarrow \begin{cases} Ax^* = AA^T (AA^T)^{-1} b = b \\ A^T A A^T (AA^T)^{-1} b - A^T b = A^T b - A^T b = \emptyset \end{cases} \Rightarrow \begin{cases} f(x) = 0 \\ \nabla f = 0 \end{cases}$$

## (MOORE-PENROSE) PSEUDO-INVVERSE

$$x^* = \arg \min_x \|Ax - b\|^2 \Leftrightarrow x^* = A^+ b$$

If  $A$  is full row rank  $\Rightarrow A^+ = A^T (AA^T)^{-1}$

If  $A$  is full column rank  $\Rightarrow A^+ = (A^T A)^{-1} A^T$

If  $A$  is not full rank  $\Rightarrow A^+$  still exists, but equations above can't be used. In practice  $A^+$  is often computed using numerical methods, such as Singular Value Decomposition (SVD).

## FROM JOINT TO CARTESIAN SPACE

Given a reference trajectory  $\dot{x}^r(t)$  for the e-e, find  $\ddot{\tau}(t)$  such that  $x(t) \approx \dot{x}^r(t)$ . We know that:

$$\ddot{x} = J \ddot{q} + \dot{J} \dot{q} \quad \ddot{x}^d = \ddot{x}^r + K_d(\dot{x}^r - \dot{x}) + K_p(x^r - x)$$

We can compute desired joint acc.  $\ddot{q}^d$  as:

$$\begin{cases} \ddot{q}^d = J^+ (\ddot{x}^d - \dot{J} \dot{q}) \\ J^d = M \ddot{q}^d + h \end{cases}$$

$$J^+ = J^T (J J^T)^{-1}$$

Assuming  $J$  is full-row rank

$$(J^*, \ddot{q}^*) = \arg \min_{\ddot{q}, J^*} \| \underbrace{J \ddot{q} + \dot{J} \dot{q} - \dot{x}^d}_{\ddot{x}} \|^2 \equiv \| \ddot{x} - \ddot{x}^d \|^2$$

subject to  $M \ddot{q} + h = \tau$

## TASK MODELS

- Generalize concept of end-effector control
- Task = control objective
- Describe task with function  $e(\cdot)$  to minimize
- Assume that  $e(\cdot)$  describes error b/w real & reference traj.:

$$e(x, u, t) = \underbrace{y(x, u)}_{\text{real}} - \underbrace{y^r(t)}_{\text{reference}} \quad \begin{aligned} x &\triangleq (q, \dot{q}) \\ u &\triangleq \tau \end{aligned}$$

- We cannot directly minimize  $e(\cdot)$  if it depends on  $q$  and  $\dot{q}$ , which are not decision variables of the LSP.

I D E A: Impose dynamics of  $e$  (e.g.,  $\dot{e}$  or  $\ddot{e}$ ) such that:

i)  $\lim_{t \rightarrow \infty} e(t) = 0$

ii) Dynamics of  $e$  is **affine function** of  $(\ddot{q}, \ddot{\tau})$

## CASE #1: VELOCITY TASK FUNCTION

$$e(x, u, t) = e(\dot{q}, t) = y(\dot{q}) - \dot{y}^r(t)$$

For instance:  
[ $y(\dot{q})$  = angular momentum]

$$\dot{e} = \ddot{y} - \ddot{y}^r = \frac{\partial y}{\partial \dot{q}} \frac{d\dot{q}}{dt} - \ddot{y}^r = J \ddot{\dot{q}} - \ddot{y}^r$$

Choose  $\dot{e}$  such that  $e \rightarrow 0$ . Let's take a linear feedback:

$$\dot{e} = -K e \quad K > 0 \quad \Rightarrow \quad \lim_{t \rightarrow \infty} e(t) = 0$$

$$J \ddot{\dot{q}} - \ddot{y}^r = -K e$$

$$\underbrace{J \ddot{\dot{q}}}_A = \underbrace{\ddot{y}^r - K e}_b \quad \Rightarrow \quad A \ddot{q} - b = 0$$

## CASE #2: CONFIGURATION TASK FUNCTION

$$e(x, u, t) = e(q, t) = y(q) - y^r(t)$$

$$\dot{e} = \dot{y} - \dot{y}^r = \frac{\partial y}{\partial q} \frac{dq}{dt} - \dot{y}^r = \bar{J} \dot{q} - \dot{y}^r$$

$$\ddot{e} = \bar{J} \ddot{q} + \bar{J} \dot{q} - \ddot{y}^r$$

Let's impose linear dynamics:  $\ddot{e} = -K_p e - K_d \dot{e}$

$$K_p \succ 0$$

$$K_d \succ 0$$

$$\bar{J} \ddot{q} + \bar{J} \dot{q} - \ddot{y}^r = -K_p e - K_d \dot{e}$$

$$\bar{J} \ddot{q} = \underbrace{\ddot{y}^r - \bar{J} \dot{q} - K_p e - K_d \dot{e}}$$

$$A \ddot{q} = b \Rightarrow \text{minimize } \|A \ddot{q} - b\|^2$$

## CASE #3: CONTROL-INPUT TASK FUNCTION

Control inputs  $u \in \mathcal{S}$  are part of the decision variables.

Therefore  $e(x, u, t) = e(u, t)$  must be **affine**:

$$e(u, t) = A_u(t) u - b(t)$$

If it's not, I cannot use it in a Q.P.-based controller.

## TASK FUNCTIONS: SUMMARY

3 kinds of task functions:

- Affine function of  $\dot{q} \equiv \tau$
- Nonlinear function of  $x$ :
  - For functions of  $\dot{q} \Rightarrow$  impose first derivative  $\dot{\epsilon}$
  - For functions of  $q \Rightarrow$  impose 2nd derivative  $\ddot{\epsilon}$

In the end we have affine functions of  $(\dot{q}, \tau)$ :

$$g(z) = \underbrace{\begin{bmatrix} A_x & A_u \end{bmatrix}}_A \begin{bmatrix} \dot{q} \\ \tau \end{bmatrix} - b$$

## TASK-SPACE INVERSE DYNAMICS

$$\begin{array}{ll}\text{minimize}_{z=(\ddot{q}, \tau)} & \frac{1}{2} \|Az - b\|^2\end{array}$$

$$\text{subject to } [M \ -I] z = -h$$

• Joint-space control:  $A = [I \ 0], \quad b = \ddot{q}^d$

$$Az - b = \ddot{q} - \ddot{q}^d$$

• End-effector control:  $A = [J \ 0], \quad b = \dot{x}^d - J\dot{q}$

$$Az - b = J\ddot{q} + \dot{J}\dot{q} - \dot{x}^d$$

## UNDER-ACTUATED SYSTEMS

Example: remove one motor from manipulator. How does it affect dynamics?

$$M\ddot{q} + h = \tilde{\tau} = \begin{bmatrix} \tilde{\tau}_0 \\ \vdots \\ \tilde{\tau}_{n-1} \end{bmatrix} = \begin{bmatrix} \emptyset \\ \tilde{\tau}_1 \\ \vdots \\ \tilde{\tau}_{n-1} \end{bmatrix} = \begin{bmatrix} \emptyset \\ \tilde{\tau}_{1:n} \end{bmatrix}$$

Introduce  $S \triangleq [\emptyset \ I]$   $M\ddot{q} + h = S^T \tilde{\tau}_{1:n}$

Examples of underactuated systems:

- legged robots (biped, quadruped)
- flying robots (most of them)
- underwater robots (most of them)
- wheeled robots

# TSID FOR UNDER ACTUATED ROBOTS

$$\text{minimize}_{\begin{matrix} z \\ \ddot{q}, \gamma \end{matrix}} \frac{1}{2} \|Az - b\|^2$$

$$\text{subject to } [M \quad -\cancel{I}] z = -h \\ S^T$$

## TSID FOR RIGID CONTACT

Rigid contacts constrain the motion:

$\zeta(q) = \text{const} \Leftrightarrow \text{contact points don't move}$

$$J \triangleq \frac{\partial \zeta(q)}{\partial q}$$

$J \dot{q} = 0 \Leftrightarrow \text{contact point velocities are null}$

$\ddot{J} \ddot{q} + \dot{J} \dot{q} = 0 \Leftrightarrow \text{contact point accelerations are null}$

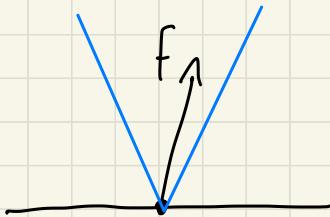
Introduce contact forces  $f$  as decision variables in QP:

$$\text{minimize } \|Az - b\|^2$$

$$z = (\ddot{q}, \tau, f)$$

subject to  $\begin{cases} M \ddot{q} + h = S^T \tau + J^T f \\ J \ddot{q} + \dot{J} \dot{q} = 0 \end{cases} \Leftrightarrow \begin{bmatrix} M & -J^T & -S^T \\ J & 0 & 0 \end{bmatrix} \begin{bmatrix} \ddot{q} \\ f \\ \tau \end{bmatrix} = \begin{bmatrix} -h \\ -\dot{J} \dot{q} \end{bmatrix}$

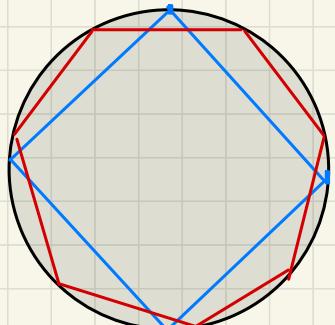
# FRIC TION CONE CONSTRAINTS



$$\| f_T \|^2 \leq \mu \| f_N \|^2$$

↑                      ↑  
Tangential            Normal  
Force                Force

$\mu$  = friction coefficient



Approximate quadratic cone with linear  
pyramidal cone:

$$B f \leq 0$$

Include this in TSID's QP.

## MULTI-TASK CONTROL

For complex redundant robots, typically one has multiple control objectives, e.g.  $N$  tasks, each defined by a task function

$$g_i(z) \triangleq \|A_i z - b_i\|^2 \quad i = 1 \dots N$$

TSID QP:

$$\begin{aligned} & \text{minimize}_{z=(\ddot{q}, \dot{r}, f)} \quad \sum_{i=1}^N w_i g_i(z) = \|A z - b\|^2 \\ & \text{subject to} \quad \begin{bmatrix} M & -J^T & -S^T \\ J & 0 & 0 \end{bmatrix} \begin{bmatrix} \ddot{q} \\ f \\ \dot{r} \end{bmatrix} = \begin{bmatrix} -h \\ -\dot{J}\dot{q} \end{bmatrix} \end{aligned}$$

$w_i \in \mathbb{R}^+$   $\triangleq$  user-defined weight

Large  $w_i \Rightarrow$  high importance to task  $i$

# OPTIMAL CONTROL

XVII century, Bernoulli (1696): ground shape to let ball go from A to B in minimum time? Brachistochrone curve

- How to accelerate car to reach goal in min time/energy with max speed

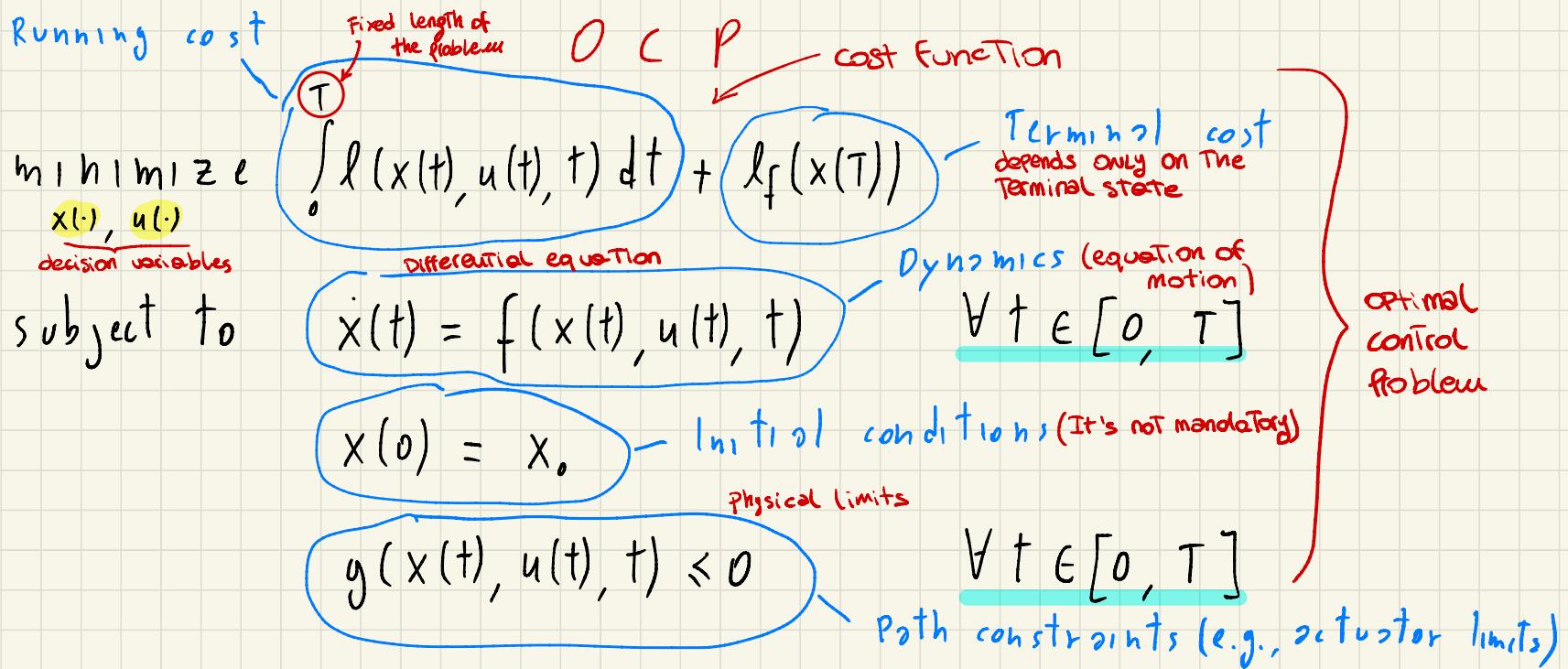
## CLASSIC CONTROL:

- need reference trajectory
- rise time (time to reach 90% of target)
- overshoot
- settling time (time to reach an equilibrium)
- hard to handle constraints



## O.C.:

- no need for ref. traj.
- objective clearly specified through cost function
- can handle constraints



- $x(\cdot), u(\cdot)$  are trajectories  $\Rightarrow$  infinite-dimensional objects
  - constraints are infinitely many! (you have infinite time instant between  $[0, T]$ )
  - $\neq$  optimization problem
- Main differences with optimization problem

# OCP EXAMPLE: PENDULUM

minimize  
 $x(\cdot), u(\cdot)$

subject to

terminal cost must be in function only  
of the terminal state

It's possible  
to use absolute  
value instead  
of square

with square you  
are more tolerant  
with small errors and  
penalize more  
big errors

$\int_0^T \left( u_1 q(t)^2 + u_2 \dot{q}(t)^2 + u_3 \ddot{q}(t)^2 \right) dt + 10^2 q(T)^2$

pendulum reach  
 $q(T)=0$  faster  
terminal condition

torque and acc. are linear dependent so we can optimize  
only one

Standard form

$I \ddot{q} = \tau + m g \sin(q)$

inertia      GRAVITY EFFECT      non-linear term

System dynamics  
of a pendulum

$\begin{bmatrix} \dot{q} \\ \ddot{q} \end{bmatrix} = \begin{bmatrix} \dot{q} \\ I^{-1}(\tau + m g \sin(q)) \end{bmatrix}$

$\forall t$

$f(x, u)$

$$|\tau(t)| \leq \tau^{\max}$$

$$\forall t \in [0, T]$$

$$q^{\min} \leq q(t) \leq q^{\max}$$

$$\forall t \in [0, T]$$

state of the system

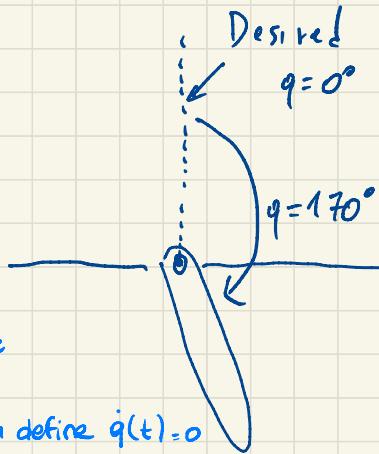
$x = (q, \dot{q})$

$u = \tau$  control input

$q(T) = \emptyset$   
model choice  
to set the final position  
to joint angle  $q=0$

If I don't specify a terminal constraint, I need to add a weight to the terminal cost. We choose to use one or the other.

↳ Same for joint velocities, we can define  $\dot{q}(t)=0$



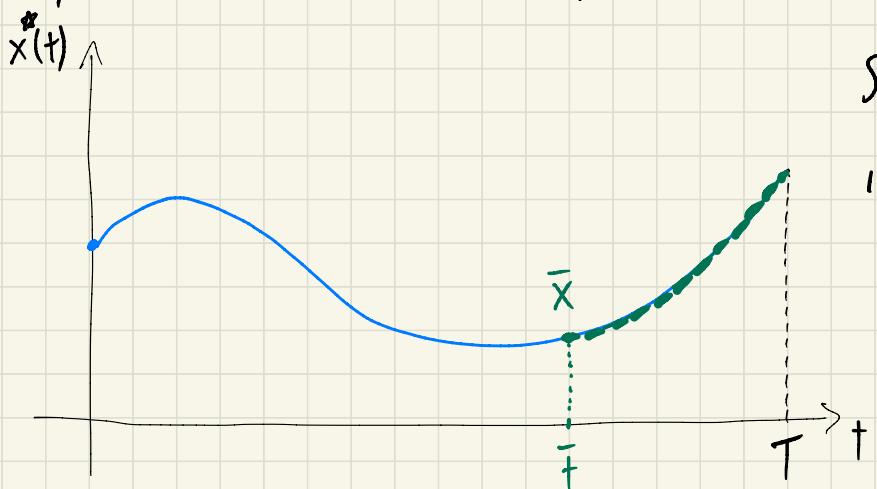
# O.C. METHOD FAMILIES

		OPTIMIZE $\Rightarrow$ DISCRETIZE	DISCRETIZE $\Rightarrow$ OPTIMIZE
		Continuous Time	Discrete Time
Global	(2)	Hamilton-Jacobi-Bellman (HJB)	Dynamic Programming (DP)
	(3)	<ul style="list-style-type: none"><li>Pontryagin Maximum Principle (PMP)</li><li>Calculus of variations</li><li>Indirect method</li></ul>	<ul style="list-style-type: none"><li>Direct method</li></ul>
Local			

NOTEBOOK numerical\_optim

# ① DP; PRINCIPLE OF OPTIMALITY (BELLMAN)

Any subarc of an optimal trajectory is optimal



Subarc on  $[\bar{t}, T]$  is optimal for initial value  $\bar{x}$ .

# DP: DISCRETE TIME OCP

$$\underset{x, u}{\text{minimize}} \quad \sum_{i=0}^{N-1} l(x_i, u_i) \quad \begin{array}{l} \text{[No terminal cost for sake} \\ \text{of simplicity]} \end{array} \quad \left| \begin{array}{l} X = (x_1, \dots, x_N) \\ U = (u_0, \dots, u_{N-1}) \end{array} \right.$$

subject to  $x_{i+1} = f(x_i, u_i) \quad i = 0 \dots N-1$

---

Take  $M$  s.t.  $0 \leq M < N$ . Then split problem cost:

$$\underset{x, u}{\text{minim.}} \quad \left[ \underbrace{\sum_{i=0}^{M-1} \left( l(x_i, u_i) + I(x_{i+1} - f(x_i, u_i)) \right)}_{\text{Indicator function}} + \underbrace{\sum_{i=M}^{N-1} \left( l(x_i, u_i) + I(x_{i+1} - f(x_i, u_i)) \right)} \right]$$

$$\underset{x, u}{\text{minimize}} \quad \left[ c_0(X_{1:M}, U_{0:M-1}) + c_M(x_M, X_{M+1:N}, U_{M:N-1}) \right]$$

$$\underset{x_{1:n}, u_{0:n-1}}{\text{minimize}} \quad \left[ c_0(X_{1:M}, U_{0:M-1}) + \underset{x_{M+1:N}, u_{M:N-1}}{\text{minimize}} \left[ c_M(x_M, X_{M+1:N}, U_{M:N-1}) \right] \right]$$

$x_M$  is treated as constant in this minimization!

# DP: DISCRETE TIME OCP

$$\underset{x_{1:n}, u_{0:n-1}}{\text{minimize}} \left[ c_0(x_{1:n}, u_{0:n-1}) + \underset{x_{M+1:N}, u_{n:n-1}}{\text{minimize}} \left[ c_M(x_M, x_{M+1:N}, u_{M:n-1}) \right] \right]$$

$V_M(x_M)$

Optimal cost of inner minimization is function of its initial state  $x_M$

$$V_0(x_0) = \underset{x_{1:n}, u_{0:n-1}}{\text{minimize}} \quad c_0(x_{1:n}, u_{0:n-1}) + V_M(x_M) \quad \text{— VALUE FUNCTION or OPTIMAL COST-TO-GO}$$

$V_M(x_M)$  = cost paid starting in  $x_M$  at time step M and behaving optimally

$$V_i(x_i) = \underset{u_i}{\text{minimize}} \left[ l(x_i, u_i) + V_{i+1} \left( \underbrace{f(x_i, u_i)}_{x_{i+1}} \right) \right]$$

RECURSIVE FORMULATION  
M=1

# DP : ALGORITHM

Given discrete-time finite-horizon DCP

$$V_N(x_N) = l_f(x_N)$$

Use recursive optimality principle (backward in time):

$$V_i(z) = \underset{u}{\text{minimize}} \quad l(z, u) + V_{i+1}(f(z, u)) \quad i = N-1, \dots, 0$$

Once you have  $V_i(\cdot) \forall i \in [0, N]$  compute optimal control as:

$$u_i^*(x) = \arg \underset{u}{\min} \quad l(x, u) + V_{i+1}(f(x, u))$$

OPTIMAL FEEDBACK

Q function

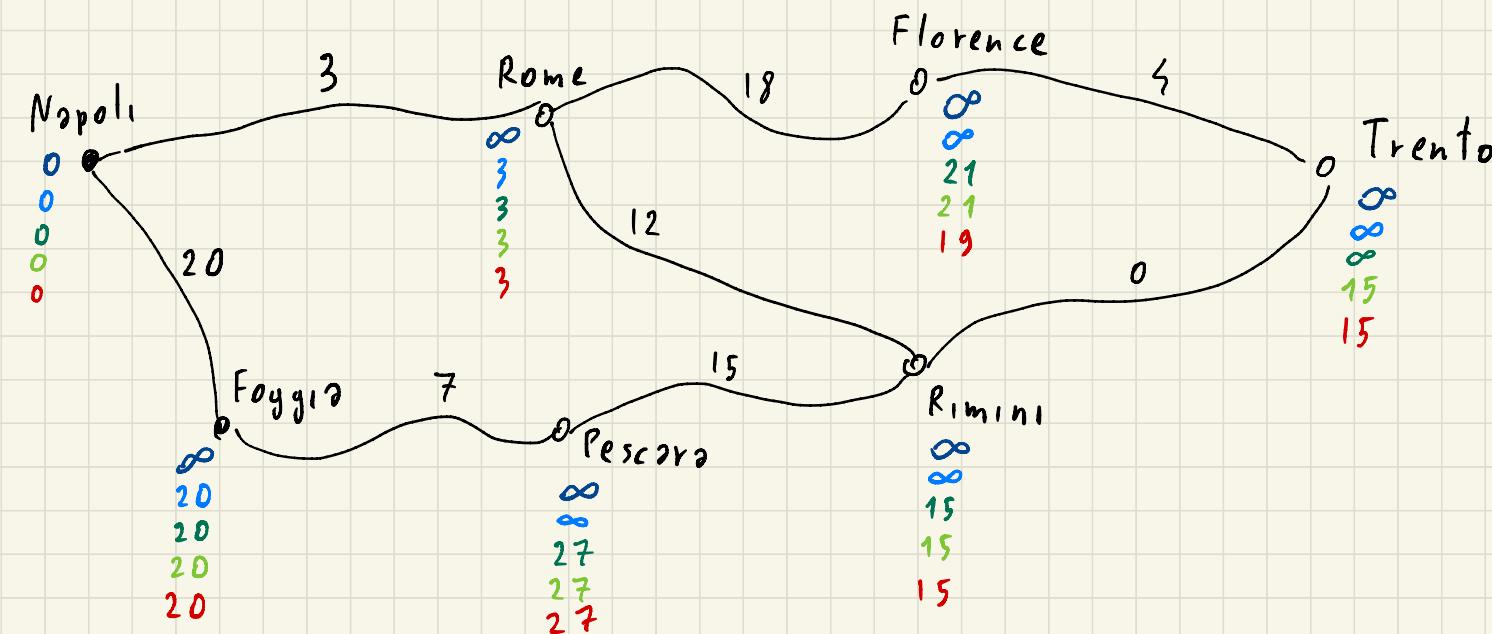
CONTROL POLICY

PARAMETRIC OPTIMIZ.  
 $\neq$   
NUMERICAL OPTIMIZAT.  
Solution is a function  
not a value

# DP EXAMPLE: ROAD TRIP

$V(K, \text{city}) \triangleq$  optimal cost-to-go from "city" on day  $K$  (10 days max)

$V(10, \cdot) \rightarrow V(9, \cdot) \rightarrow V(8, \cdot) \rightarrow V(7, \cdot) \rightarrow V(6, \cdot)$



## DP: DISCUSSION

- Global solution for discrete-time OCP
- Feedback solution (closed-loop)
- Applicable only in specific cases
  - ① discrete state-ctrl space  $\Rightarrow$  parametric optimization replaced by numerical optimization applied to all possible state values, results stored in lookup table
    - ⊖ Curse of dimensionality  $\Rightarrow$  applicable up to 2-4 states-controls
    - ⊕ Often applied to games (chess) or small continuous systems
    - ⊕ Easy to handle integer variables
  - ② Linear dynamics & quadratic cost  $\Rightarrow$  LQR

# O.C. METHOD FAMILIES

	Continuous Time	Discrete Time
Global	Hamilton-Jacobi-Bellman (HJB)	Dynamic Programming (DP)
Local	<ul style="list-style-type: none"><li>Pontryagin Maximum Principle (PMP)</li><li>Calculus of variations</li><li>Indirect method</li></ul>	<ul style="list-style-type: none"><li>Direct method</li></ul>

FROM ① D P  $\Rightarrow$  TO ② H J B

$$V(t_i, x) = \min_u l_d(x, u) + V(t_{i+1}, f_d(x, u))$$

$$l_d(x, u) = l(x, u) \delta$$

For small time step  $\delta \rightarrow 0$  we can approximate  $V$  with Taylor expansion:

$$\cancel{V(t_i, x)} = \min_u l_d(x, u) + \cancel{V(t_i, x)} + \frac{\partial V}{\partial t} \delta + \frac{\partial V}{\partial x} \frac{dx}{dt} \delta$$

$$0 = \min_u l(x, u) \delta + \dot{V} \delta + \nabla_x V f(x, u) \delta$$

$$\dot{V} \triangleq \frac{\partial V}{\partial t}$$

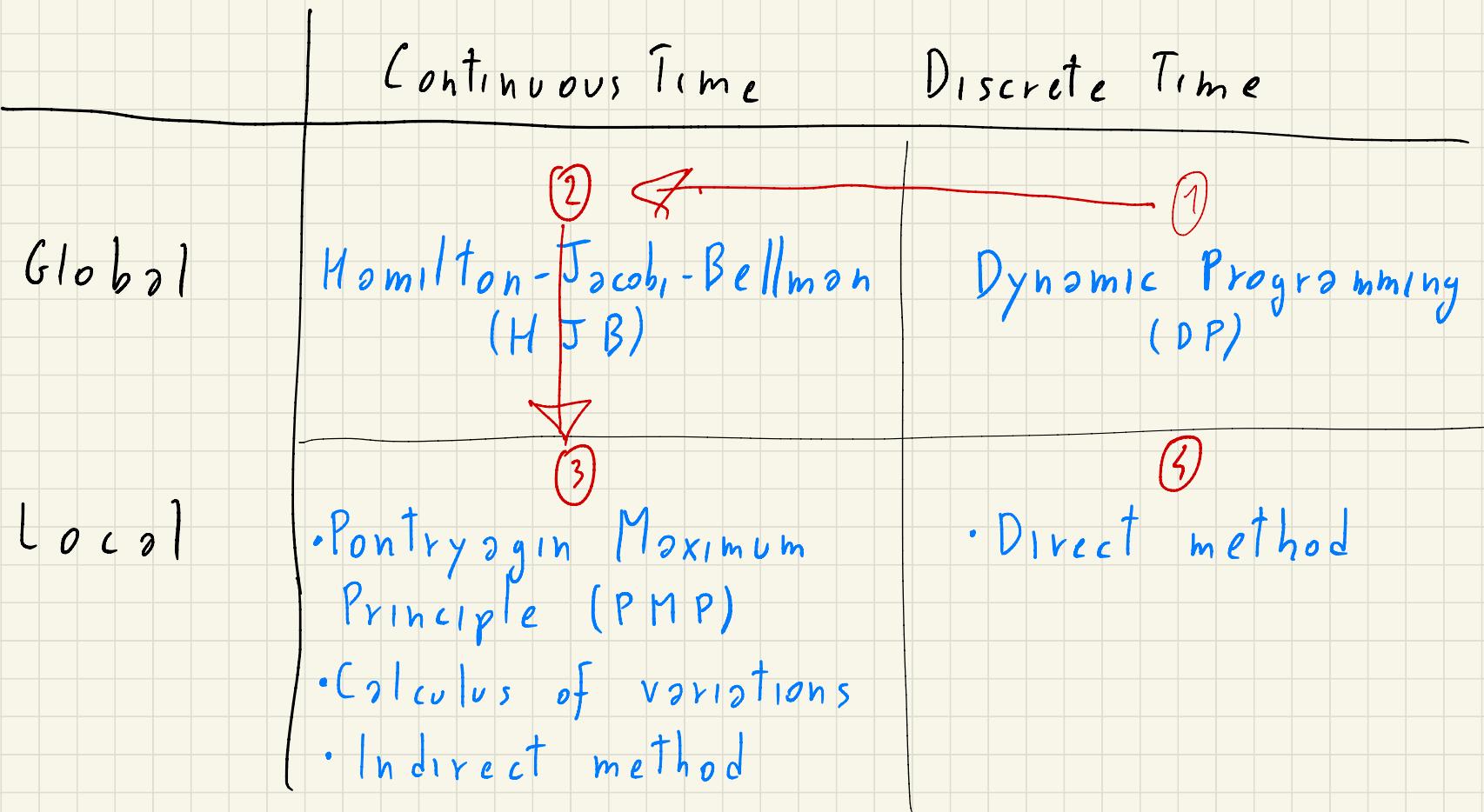
$$-\dot{V} = \min_u l(x, u) + \nabla_x V f(x, u)$$

$$V(t_f, x) = l_f(x)$$

Partial Differential Equation  $\Rightarrow$  Hard to solve!

Analytic solution for linear dynamics and quadratic cost

# O.C. METHOD FAMILIES



PMP



$$\text{HJB: } -\dot{V} = \min_u l(x, u) + \nabla_x V f(x, u)$$



OBSERVATION:  $u^*$  depends only on  $\nabla_x V$ , not on  $V$

IDEA: Introduce adjoint variables:

$$\lambda(t) \triangleq \nabla_x V(t, x(t))^T$$

Then compute  $u^*$  as:

$$u^*(t, x, \lambda) = \arg \min_u (l(x, u) + \lambda^T f(x, u)) \quad \text{Hamiltonian } H(x, u, \lambda)$$

How do we compute  $\lambda(t)$ ? Differentiate HJB w.r.t.  $x$ :

$$-\nabla_x \dot{V} = \nabla_x \left[ \min_u H(x, u, \nabla_x V) \right]$$

$$-\dot{\lambda} = \nabla_x H(x, u^*, \lambda) = \nabla_x l(x, u^*) + \lambda^T \nabla_x f(x, u^*)$$

Likewise, differentiate  $V(t_f, x) = l_f(x)$  to get:

$$\nabla_x V(t_f, x) = \nabla_x l_f$$

$$\lambda(t_f) = \nabla_x l_f \quad \leftarrow \text{TERMINAL CONDITIONS}$$

To summarize optimality conditions as boundary value problem:

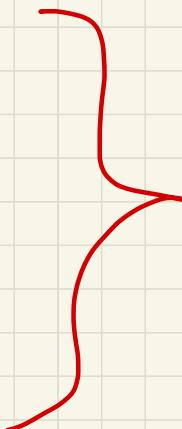
$$\dot{x} = f(x, u)$$

$$-\dot{\lambda} = \nabla_x H(x, u^*, \lambda)$$

$$u^* = \arg \min_u H(x, u, \lambda)$$

$$x(0) = x_0$$

$$\lambda(t_f) = \nabla_x l_f(x(t_f))$$



BVP

HOW DO WE SOLVE  
THIS?

## SOLVING B. V. P.

- Assume we know  $\lambda(0)$  (even if we know  $\lambda(t_f) = \nabla_x l_f(x(t_f))$ )
- We do know  $x(0) = x_0$
- LOOP:
  - Compute  $u^* = \arg \min_u H(x, u, \lambda)$
  - Compute  $\dot{\lambda} = -\nabla_x H(x, u^*, \lambda)$
  - Compute  $\dot{x} = f(x, u^*)$
  - Integrate  $\dot{x}, \dot{\lambda}$  to get new values of  $x, \lambda$ :

$$x := x + \Delta t \dot{x}$$

$$\lambda := \lambda + \Delta t \dot{\lambda}$$

But in reality we don't know  $\lambda(0)$ !

# SOLVING B.V.P.

- Guess  $\lambda(0) = \lambda_0$ .
- Integrate forward to get  $\lambda(t_f)$
- If  $\lambda(t_f) = \nabla_x \ell_f(x(t_f)) \Rightarrow \text{DONE!}$
- Otherwise, we can see  $\lambda(t_f)$  as a function of  $\lambda_0$  and solve these nonlinear equations with Newton:

$$\underbrace{\lambda(t_f, \lambda_0) - \nabla_x \ell_f(x(t_f))}_{} = 0$$

$$r(\lambda_0) = 0$$

$$r(\lambda_0) + \nabla r \Delta \lambda_0 = 0 \Rightarrow \Delta \lambda_0 = -\nabla r^{-1} r(\lambda_0)$$

How to compute this?

Newton step

# O.D.E. SENSITIVITIES

$\nabla r = \frac{\partial \lambda(t_f, \lambda_0)}{\partial \lambda_0}$  captures how  $\lambda(t_f)$  changes for small variation of  $\lambda_0$

Many techniques to compute  $\nabla r$ :

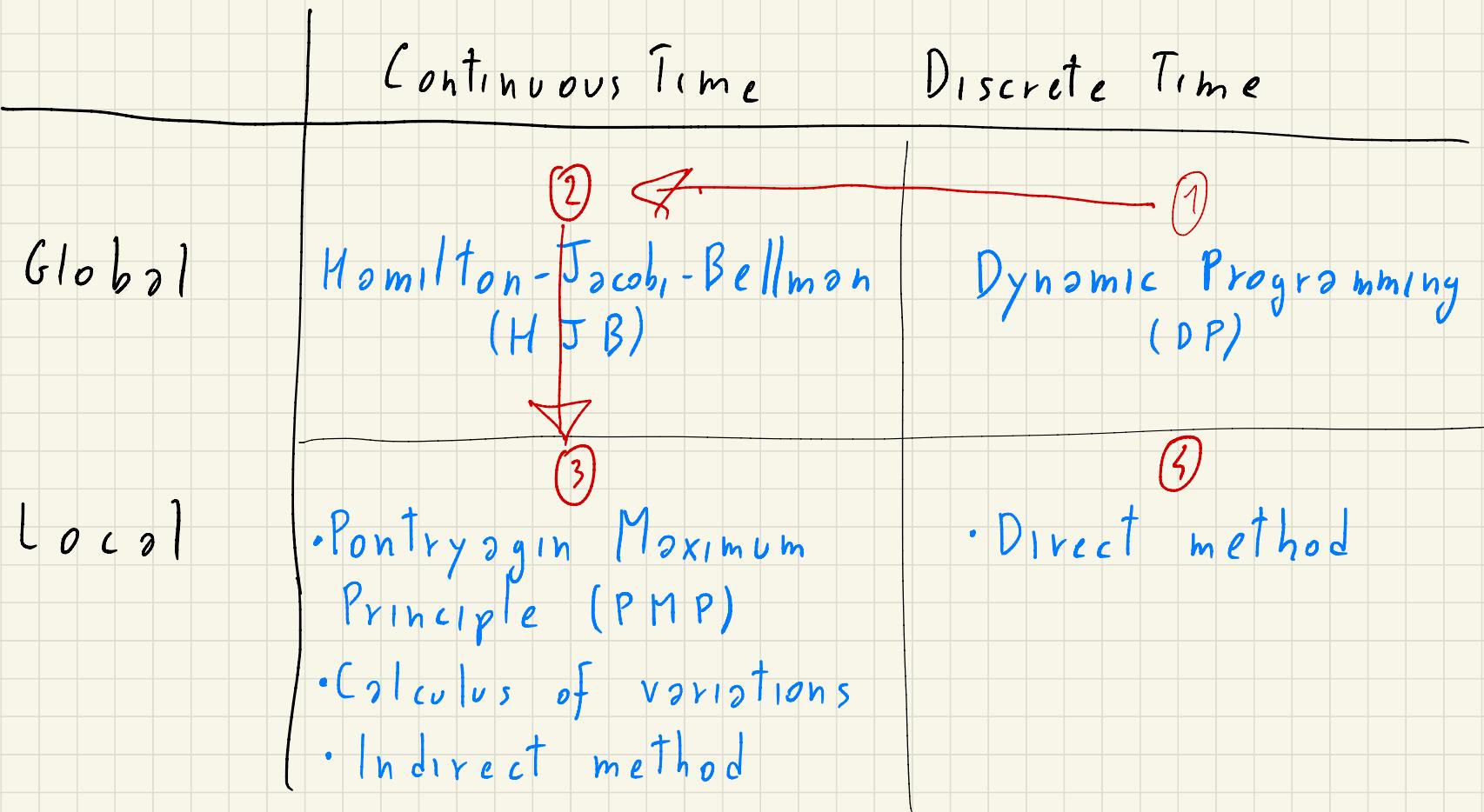
- External numerical differentiation
- Variational Differential Equations
- Automatic Differentiation
- Internal Numerical Differentiation

- perturb  $\lambda_0 \Rightarrow \bar{\lambda}_0 = \lambda_0 + \epsilon e$ :
- call integrator with  $\lambda_0 = \bar{\lambda}_0$
- compute derivatives by finite differences:

$$\frac{\lambda(t_f, \bar{\lambda}_0) - \lambda(t_f, \lambda_0)}{\epsilon}$$

- ⊕ Easy to implement
- ⊖ Computationally expensive
- ⊖ Poor accuracy

# O.C. METHOD FAMILIES



# NONLINEAR PROGRAMMING

$$\min_x f(x)$$

$$\text{s.t. } g(x) = 0 \\ h(x) \leq 0$$

Feasible set

$$\Omega = \{x \mid g(x) = 0, h(x) \leq 0\}$$

Assume  $f, g, h$  have continuous 1<sup>st</sup> and 2<sup>nd</sup> derivatives.

Let's derive the local optimality conditions.

## LOCAL OPTIMALITY

$x^*$  local minimizer  $\Leftrightarrow \exists \delta > 0$  s.t.  $f(x^*) \leq f(x) \quad \forall x \in B_\delta(x^*) \cap \Omega$

Strict local optimality implies:

$$f(x^* + d) - f(x^*) > 0 \quad \forall d, \exists \delta > 0 \text{ s.t. } x^* + d \in \Omega$$

$$d^\top \nabla f(x^*) + \frac{1}{2} d^\top \nabla^2 f(x^*) d + O(d^3) > 0$$

This can hold only if:

$$\nabla f(x^*) = 0$$

$$\nabla^2 f(x^*) \succ 0$$

UNCONSTRAINED CASE

No direction  $d$  can locally improve the cost

## UNCONSTRAINED CASE

1<sup>st</sup>-order Necessary Conditions of Optimality

$$x^* \text{ local optimum} \Rightarrow \nabla f(x^*) = \emptyset$$

2<sup>nd</sup>-order Sufficient Conditions of Optimality

$$\nabla f(x^*) = \emptyset \quad \text{and} \quad \nabla^2 f(x^*) \succ \emptyset \Rightarrow x^* \text{ strict local minimizer}$$

$$\nabla f(x^*) = \emptyset \quad \text{and} \quad \nabla^2 f(x^*) \prec \emptyset \Rightarrow x^* \text{ strict local maximizer}$$

$\nabla^2 f(x^*)$  indefinite  $\Rightarrow$  Nothing can be said

## EQUALITY CONSTRAINTS

Feasible cone:  $\mathcal{F}(x) = \{ d \mid \nabla g(x)^T d = 0 \}$  Null Space of linear approximation of equality constraints

Set of descent directions:  $\mathcal{D}(x) = \{ d \mid d^T \nabla f(x) < 0 \}$  objective function decreases

Local minimum  $x \Rightarrow \mathcal{D}(x) \cap \mathcal{F}(x) = \emptyset$

Karush-Kuhn-Tucker (KKT) conditions:

$$\begin{cases} g(x) = 0 \\ d^T \nabla f(x) \geq 0 \quad \forall d \in \mathcal{F}(x) \end{cases} \Leftrightarrow d \in \text{null}(\nabla g(x)^T)$$

Not practical to check every  $d \in \mathcal{F}(x)$

$\nabla g(x)^T = J_{g(x)} \cdot \frac{\partial}{\partial x_i}$

Multidimensional Function

More practice to check

$$\Rightarrow \nabla f(x) \in \text{range}(\nabla g(x)) \Leftrightarrow \nabla f(x) = \nabla g(x) \lambda$$

Define Lagrangian:  $\mathcal{L}(x, \lambda) = f(x) + \lambda^T g(x)$

KKT conditions can be written as:

$$\nabla \mathcal{L}(x, \lambda) = 0 \iff \begin{cases} \nabla f(x) + \nabla g(x) \lambda = 0 & \rightarrow \text{Derivation wrt } x \\ g(x) = 0 & \rightarrow \text{Derivation wrt } \lambda \end{cases}$$

# INEQUALITY CONSTRAINTS

Inequality constraints are important only at the boundaries

We can move in any direction if it is inactive. Locally you can move

Define set of active constraints  $A(x) = \{i \mid h_i(x) = 0\}$

Feasible cone  $\tilde{F}(x) = \{d \mid \nabla g(x)^T d = 0, \nabla h_i(x)^T d \leq 0 \text{ } i \in A(x)\}$

## KKT CONDITIONS:

$$\begin{cases} g(x) = 0, \quad h(x) \leq 0 \\ d^T \nabla f(x) \geq 0 \quad \forall d \in \tilde{F}(x) \end{cases}$$

For every direction in feasible cone, my cost function must increase? or stay the same

$$\Leftrightarrow \nabla f(x) + \underbrace{\nabla g(x)}_{\text{Jacobian of equality constr.}} \lambda + \underbrace{\nabla h_{A(x)}}_{\text{Jacobian of active inequality constraints}} \mu_A = 0$$

?

considering only active constraints, by moving in direction  $d$ , <sup>active</sup>constraints stay constant or decrease

$$\Leftrightarrow \begin{cases} \nabla f(x) + \nabla g(x)\lambda + \nabla h(x)\underbrace{\mu}_{\text{One multiplier for every constraint}} = 0 \\ \mu \geq 0 \\ \underbrace{\mu^T h(x) = 0}_{\text{COMPLEMENTARITY CONDITION}} \end{cases}$$

Either  $\mu_i = 0$  or  $h_i(x) = 0$  or both  $\Rightarrow$  if a in. constr. is active, then  $\mu_i$  can be positive, if an in. constr. is inactive, the  $\mu_i$  must be zero.

↳

if  $h(x) = 0 \Rightarrow \mu_i$  can be  $> 0$   
 if  $h(x) < 0 \Rightarrow \mu_i = 0$

# SOLVING KKT CONDITIONS

$$\begin{array}{ll} \min_x & f(x) \\ \text{s.t.} & g(x) = 0 \\ & h(x) \leq 0 \end{array}$$

case of eq. constraints

$\Rightarrow$

KKT system

$$\left\{ \begin{array}{l} \nabla_x \mathcal{L}(x, \lambda, \mu) = 0 \\ g(x) = 0 \\ h(x) \leq 0 \\ \lambda \geq 0 \\ \mu^T h(x) = 0 \end{array} \right.$$

extra conditions for ineq. constraints

$$\mathcal{L}(x, \lambda, \mu) = f(x) + \lambda^T g(x) + \mu^T h(x)$$

Let's focus on the case w/o inequalities.

## NEWTON METHOD

$$r(x) = \emptyset \quad \text{Current guess } \bar{x}$$

$$r(\bar{x} + \Delta x) \simeq r(\bar{x}) + \nabla r(\bar{x})^T \Delta x = \emptyset$$

$$\Delta x = -\nabla r(\bar{x})^{-T} r(\bar{x}) \quad \Leftarrow \text{NEWTON STEP}$$

Apply Newton step: INVERTIBLE?

$$x^{(k+1)} = x^{(k)} + \Delta x$$

Repeat until  $\|r(x)\| < \text{threshold}$

# SOLVE KKT WITH NEWTON

$$\begin{cases} \nabla f + \nabla g \lambda = \emptyset \\ g(x) = 0 \end{cases} \Leftrightarrow r(x, \lambda) = \begin{bmatrix} \nabla_x \mathcal{L}(x, \lambda) \\ g(x) \end{bmatrix} = \emptyset$$

$$r(\bar{x} + \Delta x, \bar{\lambda} + \Delta \lambda) \approx r(\bar{x}, \bar{\lambda}) + \nabla r^T \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \emptyset$$

$$\begin{bmatrix} \nabla_{xx}^2 \mathcal{L}(\bar{x}, \bar{\lambda}) & \nabla g(\bar{x}) \\ \nabla g(\bar{x})^T & 0 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta \lambda \end{bmatrix} = \begin{bmatrix} -\nabla_x \mathcal{L}(\bar{x}, \bar{\lambda}) \\ -g(\bar{x}) \end{bmatrix}$$

KKT matrix

Symmetric indefinite

## DESCENT DIRECTION & REGULARIZATION

Consider unconstrained problem:  $\min_x f(x)$

Direction  $d$  is descent direction  $\Leftrightarrow \nabla f(x)^T d < 0$

Newton step is:  $\Delta x = - H^{-1} \nabla f(x)$        $H \triangleq \nabla^2 f(x)$

Is  $\Delta x$  a descent direction?      ↑  
Hessian

$$-\nabla f(x)^T H^{-1} \nabla f(x) < 0 ?$$

Yes if  $H \succ 0$ . If  $\nabla^2 f(x)$  is not P.D.  $\Rightarrow$  regularize:

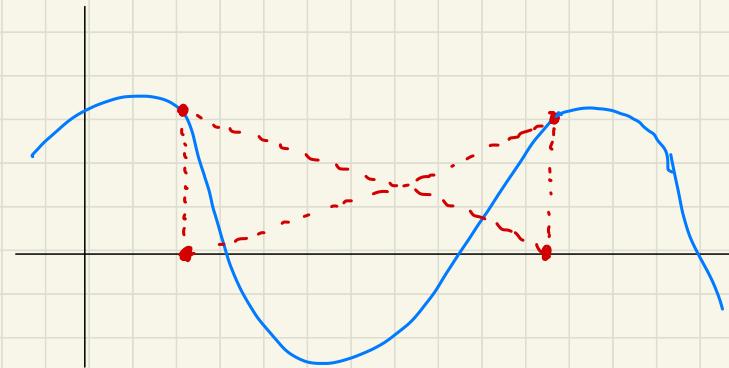
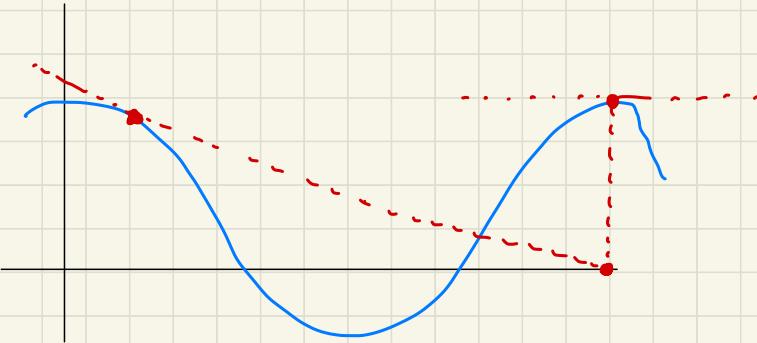
$$H = \nabla^2 f(x) + \lambda I$$

Choose  $\lambda > 0$  so that  $H \succ 0$

## ISSUES WITH VANILLA NEWTON

① You can get stuck at points of zero slope

② You can stuck in cycles or diverge



## BACKTRACKING LINE SEARCH

Step size  $\epsilon \in [0, 1]$

① Consider taking a fraction of the Newton step:  $x^{(i+1)} = x^{(i)} + \alpha \Delta x$

② Observation: if  $f(x)$  were linear  $\underbrace{f^{(i)} - f^{(i+1)}}_{\text{reduction of cost}} = \underbrace{\alpha \nabla f_i^T \Delta x}_{\substack{\text{expected from} \\ \text{linear approximation}}}$

③ In practice  $f(x)$  is nonlinear  $\Rightarrow$  we'd be happy if cost is reduced by a fraction  $\gamma$  of what linearization predicts

GOALS

$$f^{(i)} - f^{(i+1)} \geq \gamma \alpha \nabla f_i^T \Delta x$$

ARMITJO'S

CONDITION

Typically  $\gamma$  is very small, e.g.,  $\gamma < 10^{-2}$

# PSEUDO-CODE: NEWTON STEP - UNCONSTRAINED

- INPUTS:  $x^{(i)}$
- OUTPUT:  $x^{(i+1)}$

$$\textcircled{1} \quad f = f(x^{(i)})$$

$$\textcircled{2} \quad \Delta x = -(\nabla r^T + \lambda I)^{-1} r$$

$$\textcircled{3} \quad x_{\text{try}} = x^{(i)} + d \Delta x, \quad d = 1$$

$$\text{reduction} = f - f(x_{\text{try}})$$

while reduction <  $\gamma d \nabla f^T \Delta x$

$$\boxed{d = \beta d}$$

$$x_{\text{try}} = x^{(i)} + d \Delta x$$

$$\boxed{\text{reduction} = f - f(x_{\text{try}})}$$

$$x^{(i+1)} = x_{\text{try}}$$

Hyperparameters:

$$\beta \in [0, 1]$$

$$\gamma \in [0, 1]$$

$$\lambda > 0$$

## LINE SEARCH WITH EQUALITY CONSTRAINTS

Two requirements on  $\Delta x$ : i) reduce  $f(x)$   
ii) reduce  $\|g(x)\|$

Introduce  $\ell_1$  merit function:

$$M_1(x) = f(x) + \sigma \|g(x)\|_1$$

Exact merit function if:

$x^*, \lambda^*$  satisfy KKT conditions  $\Leftrightarrow x^*, \lambda^*$  is local min. of  $M_1(x)$

$M_1(x)$  is exact iff  $\sigma > \|\lambda\|_\infty$

$\sigma$  adapted at each iteration

## INEQUALITY CONSTRAINTS

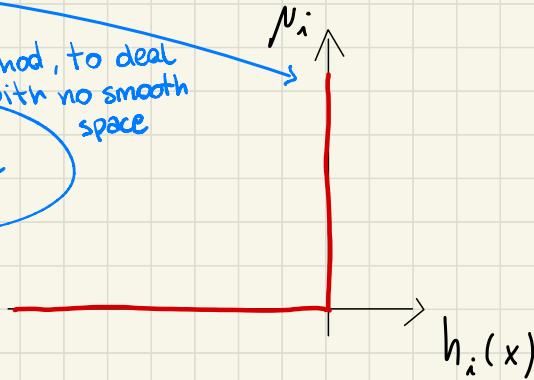
Newton method cannot be directly used because solution space generated by complementarity constraints is not smooth:

$$\mu^T h(x) = 0 \Rightarrow \mu^T h(x) \leq \epsilon$$

Interior point method, to deal with no smooth space

$$h(x) \leq 0$$

$$\mu \geq 0$$



Two main approaches exist:

- Active set → Transform active constr. into equality constr. and remove inactive constr.
- Interior point

## DIRECT METHODS

KEY IDEA: Discretize OCP  $\Rightarrow$  Use NLP solver to find local optimum

The finer the discretization  $\Rightarrow$  better approximation  $(+)\ominus$

Discretization:

- ① Parametrize state/control trajectories (e.g. with polynomials)
- ② Enforce constraints (dynamics + path inequalities) on a time grid  $(+)\oplus$

$$t_0 < t_1 < t_2 < \dots < t_N$$

Tons of theory on how to choose traj. parametrization and time grid so that NLP is:

- i) good approximation of original OCP
- ii) well conditioned  $\rightarrow$  important to solve linear eq.  $\rightarrow$  Newton method rely on linear eq.

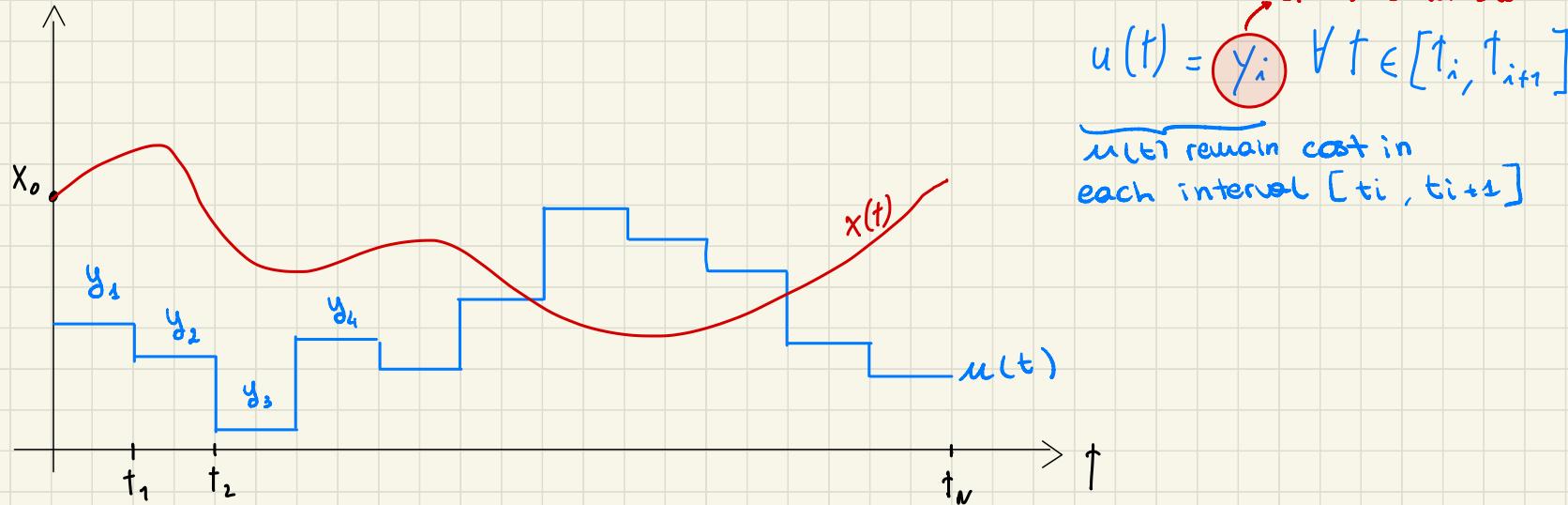
## DIRECT METHODS; FAMILIES

trajectory

- ① SINGLE SHOOTING;  
(SEQUENTIAL)
  - Discretize only  $u(t)$
  - Compute  $x(t)$  <sup>state trajectory</sup> integrating dynamics
  - No need to have dynamics constraints
- ② COLLOCATION;  
(SIMULTANEOUS)
  - Discretize both  $x(t)$  and  $u(t)$
  - Enforce dynamics on time grid
    - ↳ To ensure that  $x(t)$  and  $u(t)$  are coherent with dynamics
- ③ MULTIPLE SHOOTING
  - Discretize all  $u(t)$
  - Discretize  $x(t)$  at few points in time
  - Compute intermediate values of  $x$  by numerical integration

# ① SINGLE SHOOTING

Discretize  $u(t)$  on fixed grid  $\theta = t_0 < t_1 < \dots < t_N = t_f$



Compute  $x(t)$  from  $u(t)$  and  $x_0$  by integrating dynamics:

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

$$x(0) = x_0 \quad \text{Initial state is fixed and known}$$

# SINGLE SHOOTING - NLP

$$\underset{y}{\text{minimize}} \quad \int_0^{t_f} l(x(t; y), u(t; y), t) dt + l_F(x(t_N; y))$$

*Both state and control depends on y*

$$\text{subject to } g(x(t_i; y), u(t_i; y), t_i) \leq 0 \quad i = 0 \dots N$$

*Discretized path constraints*

Running cost integral can be computed when integrating dynamics:

$$c(t) = \int_0^t l(\cdot, \cdot) dt \Rightarrow \bar{x} = (x, c)$$

$$\begin{cases} \dot{x}(t) = f(x, u, t) \\ \dot{c}(t) = l(x, u) \\ c(0) = 0 \end{cases}$$

$$\dot{\bar{x}} = \begin{bmatrix} f(x, u, t) \\ l(x, u) \end{bmatrix}$$

# NUMERICAL INTEGRATION

Given ODE:  $\dot{x}(t) = f(x(t), t)$ ,  $x(0) = x_0$

Compute  $x(t) \forall t \in [0, T]$ .

If  $f$  is Lipschitz continuous  $\Rightarrow$  unique solution  $x$ .  
 ≈ bounded first derivatives

Neglect  $u$   
 because  $u$  depends  
 just on  $t$  and

## EXPLICIT EULER

$$\dot{x} = \lim_{h \rightarrow 0} \frac{x(t+h) - x(t)}{h}$$

Definition of derivative

Take  $h$  small

$\Rightarrow$  compromise between  
 computation and accuracy:

smaller  $h \Rightarrow$  ↑ comp. & accuracy  
 bigger  $h \Rightarrow$  ↓ comp. & accuracy

$$h f(x, t) \simeq x(t+h) - x(t) \Rightarrow x(t+h) = x(t) + h f(x, t)$$

Oh if  $h$  is sufficiently small  $\Rightarrow$  computationally expensive!  
 Not very accurate (order 1)

# RUNGE-KUTTA METHODS

$$x_{n+1} = x_n + h \sum_{i=1}^q b_i k_i$$

Approximate average  $\dot{x}$  with weighted average of  $k_i$ , with weights  $b_i$ , with  $\sum b_i = 1$ ,  $b_i \geq 0$

$$k_i = f(x_n + h \sum_{j=1}^q a_{ij} k_j, t_n + c_i h)$$

Approximate value of  $\dot{x}$  at time  $t_n + c_i h$ , with  $c_i \in [0, 1]$

- $q$  is the order of the method
- For  $q=1$  we recover Euler ( $b_1=1$ ,  $a_{11}=\emptyset$ ,  $c_1=\emptyset$ )
- If  $a_{ij} = \emptyset \quad \forall j \geq i \Rightarrow$  **Explicit** method

otherwise

$\Rightarrow$  **Implicit** method  $\Rightarrow$  need to solve system of nonlinear equations

more computationally demanding, but more stable

- For  $q > 1$  there exists many versions of same **order p**:

$$\lim_{h \rightarrow 0} [x(t) - \hat{x}(t; t-h, x(t-h))] = O(h^{p+1})$$

↑  
(consistency)

Exact trajectory Integrator output

# BUTCHER TABLEAU

$c_1$	$a_{11} \dots a_{1q}$
$\vdots$	$\ddots \ddots \vdots$
$\vdots$	$\ddots \ddots \vdots$
$c_q$	$a_{q1} \dots a_{qq}$
	<hr/>
	$b_1 \dots b_q$

- Uniquely defines an RK method
- Most common RK method is:

$$k_1 = f(x_n, t_n)$$

EXAMPLE OF METHOD  
FOR ORDER 4

0	0	0	0	
$\frac{1}{2}$	$\frac{1}{2}$	0	0	
$\frac{1}{2}$	0	$\frac{1}{2}$	0	
1	0	0	1	
	<hr/>	<hr/>	<hr/>	
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

order 4  $\rightarrow$

IF they are  
non-zero you  
have an  
implicit  
method

$$k_2 = f(x_n + \frac{1}{2}h k_1, t_n + \frac{h}{2})$$

$$k_3 = f(x_n + \frac{1}{2}h k_2, t_n + \frac{h}{2})$$

$$k_4 = f(x_n + h k_3, t_n + h)$$

$$x_{n+1} = x_n + \frac{h}{6} (k_1 + 2k_2 + 2k_3 + k_4)$$

For  $p \geq 5$   $\exists$  method with  $p = q$ !

→ with order 5 we  
have to do 6 eval.  
of dynamics and so on

Evaluate dynamics  
4 times for every  
step of integration  
↓

very accurate but very  
comp. demanding

## SUMMARY - SINGLE SHOOTING

- Remove  $x$  from decision variables
  - Discretize  $u(t) = y_i \quad \forall t \in [t_i, t_{i+1}]$
  - Compute  $x$  by integrating dynamics  $\dot{x} = f(x, u, t)$
  - Remove dynamics from constraints
- I remove  $x$  integrating
- High-order integration schemes (typically) more efficient
  - Differentiation of integration scheme often computed using automatic differentiation libraries.
  - Implicit schemes typically needed for stiff dynamics

## ② COLLOCATION

Discretize  $x(t)$  with polynomials (e.g. order  $k=0$ ) on fine grid:

$$x(t) = s_i \quad \forall t \in [t_i, t_{i+1}] \quad \begin{matrix} \text{One polynomial for each} \\ \text{time interval} \end{matrix}$$

Replace dynamics  $\dot{x} = f(x, u)$  with:

APPS. of dynamic

$$\frac{s_{i+1} - s_i}{t_{i+1} - t_i}$$

Dynamics

$$f(s_i, y_i) = 0$$

I approximate  $\dot{x}$

$$i = 0 \dots N-1$$

APPOX. OF dynamics must be equal to dynamics

$$c_i(s_i, s_{i+1}, y_i) = 0$$

Best approximation  
possible

Approximate cost integral:

$$\int_{t_i}^{t_{i+1}} l(x(t), u(t)) dt \approx l(s_i, y_i)(t_{i+1} - t_i) \triangleq l_i(s_i, y_i)$$

# COLLOCATION NLP (SPARSE)

$$\begin{matrix} \text{minimize}_s \\ s, y \end{matrix}$$

$$\sum_{i=0}^{N-1} l_i(s_i, y_i) + l_f(s_N)$$

OPT. problem

$$\text{subject to } s_0 - x_0 = 0$$

initial conditions

$$\text{discretized dynamics} \rightarrow c_i(s_i, s_{i+1}, y_i) = 0 \quad i = 0 \dots N-1$$

$$\begin{matrix} \text{discretized path} \\ \text{constraints} \end{matrix} \rightarrow g_i(s_i, y_i) \leq 0 \quad \begin{matrix} \text{only check that} \\ i = 0 \dots N-1 \end{matrix}$$

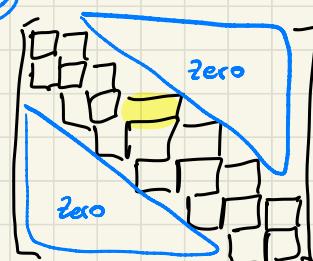
at discrete time instants

Why sparse? Because  $H H^T$  matrix is sparse!

$$H_{ij} = \frac{\partial}{\partial s_i} \left( \frac{\partial}{\partial s_j} L \right) = 0 \quad \forall i \notin \{j, j+1, j-1\}$$

computation became more efficient

Most of the entries of the matrix are zero



### ③ MULTIPLE SHOOTING

Discretize control  $u(t)$  on coarse grid  $t_0 < t_1 < \dots < t_n$

$$u(t) = y_i \quad \forall t \in [t_i, t_{i+1}]$$

Integrate numerically ODE in each interval  $[t_i, t_{i+1}]$   
state variable  $s_i$ :

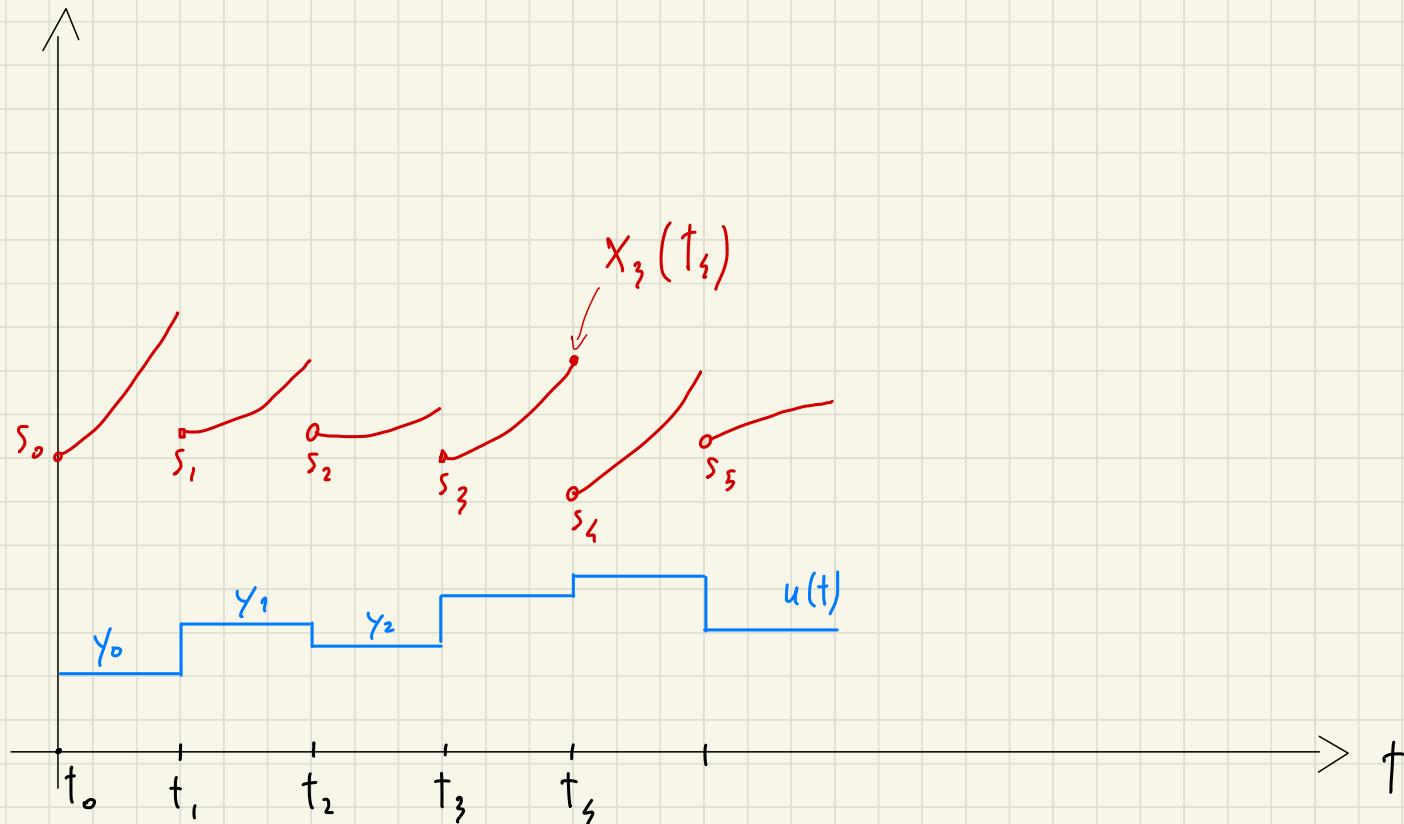
$$\dot{x}_i(t) = f(x_i(t), y_i) \quad t \in [t_i, t_{i+1}]$$

$$x_i(t_i) = s_i$$

Obtain  $x_i(t)$ . Numerically compute integrals:

$$l_i(s_i, y_i) = \int_{t_i}^{t_{i+1}} l(x_i(t), y_i) dt$$

# MULTIPLE SHOOTING SKETCH



# MULTIPLE SHOOTING NLP

$$\underset{s, y}{\text{minimize}} \quad \sum_{i=0}^{N-1} l_i(s_i, y_i) + l_f(x_N)$$

subject to  $s_0 - x_0 = 0$  initial conditions

$$s_{i+1} - x_i(t_{i+1}; s_i, y_i) = 0 \quad \text{continuity constr.}$$

$$g_i(s_i, y_i) \leq 0 \quad \text{path constraints}$$

Result of numerical  
integration from  $s_i$

# DIRECT METHODS - SUMMARY

## ① SINGLE SHOOTING

⊕ Small problem

⊕ Need initial guess  
only for  $x(t)$

⊕ Can use off-the-shelf  
ODE solvers

⊖ Cannot exploit knowledge  
of  $x(t)$  in initialization

⊖ Numerical issues for  
unstable systems

## ③ MULTIPLE SHOOTING

⊕ Medium problem

⊕ Sparser than ①

⊖ Less sparse than ②

⊕ Unstable OK

⊕ Initialize  $x(t)$

⊕ Can adapt time grid

## ② COLLOCATION

⊖ Large problem

⊖ Cannot adopt  
time grid

⊕ Sparse problem

⊕ Work well with  
unstable systems

⊕ Can initialize  $x(t)$

# MODEL PREDICTIVE CONTROL

- Use optimal control for controlling system (robot)
- Solve  $\rightarrow$  finite-horizon OCP using **current state** as initial state.

$$x^*, u^* = \underset{x, u}{\arg \min}$$

$$\sum_{i=0}^{N-1} l(x_i, u_i)$$

subject to  $x_{i+1} = f(x_i, u_i, k)$  limits on state-var.

$$i = 0 \dots N-1$$

$$x_{i+1} \in \mathcal{X}, \quad u_i \in \mathcal{U} \quad \begin{matrix} \text{limits on} \\ \text{control} \end{matrix} \quad i = 0 \dots N-1$$

$$x_0 = x^{\text{mess}}$$

- Apply computed control **for this time step:**  $\tau = u_0^*$

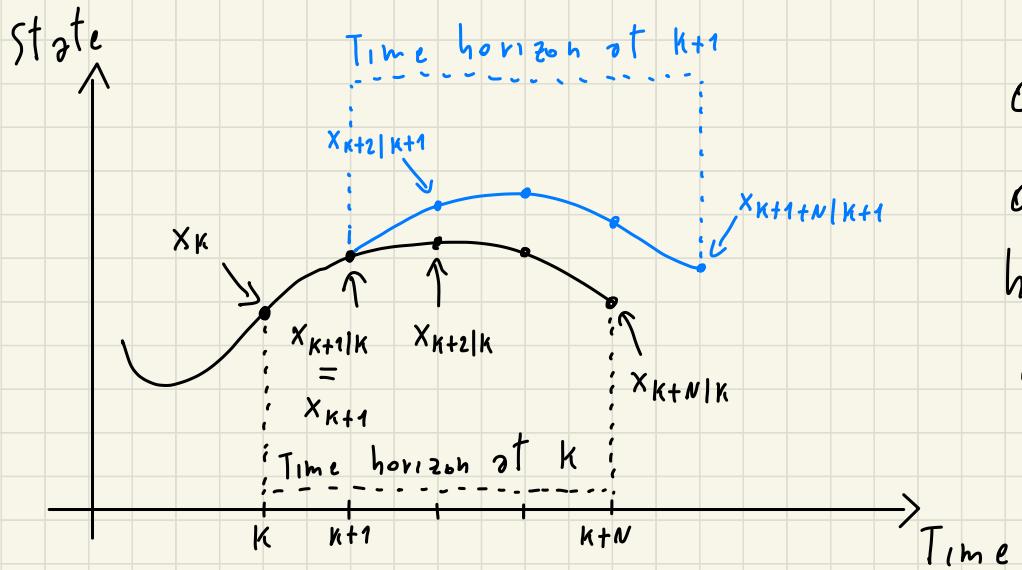
- Do not use the rest of  $u^*$ , nor  $x^*$

- Repeat the same process at next control loop

Remain optimal  
despite changes  
in environment

$\rightarrow$  to keep being  
updated with  
environment

## PREDICTED VS ACTUAL TRAJECTORIES



OCP at time  $k$  and  $k+1$   
optimize over different  
horizons so can result in  
different trajectories

Even without disturbances, predicted and actual trajectories can be different!

## CHALLENGES IN MPC

### ① FEASIBILITY

- Can I ensure the OCP is always feasible?
- What can I do if it's not feasible?

### ② STABILITY

- Can I ensure MPC stabilizes the system?

### ③ COMPUTATION TIME

- Can I solve the OCP sufficiently fast?

## INFINITE-HORIZON MPC

- If  $N = \infty$  the time horizon seen by the solved OCP's is the same
- In this case predicted and actual trajectories are the same (assuming no disturbances)
- This follows from Bellman's optimality principle:

Given optimal traj.  $x^*(0) \dots x^*(N)$ , the subtraj.  $x^*(k) \dots x^*(N)$  is optimal for OCP on horizon  $[k, N]$  starting from  $x(k)$

## INFINITE HORIZON STABILITY & FEASIBILITY

With  $N = \infty$ , if cost  $\ell(x, u) \geq d \|x\|$   $\forall x, u$  for some  $d > 0$ . Then having a finite cost implies stability.

Moreover, since predicted and actual traj. are the same, feasibility of the first OCP implies feasibility of all successive OCP's (recursive feasibility)

## IDEAS:

Add terminal cost and constraints to finite-horizon OCP to mimic an infinite horizon.

Reminiscent of VALUE FUNCTION

## TERMINAL COST & CONSTRAINTS

$$V_N^*(\bar{x}) = \underset{\mathcal{U}}{\text{minimize}} \sum_{n=0}^{N-1} l(x_n, u_n) + l_f(x_N)$$

TERMINAL  
COST

$$\text{subject to } x_{n+1} = f(x_n, u_n) \quad k = 0 \dots N-1$$

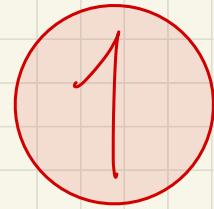
$$x_{n+1} \in \mathcal{X}, \quad u_n \in \mathcal{U} \quad k = 0 \dots N-1$$

$$x_0 = \bar{x}$$

$$x_N \in \mathcal{X}_F$$

TERMINAL  
CONSTRAINT

How do we choose  $l_f(\cdot)$ ,  $\mathcal{X}_F$ ?



FEASIBILITY

## FEASIBILITY

Is the OCP going to be feasible at all sampling instants +?

- INPUT CONSTRAINTS ONLY  $\Rightarrow$  Always feasible
- HARD STATE CONSTRAINTS  $\Rightarrow$  If  $N < \infty$  there is no guarantee that OCP remains feasible, even in the nominal case.

Maximum output admissible set theory:

- $N < \infty$  is enough
- Hard to know which value of  $N$  is enough
- Computation time may be too high for such  $N$

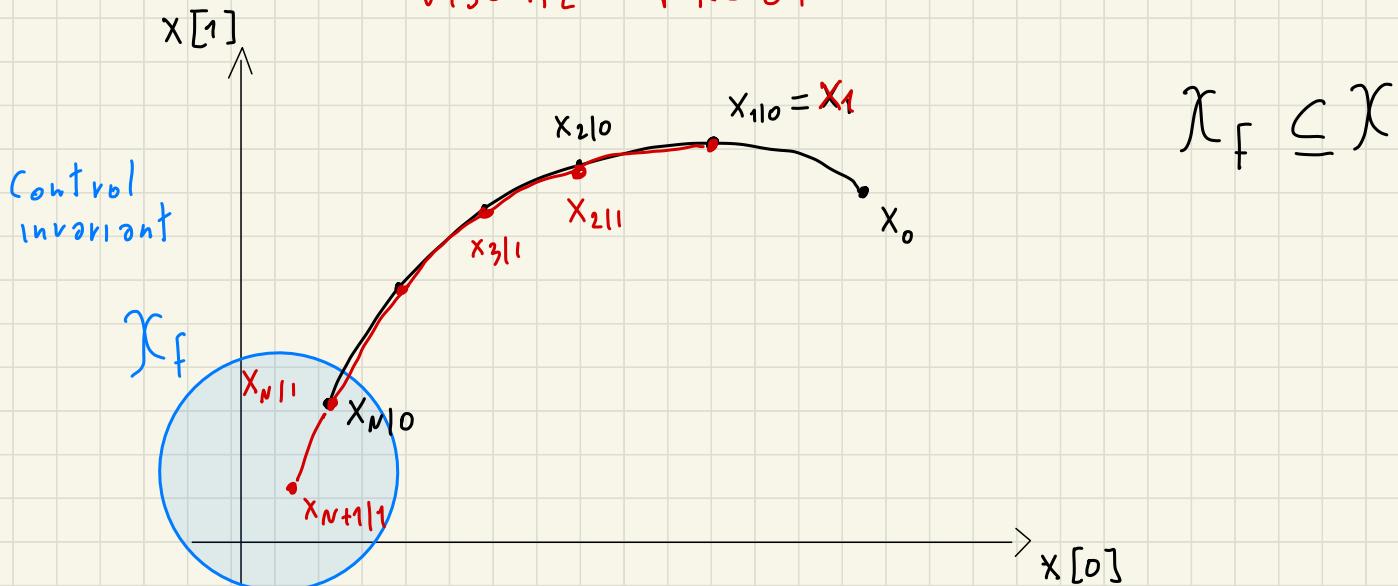
## THEOREM

If  $X_f$  is control-invariant  $\Rightarrow$  MPC is recursively feasible

### DEFINITION

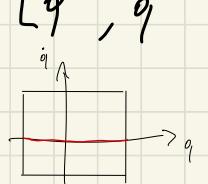
- A set  $S$  is control invariant if  $\forall x \in S, \exists u \in U$  s.t.  $f(x, u) \in S$   
i.e. if you start in  $S$ , you can remain in  $S$ .

### VISUAL PROOF



## CONTROL INVARIANT SETS

Example: for a **manipulator**, set of zero velocity states is control invariant if  $|g(q)| \leq \tau^{\max}$  if  $q \in [q^{\min}, q^{\max}]$

$$S = \{(q, \dot{q}) \mid q^{\min} \leq q \leq q^{\max}, \dot{q} \in \mathbb{R}^n\}$$


How do I compute C.I. sets in general?

- Hard problem for nonlinear systems
- Set of equilibria  $S$  is always C.I., but small:  
$$S = \{x \in X \mid \exists u \in U : x = f(x, u)\}$$
- Can use backward reachable set of  $S$ , i.e. states starting from which you can reach  $S$   
→ Aggiungi schema

## COMPUTING N-STEP BACKWARD REACHABLE SETS

- We can easily check whether a state  $\bar{x} \in \mathcal{B}_n(S)$  ( $N$ -step backward reachable set of  $S$ ) by solving this OCP:

minimize  
 $x, u$

subject to     $x_{i+1} = f(x_i, u_i) \quad i = 0 \dots N$   
 $x_{i+1} \in \mathcal{X}, \quad u_i \in \mathcal{U} \quad i = 0 \dots N$       (1)

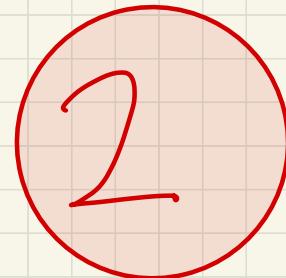
$$x_N \in S \quad \xleftarrow{\hspace{1cm}} \quad x_N = x_{N+1}$$

$$x_0 = \bar{x}$$

- If a solution exists  $\Rightarrow \bar{x} \in \mathcal{B}_n(S)$
- Otherwise  $\Rightarrow \bar{x} \notin \mathcal{B}_n(S)$

- Use function approximation techniques to learn an approximate representation of  $B_N(s)$  from examples
- To build dataset, repeat the following steps:
  - Sample state  $\bar{x}$
  - If  $\bar{x} \notin \mathcal{X} \Rightarrow$  Add data point  $(\bar{x}, \emptyset)$  to dataset
  - Else, try to solve OCP (1)
  - If solution found  $\Rightarrow$  Add  $(\bar{x}, 1)$  to dataset
  - Else  $\Rightarrow$  Add  $(\bar{x}, 0)$  to dataset
- When dataset is built, train classifier (e.g. Neural network) to predict label (0 or 1) given state.
- Finally, use classifier as terminal constraint in MPC

FINE



STABILITY

# LYAPUNOV FUNCTIONS

Def:  $V: \mathbb{R}^n \rightarrow \mathbb{R}$  is exponential Lyapunov function if  $\exists d_1, d_2, d_3 > 0$   
s.t.  $\forall x \in \mathcal{X}$ :

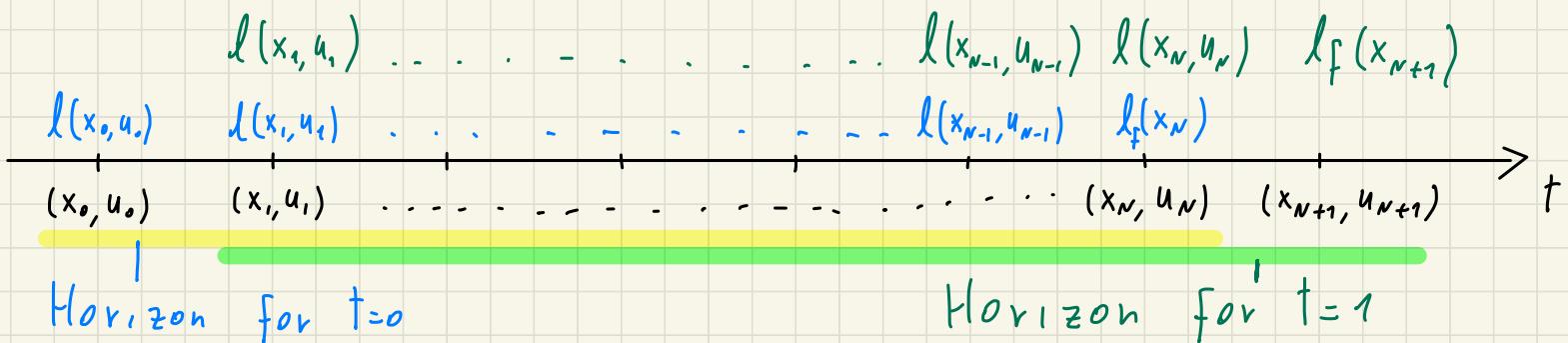
$$d_1 |x| \leq V(x) \leq d_2 |x| \Rightarrow \begin{cases} V(x) > 0 & \forall x \neq 0 \\ V(0) = 0 \end{cases}$$

$$V(f(x)) - V(x) \leq -d_3 |x| \quad \Rightarrow \quad V(x_s) \leq V(x_0) - d_3(x_0)$$

Th: If  $\exists$  Lyap. func.  $\Rightarrow$  origin is exponentially stable in  $\mathcal{X}$ .

$$\|x_n\| \leq c \gamma^n \|x_0\| \text{ for some } c > 0, \gamma \in [0, 1]$$

# VALUE FUNCTION AS LYAPUNOV FUNCTION



$$V_o^*(x_o) = \sum_{i=0}^{N-1} l(x_i, u_i) + l_f(x_N)$$

$$V_1(x_1) = \sum_{i=1}^N l(x_i, u_i) + l_f(x_{N+1})$$

$$V_1(x_1) - V_o^*(x_o) = l(x_N, u_N) + l_f(x_{N+1}) - l_f(x_N) - l(x_o, u_o) \leq -d_3 |x_o| \quad \forall x_o \in X_f$$

Sufficient but not necessary because  $V_1(x_1)$  is not optimal!

## STABILITY ASSUMPTIONS

Ass #1:  $f(\theta, \theta) = \theta, l(\theta, \theta) = \theta, l_f(\theta) = \theta$

Ass. #2:  $\forall x \in \mathcal{X}_f, \exists u \in \mathcal{U}$  s.t.:

$$l_f(f(x, u)) - l_f(x) \leq -l(x, u) \quad \text{Terminal cost decrease}$$

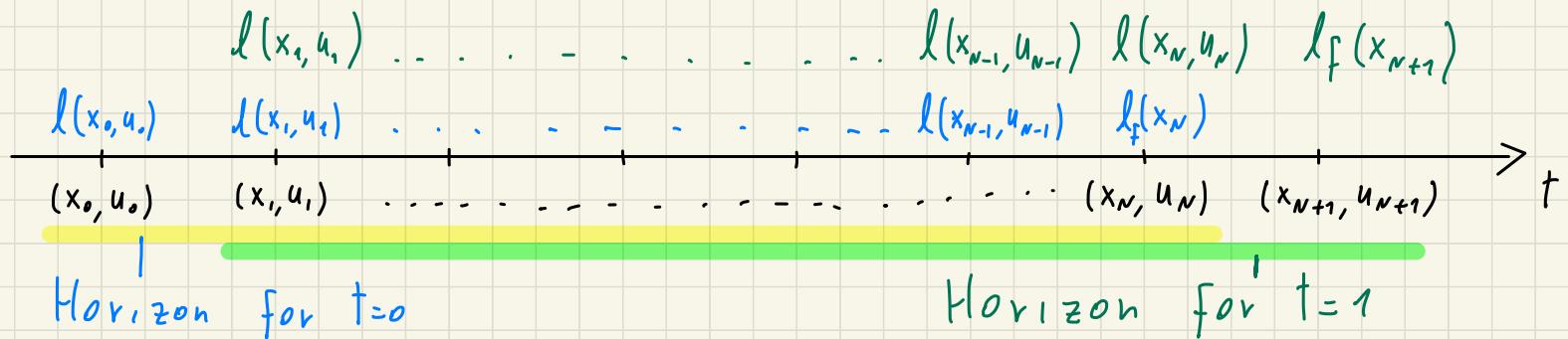
Ass. #3:  $\exists d_1 > 0$  s.t.  $\forall x \in \mathcal{X}_{N_1} \forall u \in \mathcal{U}$

Set of states from which O.C.P. has a solution

$$l(x, u) \geq d_1 \|x\| \quad \text{Cost lower bound}$$

If #1, #2, #3 hold  $\Rightarrow V^*(\cdot)$  is an exp. Lyapunov function

# VALUE FUNCTION AS LYAPUNOV FUNCTION



$$V_o^*(x_o) = \sum_{i=0}^{N-1} l(x_i, u_i) + l_f(x_N)$$

$$V_1(x_1) = \sum_{i=1}^N l(x_i, u_i) + l_f(x_{N+1})$$

$$V_1(x_1) - V_o^*(x_o) = \underbrace{l(x_N, u_N) + l_f(x_{N+1})}_{\leq 0} - \underbrace{l_f(x_N) - l(x_o, u_o)}_{\leq -d_1 |x_o|} \leq -d_3 |x_o| \quad \forall x_N \in X_f$$

by Ass. #2                                  by Ass. #3

## STABILITY ASSUMPTIONS

Note that assuming  $f(0,0)=0$  is not restrictive, if we want to stabilize  $x^* \neq 0$ , s.t.  $f(x^*, u^*) = x^*$ , i.e.  $x^*$  must be an equilibrium  $\Rightarrow$  Define new state  $\bar{x} = x - x^*$ , control  $\bar{u} = u - u^*$ , dynamics:

$$\left. \begin{array}{l} x^+ = f(x, u) \\ x^* = f(x^*, u^*) \end{array} \right\} \quad \begin{aligned} x^+ - x^* &= f(\bar{x} + x^*, \bar{u} + u^*) - x^* \\ \bar{x}^+ &= \bar{f}(\bar{x}, \bar{u}) \end{aligned}$$

So that  $\bar{f}(0, 0) = f(0 + x^*, 0 + u^*) - x^* = 0$

# STABILITY OPTIONS

① Pick  $\mathcal{X}_f$  and  $l_f$  (e.g. using sampling-based approach and function approximators) such that:

-  $\mathcal{X}_f$  is control invariant

-  $\forall x \in \mathcal{X}_f \exists u \in \mathcal{U}$  s.t.  $l_f(f(x, u)) - l_f(x) \leq -\ell(x, u)$

② Pick  $\mathcal{X}_f = \{0\} \Rightarrow$  small basin of attraction

② is special case of ①

③ Pick  $N$  "large enough" and set  $l_f(x) = 0$ ,  $\mathcal{X}_f = \mathcal{X}$



Most common in practice, especially for nonlinear systems!

Downside: Large  $N \Rightarrow$  large computation time!

## STABILITY EXTENSIONS

- What if system is time varying?
  - e.g., trajectory tracking can be cast as regulation of time-varying system
- Stability + recursive feasibility can still be ensured with time-varying  $\chi_f$ ,  $l_f(\cdot)$
- What if cost measures some kind of energy consumption  
 $\Rightarrow$  does not satisfy  $l(x, u) \geq d|x|$   
 $\Rightarrow$  Economic MPC

### ③ COMPUTATION TIME

- KEY IDEA: **warm start**, i.e. use solution of previous OCP as initial guess for current OCP
- Shift trajectory back by 1 time step:  $u_n^{guess} \leftarrow u_{n+1}^*$
- Use zero as initial guess for last time step:  $u_{n-1}^{guess} \leftarrow 0$
- **Don't iterate until convergence**, just do 1 Newton step
  - Possibly skip line search
- Use hardware accelerators (e.g. GPU)

# REINFORCEMENT LEARNING

- Born in Computer Science community
- Many similarities with Optimal Control, but RL
  - tries to find **globally optimal policy**
  - assumes dynamics is unknown **MAIN DIFFERENCE**
  - initially focused on **finite size state/control spaces**
  - uses different terminology → same concepts but with different names
  - typically assumes dynamics is **stochastic**
    - we will not in this course, to simplify equations
  - typically uses **infinite horizon**

# TERMINOLOGY

RL

OC

State s

Action a

Environment

Reward

Return

Maximize

Value function

Optimal Value function

State x

Control u

Plant

Cost

Cost to go

Minimize

Cost to go (of a policy)

Value function

# MARKOV DECISION PROCESSES (MDP)

- Describe environment for RL problem
- Fully observable
- Markov property: Future is independent of the past given the present

$$P(x_{t+1} | x_t) = P(x_{t+1} | x_1, \dots, x_t)$$

- State transition probability:  $P_{xx'} = P(x_{t+1} = x' | x_t = x)$

- State transition matrix:

• Each row sums to 1  
total probability

• In deterministic systems

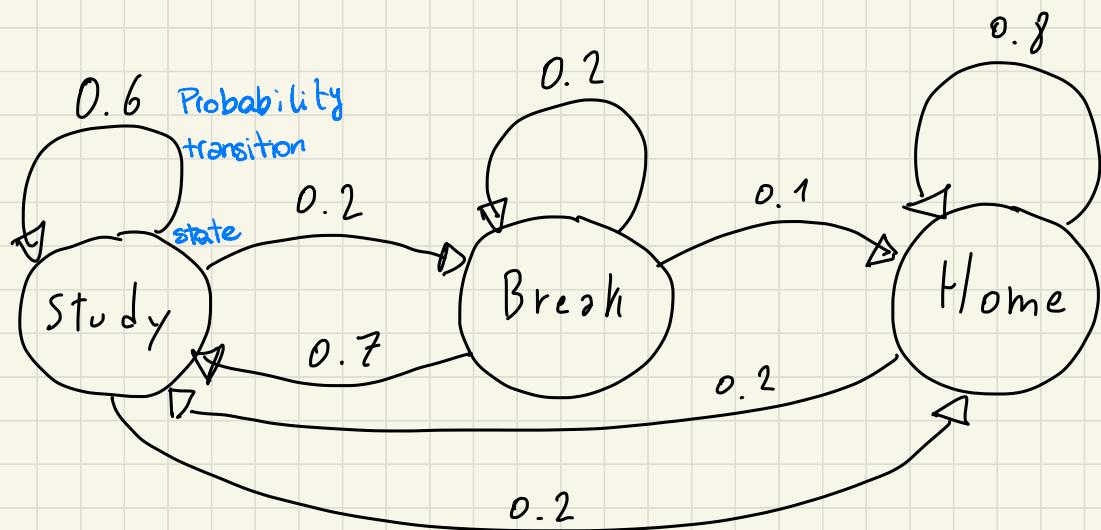
$P$  contains only 0 and 1

$$P = \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}$$

# MARKOV PROCESS

- Tuple  $\langle X, P \rangle$ 
  - $X$  = (finite) set of states
  - $P$  = state transition probability matrix
- A.K.A. Markov Chain

$$\begin{matrix} & S & B & H \\ S & \left[ \begin{matrix} 0.6 & 0.2 & 0.2 \\ 0.7 & 0.2 & 0.1 \\ 0.2 & 0 & 0.8 \end{matrix} \right] = P \end{matrix}$$



# STATE TRANSITION PROBABILITY MATRIX

- Since  $P$  has all nonnegative entries, and it's irreducible, by Perron-Frobenius theorem we know that:  
Associated to a strongly connected graph

PROPERTY 1:

- The largest (in norm) eigenvalue  $r$  of  $P$  satisfies:

$$\min_i \sum_j P_{ij} \leq r \leq \max_i \sum_j P_{ij} \Rightarrow r = 1$$

PROPERTY 2:

- All eigenvalues of  $P$  are smaller than 1:

$$\lambda_i(P) < 1 \quad \forall i$$

IMPORTANT

$$\|Pv\| \leq \|v\|$$

Directly consequence of Property 1.

$Pv$  can only

restrict the vector

# PROBABILITY MATRIX vs DYNAMICS FUNCTION

In OCP for deterministic systems dynamics was encoded as:

CONTINUOUS OR DISCRETE STATE SPACE  $\xleftarrow{\text{PLUS}} x^+ = f(x)$

$x^+ = f(x, w)$   $w \sim P$   $\xrightarrow{\text{uncertainty that follow}}$   
 $\xrightarrow{\text{an uncertainty prob. (usually known)}}$   
 $\xrightarrow{\text{uncertainty}}$

In MP dynamics is encoded as probability density func.:

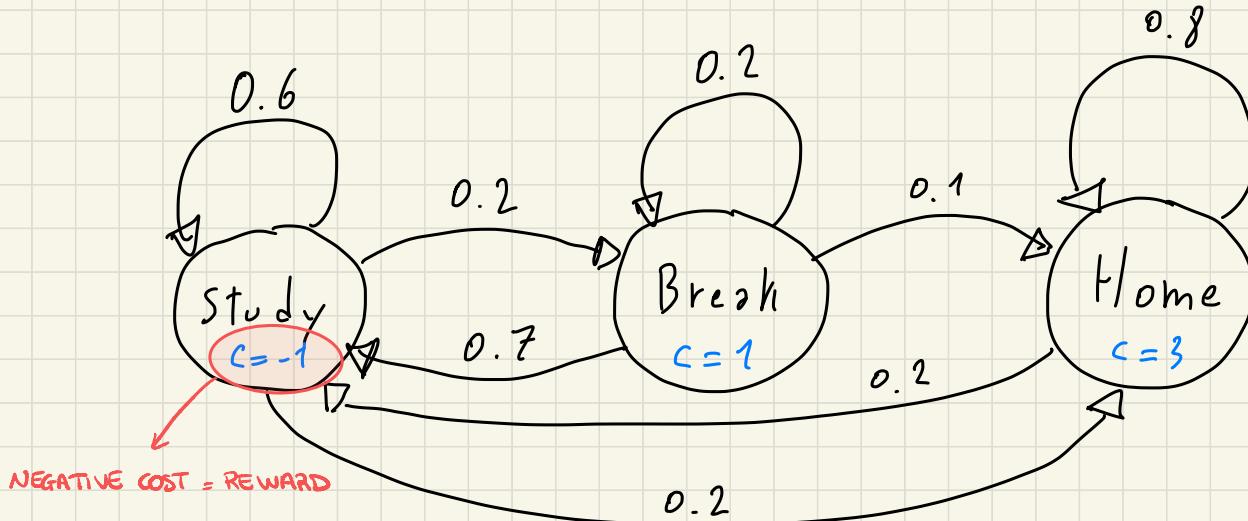
ONLY WITH DISCRETE STATE SPACE  $\xleftarrow{\text{IP}} \text{IP}(x^+ | x)$  Probability of transitioning from state  $x$  to next state  $x^+$

If system is deterministic, both approaches can be used.

We will favor the first one whenever possible.

# MARKOV REWARD PROCESS

- Tuple  $\langle \mathcal{X}, P, C, \gamma \rangle$ 
  - $C$  = cost function  $C_x = \mathbb{E} [l_{t+1} | x_t = x]$
  - $\gamma$  = discount factor  $\gamma \in [0, 1]$  PART OF PROBLEM DEFINITION  
↳ Since we work with  $\infty$  horizon, it prevents cost to go to  $\infty$



## COST TO GO (RETURN)

- Total discounted cost from time  $t$  to  $\infty$ :

$$J_t = \underset{\text{cost}}{l_t} + \gamma \underset{\text{cost of next state}}{l_{t+1}} + \dots = \sum_{k=0}^{\infty} \gamma^k l_{t+k}$$

IF  $\gamma \rightarrow 1$  present and future are almost the same

- $\gamma$  represents preference for later costs over immediate costs
  - $\gamma$  close to 0  $\Rightarrow$  myopic evaluation
  - $\gamma$  close to 1  $\Rightarrow$  far-sighted evaluation
- Why discount?
  - Avoid infinite cost
  - Uncertainty about the future
  - Immediate cost could require interest

$\uparrow\downarrow\gamma$   $\downarrow\downarrow$  convergence

$\gamma$  = discount factor

$0 \leq \gamma \leq 1$

# VALUE FUNCTION (COST-TO-GO)

- Cost starting from state  $x$ :

$$V(x) = J_t (x_t = x) = l_t + \gamma l_{t+1} + \gamma^2 l_{t+2} + \dots$$

- Can be decomposed in two parts:

$$l_t + \underbrace{\gamma (l_{t+1} + \gamma l_{t+2} + \dots)}_{V(x_{t+1})}$$

$$V(x) = l(x) + \gamma \underbrace{V(f(x))}_{\text{same thing written in different ways}}$$

- This is the Bellman Equation

- Can be expressed in matrix form

$$V = C + \gamma P V$$

Method to compute the value of next state  $V(F(x))$

$$\begin{bmatrix} V(1) \\ \vdots \\ V(h) \end{bmatrix} = \begin{bmatrix} l(1) \\ \vdots \\ l(h) \end{bmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1h} \\ \vdots & \ddots & \vdots \\ P_{h1} & \dots & P_{hh} \end{bmatrix} \begin{bmatrix} V(1) \\ \vdots \\ V(h) \end{bmatrix}$$

Compute the expected value to go from state 1 to other n states

Weighted average in which weight is the probability

- Linear equation in  $V$ :

$$V = (I - \gamma P)^{-1} C$$

- Direct solution only possible for small MRP's
- Iterative methods used for large MRP's

# MAROV DECISION PROCESS (MDP)

- Tuple  $\langle X, \mathcal{U}, P, C, \gamma \rangle$  We introduce control input (Action  $a$ )
- $\mathcal{U}$  = finite set of control inputs
- $P_{xx'}^u = \Pr(X_{t+1} = x' \mid X_t = x, U_t = u)$
- $C_x^u = l_t(X_t = x, U_t = u)$
- Policy = distribution over control inputs given states.  
↳ Stochastic map from state to control  
 $\Pi(u|x) = \Pr(U_t = u \mid X_t = x)$
- We will assume deterministic policies:  $u = \Pi(x)$

We can use stochastic policy during policy search process

# ACTION-VALUE FUNCTION

Cost starting from state  $x$ , applying input  $u$ , and then following policy  $\pi$

$$Q^\pi(x, u) = J + \left( x_1 = x, u_1 = u, u_{k>1} \sim \pi \right)$$

Action and policy are independent cost starting from a given state, applying a given action

You don't know the future cost, it depends on the next action you decide to take

Can be decomposed as:

$$\begin{aligned} Q^\pi(x, u) &= l(x, u) + \gamma Q^\pi(f(x, u), \pi(f(x, u))) \\ &= l(x, u) + \gamma V^\pi(f(x, u)) \end{aligned}$$

$\pi$  evaluated at the next state

First action is the one that doesn't come from action set  $A$

where:

$$V^\pi(x) = Q^\pi(x, \pi(x))$$

Relation between value function and action-value function

# OPTIMAL VALUE FUNCTION

Policy that minimize the return

- Minimum Value function over all policies:

$$V^*(x) = \min_{\pi} V^\pi(x)$$

- Optimal action-value function:

$$Q^*(x, u) = \min_{\pi} Q^\pi(x, u)$$

- Optimal policy:

$$\pi^* \leq \pi \quad \forall \pi$$

[where  $\pi \leq \pi'$  if  $V^{\pi'}(x) \leq V^{\pi}(x) \quad \forall x$ ]

# OPTIMAL POLICY

- Minimize  $Q^*$  over  $u$ :

$$\pi^*(x) = \arg \min_{z \in V} Q^*(x, z)$$

- If we know  $Q^*(x, u)$  we immediately have optimal policy

$$\pi^*(x) = \underset{u}{\min} l(x, u) + \gamma V(f(x, u))$$

## BELLMAN OPTIMALITY EQUATIONS

$$V^*(x) = \min_u Q^*(x, u)$$

$$Q^*(x, u) = l(x, u) + \gamma V^*(f(x, u))$$

$$\left. \begin{array}{l} V^*(x) = \min_u l(x, u) + \gamma V^*(f(x, u)) \\ Q^*(x, u) = l(x, u) + \gamma \min_{u'} Q^*(f(x, u), u') \end{array} \right\} \begin{array}{l} \text{• Nonlinear} \\ \text{• No closed-form solution} \\ \text{• Iterative algorithms} \end{array}$$

# DYNAMIC PROGRAMMING

- Assume full knowledge of an MDP
- Two classes of problems:
  - ① PREDICTION:
    - Input: MDP  $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, \mathcal{C}, \gamma \rangle$  and policy  $\pi$
    - Output: Value function  $V^\pi$
  - ② CONTROL:
    - Input: MDP
    - Output: Optimal value function  $V^*$  and Optimal policy  $\pi^*$

# FINITE VS INFINITE HORIZON

- Finite horizon:

- theory is simpler
- mostly used in optimal control
- policy and Value depend on time

- Infinite horizon:

- theory is more complex, but elegant
- mostly used in RL
- policy and Value are stationary
- good approximation of problems with finite but long horizon

# BELLMAN OPERATORS

- Bellman (expectation backup) operator  $\bar{T}^\pi$ :

$$\bar{T}^\pi(V) = C^\pi + \gamma P^\pi V$$

- Bellman optimality (backup) operator  $\bar{T}$ :

$$\bar{T}(V) = \min_{u \in \mathcal{U}} C^u + \gamma P^u V$$

Different  
Notation  $\rightarrow$

$$(\bar{T}V)(x) = \min_{u \in \mathcal{U}} l(x, u) + \gamma V(f(x, u))$$

## PREDICTION: ITERATIVE POLICY EVALUATION

- Problem: evaluate a given policy  $\pi$
- Solution: apply iteratively Bellman expectation operator:

$$V_{n+1}(x) = l(x, \pi(x)) + \gamma V_n(f(x, \pi(x))) \quad \forall x \in \mathcal{X}$$

- In short:  $V_{n+1} = T^\pi V_n$
- Start from arbitrary Value function  $V_0$
- Guarantee to converge to  $V_\pi$ :  $\lim_{n \rightarrow \infty} V_n = V_\pi$

# CONVERGENCE PROOF OF POLICY EVALUATION

$$V_{k+1} = T^\pi V_k$$

Prove  $T^\pi$  is contracting

$$\|T^\pi V - T^\pi Z\|_\infty = \|\cancel{C} + \gamma P^\pi V - \cancel{C} - \gamma P^\pi Z\|_\infty = \gamma \|P^\pi(V - Z)\|_\infty$$

$$\leq \gamma \|V - Z\|_\infty$$

Because each row of  $P^\pi$  sums to 1

Given  $V_0$  and  $V_\pi$  show  $V_k \rightarrow V_\pi$  as  $k \rightarrow \infty$ :

$$\|V_k - V_\pi\|_\infty = \|T^\pi V_{k-1} - T^\pi V_\pi\|_\infty \leq \gamma \|V_{k-1} - V_\pi\|_\infty \leq \gamma^2 \|V_{k-2} - V_\pi\|_\infty$$

$$\leq \dots \leq \gamma^k \|V_0 - V_\pi\|_\infty$$

$V_k$  converges to  $V_\pi$  at geometric rate.

Prove  $V_\pi$  is unique fixed point of  $T^\pi$ . Assume  $V$  and  $Z$  are both fixed points of  $T^\pi$ :

$$T^\pi V = V \quad T^\pi Z = Z$$

Prove  $V = Z$ :

$$\|T^\pi V - T^\pi Z\|_\infty \leq \gamma \|V - Z\|_\infty \leftarrow \text{Because } T^\pi \text{ is contracting}$$

$$\|T^\pi V - T^\pi Z\|_\infty = \|V - Z\|_\infty \leftarrow \text{Because } V, Z \text{ are fixed points}$$

$$\gamma \|V - Z\|_\infty \geq \|V - Z\|_\infty \Rightarrow \|V - Z\|_\infty = 0 \quad \text{Because } \gamma < 1$$

## CONTROL: POLICY ITERATION

- Problem: Find optimal policy  $\pi^*$
- Solution:
  - Start with arbitrary policy  $\pi_0$ .
  - For  $k = 0 \dots \infty$ 
    - Evaluate  $\pi_k \Rightarrow V^{\pi_k}$
    - Improve policy acting greedily w.r.t.  $V^{\pi_k}$ :
$$\pi_{k+1}(x) = \arg \min_u l(x, u) + \gamma V^{\pi_k}(f(x, u))$$
- Always converges to  $\pi^*$ :  $\lim_{k \rightarrow \infty} \pi_k = \pi^*$

## CONVERGENCE OF POLICY ITERATION

Show that  $\pi_{n+1} \leq \pi_n \quad \forall n \iff \pi' \leq \pi$

$$\pi'(x) = \arg \min_u [l(x, u) + \gamma V^\pi(f(x, u))]$$

$$\underbrace{\min_u [l(x, u) + \gamma V^\pi(f(x, u))] \leq \underbrace{l(x, \pi(x)) + \gamma V^\pi(f(x, \pi(x)))}$$

Cost we get by following  $\pi'$  for first step and then following  $\pi$

$$Q^\pi(x, \pi') = \min_u [Q^\pi(x, u)] \leq Q^\pi(x, \pi(x)) = V^\pi(x)$$

$$V^\pi(x) \geq \min_u [l(x, u) + \gamma V^\pi(f(x, u))]$$

$$\geq l(x, \pi') + \gamma Q^\pi(x', \pi') = \quad x' \triangleq f(x, \pi')$$

$$= l(x, \pi') + \gamma l(x', \pi') + \gamma^2 V^\pi(x'') \quad \begin{cases} \text{Cost following } \pi' \text{ for 2} \\ \text{steps and then } \pi \end{cases}$$

$$\geq l(x, \pi') + \gamma l(x', \pi') + \gamma^2 Q^\pi(x'', \pi') \quad \begin{cases} \pi' \text{ for 3 steps} \end{cases}$$

$$\geq \sum_{i=0}^{\infty} \gamma^i l(x^{(i)}, \pi') = V^{\pi'}(x) \Rightarrow V^{\pi'} \geq V^\pi \iff \pi' \leq \pi$$

## MODIFIED POLICY ITERATION

- Evaluating exactly a policy can take many iterations
- What if we compute an approximation of  $V^{\pi_k}$ 
  - Use  $m_k$  policy evaluation iterations to compute  $V^{\pi_k}$
- Under some mild assumptions it converges to  $V^*$
- For  $m_k = \infty$  we recover Policy Iteration
- For  $m_k = 1$  we get Value Iteration
  - every time you update  $V$  you use a policy that's greedy w.r.t.  $V$

# CONTROL: VALUE ITERATION

Algorithm:

- Start with arbitrary value function  $V_0$
  - For  $k = 0 \dots \infty$ 
    - For any  $x \in \mathcal{X}$ 
      - $V_{k+1}(x) = \min_{u \in \mathcal{U}} l(x, u) + \gamma V_k(f(x, u))$
- $\left. \begin{matrix} \\ \\ \end{matrix} \right\} V_{k+1} = T V_k$

•  $V_k$  converges to  $V^*$ :  $\lim_{k \rightarrow \infty} V_k = V^*$

• Optimal policy  $\pi^*$  is computed at the end from  $V^*$

$$\pi^* = \arg \min_u l + \gamma V^*$$

# CONVERGENCE OF VALUE ITERATION

Three-step prove:

① Prove  $T$  is contracting

② Prove  $V_{k+1} = TV_k$  converges to fixed point  $V = TV$

③ Prove  $T$  has unique fixed point  $V$

① Prove  $T$  is contracting, i.e.  $\|Tz - TV\|_\infty \leq \gamma \|z - V\|_\infty$

Remember definition of  $T$ :  $TV = \min_u [C^u + \gamma P^u V]$

$$\|Tz - TV\|_\infty = \left\| \min_u [C^u + \gamma P^u z] - \min_u [C^u + \gamma P^u V] \right\|_\infty =$$

Exploit fact that:  $\min_y f(y) = -\max_y [-f(y)]$

$$= \left\| -\max_u (-Q^z) + \max_w (-Q^V) \right\|_\infty =$$

$$= \left\| \max_w (-Q^V) - \max_u (-Q^z) \right\|_\infty =$$

Exploit that:  $|\max_y f(y) - \max_w g(w)| \leq \max_y |f(y) - g(y)|$

$$\leq \left\| \max_u [-C^u - \gamma P^u V + C^u + \gamma P^u z] \right\|_\infty =$$

$$= \gamma \left\| \max_u P^u (z - V) \right\|_\infty \leq \gamma \left\| \max_u P^u |z - V| \right\|_\infty =$$

The infinity norm takes the max over the states, so:

$$= \gamma \max_{x,u} P_x^u |z - V|$$

Exploit the fact that:  $|z(x) - V(x)| \leq \|z - V\|_\infty \quad \forall x \in X$

$$\leq \gamma \max_{x,u} P_x^u \mathbf{1} \|z - V\|_\infty$$

$\mathbf{1}$  = vector with  
all 1's

Since each row of  $P^u$  sums to 1:

$$= \gamma \|z - V\|_\infty$$

② Assume  $V_{n+1} = T V_n$ , and  $V$  is fixed point of  $T$ :  $T V = V$

$$\|V_n - V\| = \|TV_{n-1} - TV\|$$

$$\leq \gamma \|V_{n-1} - V\| = \gamma \|TV_{n-2} - TV\|$$

$$\leq \gamma^2 \|V_{n-2} - V\| = \dots$$

$$\leq \gamma^n \|V_0 - V\| \Rightarrow \lim_{n \rightarrow \infty} \|V_n - V\| = 0$$

So  $V_n$  converges to a fixed point of  $T$  at a geometric rate.

Finally let's show that  $T$  has a unique fixed point.  
Let's assume  $V$  and  $Z$  are both fixed points of  $\bar{T}$ :

$$\|V - Z\| = \|\bar{T}V - \bar{T}Z\| \leq \gamma \|V - Z\| \Rightarrow (1 - \gamma) \|V - Z\| \leq 0$$

$$\Rightarrow \|V - Z\| = 0 \Rightarrow V = Z$$

## MODEL-FREE PREDICTION

Problem: Find Value function of unknown MDP by acting and observing

- Monte Carlo (MC) methods can solve this problem
- Simple idea: Value = mean return
- Constraint: can only be applied to episodic MDP
  - All episodes must end

# MONTE CARLO POLICY EVALUATION

Goal: learn  $V_{\pi}$  from episodes of experience under policy  $\pi$

$$x_1, u_1, l_1, \dots, x_k \sim \pi$$

Cost-to-go is total discounted cost:

$$J_t = l_t + \gamma l_{t+1} + \dots + \gamma^{T-1} l_{T-1}$$

Value function is expected cost-to-go under policy  $\pi$ :

$$V^{\pi}(x) = \mathbb{E}_{\pi}[J_t \mid x_t = x] \leftarrow \begin{array}{l} \text{implicitly assume some} \\ \text{stochasticity (dynamics, policy, cost)} \end{array}$$

MC policy evaluation estimates  $V^{\pi}$  as the average cost-to-go.

## FIRST-VISIT MC POLICY EVALUATION

The first time that a state  $x$  is visited in an episode:

- increment counter  $N(x) := N(x) + 1$
- increment total cost  $C(x) := C(x) + J(x)$
- estimate value as  $V(x) = C(x) / N(x)$

By law of large numbers  $V(x) \rightarrow V^\pi(x)$  as  $N(x) \rightarrow \infty$

Another version exists where counter and total cost are incremented every time a state is visited (every visit).

## INCREMENTAL MC UPDATES

IDEA: Compute mean incrementally:

$$V_N = \frac{1}{N} \sum_{j=1}^N J_j = \frac{J_N + \sum_{j=1}^{N-1} J_j}{N} = \frac{J_N + (N-1)V_{N-1}}{N} = V_{N-1} + \frac{1}{N}(J_N - V_{N-1})$$

In non-stationary problems it can be useful to track a running mean, i.e. forget old episodes

$$V(x) := V(x) + \alpha (J - V(x))$$

## TEMPORAL DIFFERENCE LEARNING - TD( $\theta$ )

- Alternative to MC
- Estimate cost-to-go by 1-step lookahead (TD( $\theta$ )):

$$V(x_t) := V(x_t) + \alpha_t (\underbrace{r_t + \gamma V(x_{t+1}) - V(x_t)}_{\text{TD error}})$$

$r_t + \gamma V(x_{t+1})$  TD target = estimated return

- Compare to MC incremental update:

$$V(x_t) := V(x_t) + \alpha_t (J_t - V(x_t))$$

- Policy evaluation is TD( $\theta$ ) with  $\alpha_t=1$  and sampling all states uniformly

## TD( $\theta$ ) PROPERTIES

- TD( $\theta$ ) is stochastic approximation algorithm
- Since Bellman operator has unique fixed point  $\Rightarrow$  If TD( $\theta$ ) converges (and all states are sampled infinitely often) then it must converge to  $V^*$ .
- Convergence is guaranteed if step-size sequence satisfies Robbins-Monro (RM) conditions

$$\sum_{t=0}^{\infty} \alpha_t = \infty$$

$$\sum_{t=0}^{\infty} \alpha_t^2 < \infty$$

[Ex.  $\alpha_t = \frac{\bar{\alpha}}{t}$ ]

- In practice people often use constant step sizes

## MC vs TD

- TD can learn after every step
  - MC must wait until end of episode to know cost
  - TD works in non-terminating environments
  - MC only works for episodic (terminating) environments
  - TD target is biased estimate of  $V_{\pi}(x_t)$
  - Cost-to-go  $J_t$  is unbiased estimate of  $V_{\pi}(x_t)$
  - TD Target is lower variance than  $J_t$ :
    - $J_t$  depends on many random transitions
    - TD target depends on one random transition
- Monte Carlo target*

## MC vs TD

- MC has high variance, zero bias
  - good convergence properties
  - (even with function approximation)
  - not very sensitive to initialization
  - simple to understand and use
- TD has low variance, some bias
  - usually more efficient than MC
  - $\text{TD}(0)$  converges to  $V_\pi(x_t)$
  - (but not always with function approximation)
  - more sensitive to initial value

# AB EXAMPLE

Two states, (A, B), no discounting, 8 episodes

A, 0, B, 0

B, 1

B, 1

B, 1

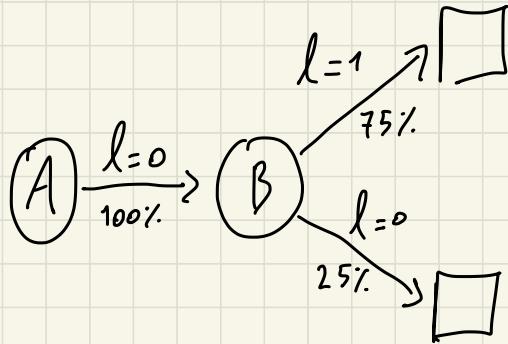
B, 1

B, 1

B, 1

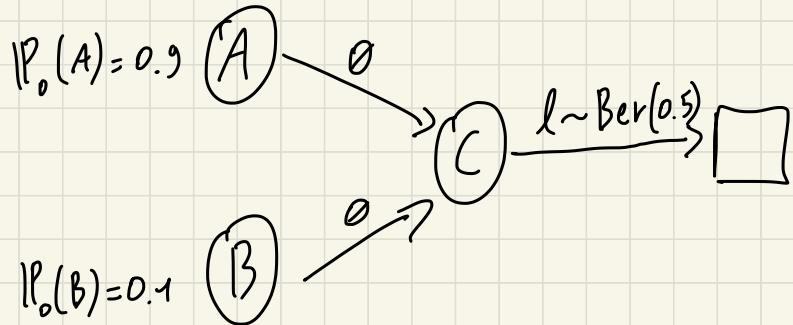
B, 0

$$\begin{cases} V(B) = 0.75 \\ V(A) = \begin{matrix} \nearrow 0 \text{ for MC} \\ \searrow 0.75 \text{ for TD}(\theta) \end{matrix} \end{cases}$$



What's  $V(A)$ ,  $V(B)$ ?

# ABC EXAMPLE



- After  $k$  visits of  $B$ ,  $C$  has been visited  $\approx 10k$  times
- $\text{Var}[\hat{V}(C)] \approx 1/(10k)$

- With TD( $\theta$ ) the variance of the TD Target for  $V(B)$  decreases with the # of episodes

$$\hat{V}(B) := \hat{V}(B) + \alpha (\theta + \hat{V}(C) - \hat{V}(B))$$

- With MC the variance of the Target is always 0.25
  - MC is slower to converge in this example

- Assume how the cost after  $c$  is deterministically 1
- MC becomes faster because it immediately gets the true target value
- TD(0) has to wait for  $\hat{V}(c)$  to converge to 1 before the target of  $V(b)$  gets close to 1
  - Things would be even worse if we had a long chain of states with non-zero cost only at the end

# BOOTSTRAPPING & SAMPLING

Bootstrapping: updates involve an estimate

Sampling: updates sample an expectation

Bootstr. Smpl.	NO	YES
NO	Exhaustive search	DP
YES	MC	TD

# $n$ -STEP RETURN

$$n=1 \quad (\text{TD}) \quad J_t^{(1)} = l_t + \gamma V(x_{t+1})$$

$$n=2 \quad J_t^{(2)} = l_t + \gamma l_{t+1} + \gamma^2 V(x_{t+2})$$

$$\vdots \quad \vdots$$

$$n \quad J_t^{(n)} = l_t + \gamma l_{t+1} + \dots + \gamma^{n-1} l_{t+n-1} + \gamma^n V(x_{t+n})$$

$$n=\infty \quad (\text{MC}) \quad J_t^{(\infty)} = l_t + \gamma l_{t+1} + \dots + \gamma^{T-1} l_{t+T-1}$$

- $n$ -step TD learning:

$$V(x_t) := V(x_t) + \alpha_t (J_t^{(n)} - V(x_t))$$

# CONTROL LEARNING METHODS

MODEL FREE CONTROL

		LEARNED QUANTITIES	
		$V(x) / Q(x, u)$	$V(x) / Q(x, u) + \pi(x)$
MDP	Known	Value Iteration	Policy Iteration
	Unknown	Direct methods	Actor-Critic methods

- Why learn  $Q$  instead of  $V$ ?

$$\pi(x) = \arg \min_u l(x, u) + \gamma V(x') = \arg \min_u Q(x, u)$$

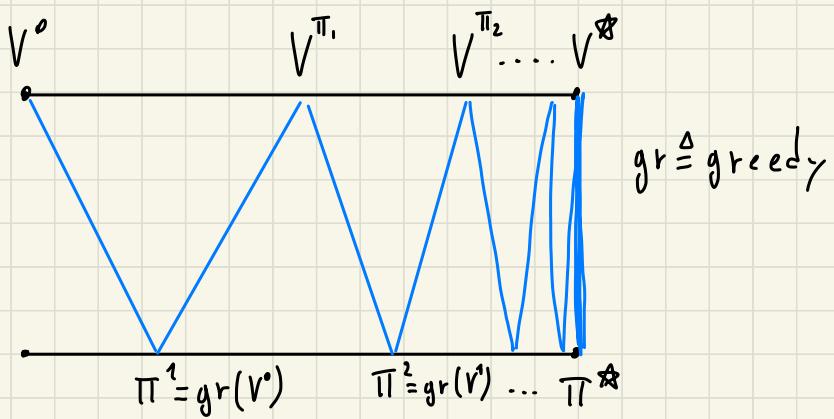
Action-value Function

MODEL FREE CONTROL

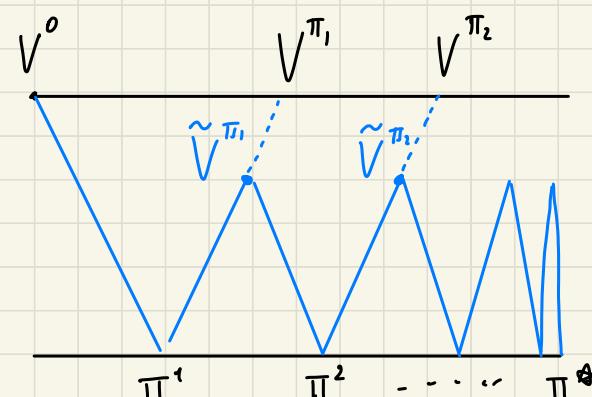
- Cannot compute  $x'$  if MDP is unknown!

# ACTOR-CRITIC METHODS: GENERALIZED POLICY ITER.

- Extension of Policy Iteration to unknown MDP
- PI alternates "policy evaluation" and "policy improvement"
- Unknown MDP  $\Rightarrow$  use MC or TD to evaluate policy
- Exact policy evaluation can take infinitely many samples
- Solution: Improve policy based on approximate value function



PI: exact policy evaluation



GPI: approximate policy evaluation

- The value function is called "critic" because it evaluates the policy, which is called "actor", in order to improve it
- The policy used to generate transition is typically not the same that is evaluated and improved
  - The "behavior policy" mixes small amount of exploration into "target policy"
- GPI can generate policies worse than previous ones

## SARSA: CRITIC

- Use  $\text{TD}(\theta)$  to learn  $Q$  from on-policy samples
- After each transition  $(x, u, l, x', u')$ , update  $Q$  as follows:  
$$\delta(Q) := l(x, u) + \gamma Q(x', u') - Q(x, u)$$
$$Q_{n+1}(x, u) := Q_n(x, u) + \alpha_n \delta(Q)$$
- If  $\pi$  was fixed  $\Rightarrow$  SARSA = TD( $\theta$ )
- As TD, SARSA can diverge in off-policy situations
- How to improve policy?

# HOW TO CHOOSE CONTROL?

- What about being greedy w.r.t.  $Q_n(\cdot)$ ?  $u(x) = \arg \min_w Q_n(x, w)$

## EXAMPLE

- You can choose one of two doors in front of you
  - Left door  $\Rightarrow l = 0 \Rightarrow V(\text{left}) = 0$
  - Right door  $\Rightarrow l = -1 \Rightarrow V(\text{right}) = -1$
  - Right door  $\Rightarrow l = -3 \Rightarrow V(\text{right}) = -2$

Are you sure the right door is the best one?

## SARSA : ACTOR

- Simplest idea: take policy that is **greedy** w.r.t. Q
  - Greedy policy can be computed on the fly  $\Rightarrow$  faster!
- **GREEDY STRATEGY:**
  - Always choose control with best estimated cost
  - Can fail to find best action  $\Rightarrow$  large loss over time
  - Good learner must choose suboptimal control to **explore**
- **$\epsilon$ -GREEDY STRATEGY:**
  - Choose **random** control with **probability  $\epsilon$**
  - Choose **greedy** control with **probability  $1-\epsilon$**
  - Simple way to ensure exploration
  - $\epsilon$  can change over time (typically decreasing to 0)

## SARSA : CONVERGENCE

- Use  $\epsilon$ -greedy as behavior policy to ensure exploration
- SARSA converges to  $Q^*(x, u)$  if:
  - ① Greedy in the Limit with Infinite Exploration (GLIE) policies:
    - All state-action pair are visited infinitely many times
    - Policy converges on greedy policy (e.g.:  $\epsilon_n = \frac{1}{n}$ )
  - ② Robbins-Monro sequence of step sizes:

$$\sum_{n=1}^{\infty} d_n = \infty$$

$$\sum_{n=1}^{\infty} \alpha_n^2 < \infty$$

What if we want to learn off-policy?

# ON-POLICY vs OFF-POLICY LEARNING

## • On-policy (e.g. SARSA)

- learn about policy  $\pi$  from experience sampled from  $\pi$
- behavior policy  $\approx$  target policy (e.g.  $\epsilon$ -greedy, with  $\epsilon$  small)
- limited exploration  $\Rightarrow$  low sample efficiency

## • Off-policy

- learn about policy  $\pi$  from experience sampled from  $\mu$
- behavior policy  $\mu$  can be arbitrarily different from target policy  $\pi$
- higher sample efficiency

## DIRECT METHODS: Q LEARNING

- Iteratively update an estimate  $Q_h(x, u)$  of  $Q^*(x, u)$
- After observing transition  $(x, u, l, x')$  update with:

$$\delta(Q) := l + \gamma \min_{w \in \mathcal{W}} Q(x', w) - Q(x, u)$$

$$Q_{h+1}(x, u) := Q_h(x, u) + \alpha_h \delta(Q_h)$$

- Same as SARSA, but uses  $\min_w Q(x', w)$  instead of  $Q(x', u')$
- Error  $\delta(Q(x, u))$  is independent of behavior policy, which is the one that generates  $u' \Rightarrow$  OFF-POLICY
- Target policy is:  $\pi(x) = \arg \min_u Q(x, u)$

## Q-LEARNING: CONVERGENCE

Need exploration!  
↓

- At stochastic equilibrium,  $V(x, u)$  visited infinitely often:

$$E[\delta(Q)] = 0 = T Q - Q \Rightarrow Q = Q^* \leftarrow \text{unique fixed point of } T$$

- $Q_n$  converges when appropriate learning rates are used (Robbins-Monro conditions) and sufficient exploration is ensured

# SUMMARY

DP

Policy Evaluation

$$V(x) = l + \gamma V(x')$$

Policy Iteration

$$V(x) = l + \gamma V(x')$$

$$\pi(x) = \arg \min_u l + \gamma V(x'(u))$$

Value Iteration

$$V(x) = \min_u l + \gamma V(x'(u))$$

TD

TD Learning

$$V(x) \leftarrow l + \gamma V(x')$$

Sarsa

$$Q(x, u) \leftarrow l + \gamma Q(x', u')$$

$$\pi(x) = \arg \min_u Q(x, u)$$

Q-Learning

$$Q(x, u) \leftarrow l + \gamma \min_{u'} Q(x', u')$$

# VALUE FUNCTION APPROXIMATION

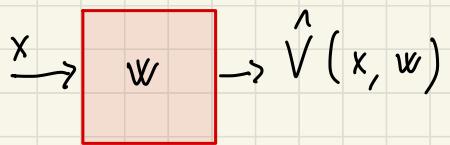
- Discrete-space algorithms don't scale to large systems:
  - Backgammon:  $10^{20}$  states
  - Go:  $10^{170}$  states
  - Robots: continuous state space
- Too many states/control to store and learn
- IDEA: Use function approximation to represent  $V(x) / Q(x, u)$

$$\hat{V}(x, w) \approx V(x)$$

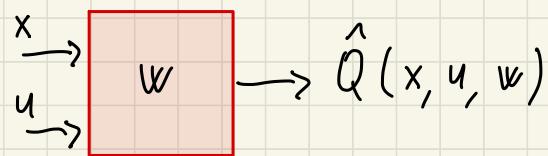
$$\hat{Q}(x, u, w) \approx Q(x, u)$$

- $w$  = parameters defining the function

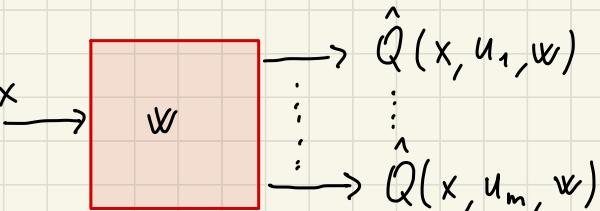
# TYPES OF VALUE FUNCTION APPROXIMATION



- Neural networks are often used as function approximation



- Any differentiable approximator can be used:



- linear combination of features
- Fourier basis
- Polynomials

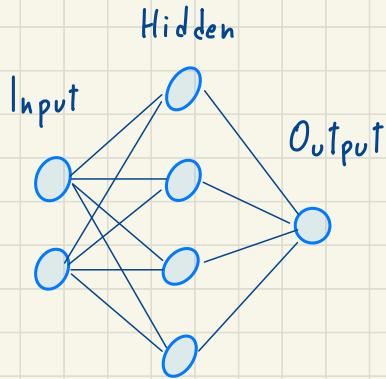
Arbitrary user-defined functions of  $x$

- Problems:

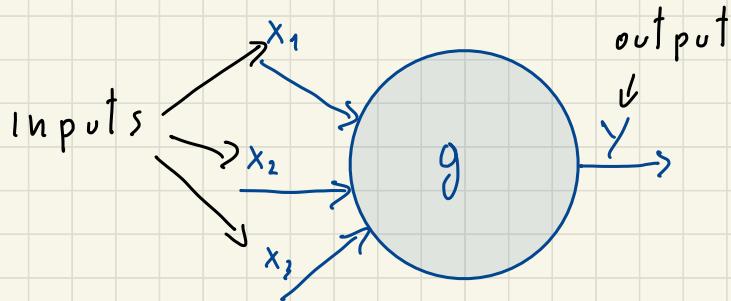
- non-stationary target  $V$
- non i.i.d. data

# NEURAL NETWORKS

- Parametric function structured in "layers":
  - 1 input layer (size = # of inputs)
  - $\geq 1$  hidden layers
  - 1 output layer (size = # of outputs)
- Each layer is composed of "neurons"
- In fully-connected networks each neurons receives as input the outputs of all neurons of the previous layer
- A neuron outputs a nonlinear function of a linear combination of its inputs



# NEURONS



$$y = g \left( \sum_i w_i x_i + w_0 \right)$$

- $g(\cdot)$  is the non-linear "activation" function, e.g.:
  - ReLU, Rectified Linear Unit:  $g(p) = \max(0, p)$
  - Hyperbolic Tangent:  $g(p) = \tanh(p) = \frac{e^p - e^{-p}}{e^p + e^{-p}}$
- $w$  is the vector of weights associated to the inputs  $x$
- $w_0$  is called the "bias"

## TRAINING A NEURAL NETWORK

- Training = Finding the values of all weights  $w$
- Typically formulated as unconstrained optimization:  
$$\underset{w}{\text{minimize}} \sum_i \|\phi(x_i; w) - y_i\|^2 \triangleq J(w)$$

*LOSS  
FUNCTION*
- Given a list of input-output pairs  $(x_i, y_i)$ , find  $w$  such that the network's output  $\phi(x; w)$  is close to  $y$ .
- Since  $\text{size}(w)$  can be very large, typically solved with some variants of gradient descent:

$$w := w - \alpha \nabla J(w)$$

$\alpha$  = step size

## STOCHASTIC GRADIENT DESCENT

- If the number of input-output pairs is large, computing  $\nabla J$  can be computationally expensive

$$\nabla_w J = \sum_{i \in D} 2 (\phi(x_i; w) - y_i) \nabla_w \phi(x_i; w)$$

- IDEA: Compute  $\nabla J$  using only a **subset** of the data
- In the limit the subset (**mini batch**) can be just one sample
- Gradient is approximated (i.e. noisy) so:
  - **more iterations** needed to converge
  - each iteration is much **cheaper to compute**

# INCREMENTAL PREDICTION ALGORITHMS

- How can we get target  $y = V(x)$ ?

- For MC, target is cost-to-go  $J_t(x)$ :

$$\Delta w = \alpha (J_t - \phi(x_t, w)) \nabla_w \phi(x_t, w)$$

- MC converges to local optimum even with nonlinear approximators

- For TD(0), target is TD target:

$$\Delta w = \alpha (l_t + \gamma \phi(x_{t+1}, w) - \phi(x_t, w)) \nabla_w \phi(x_t, w)$$

- TD converges (slow) to global optimum with linear approxm.

# INCREMENTAL CONTROL WITH FUNCTION APPROXIM.

IDEA: Generalized Policy Iteration with approximate Action-Value function:

$$\min_{\psi} \mathbb{E}_{\pi} \left[ (\phi(x, u, w) - Q_{\pi}(x, u))^2 \right]$$

- Use SGD to find local minimum
  - For MC, target is cost-to-go  $J_T(x_t, u_t)$
  - For TD(0), target is  $l_t(x_t, u_t) + \gamma \phi(x_{t+1}, u_{t+1}, w)$
- Use  $\epsilon$ -greedy policy to ensure exploration
  - Greedy control is easily computed for discrete control space

# CONVERGENCE OF CONTROL ALGORITHMS

	VALUE ENCODING		
	Table lookup	Linear	Nonlinear
MC control	✓	(✓)	X
Sarsa	✓	(✓)	X
Q-Learning	✓	X	X

(✓) = chatter around optimum

$$\text{Linear} \Rightarrow \phi(x, w) = h(x)^T w \quad h(x) = \text{nonlinear function}$$

## LIMITATIONS OF INCREMENTAL CONTROL

- Each transition  $(x, u, l, x')$  is used to perform a single gradient step on  $w$  and then thrown away  
=> POOR DATA EFFICIENCY
- Consecutive transitions are strongly correlated because the state evolves continuously  
=> DATA IS NOT INDEPENDENTLY & IDENTICALLY DISTRIBUTED

Can we overcome these issues?

## BATCH LEARNING

- IDEA: Instead of learning incrementally, learn after having collected a batch of training data
- Least-squares predictions

$$\min_w \sum_{t=1}^T (V_t - \phi(x_t, w))^2$$

- How to achieve it? Experience replay:
  - store data  $\langle x_t, V_t \rangle$  in replay buffer
  - sample minibatch from buffer
  - apply SGD to update  $w$
  - converge to LS solution

## DEEP Q-NETWORK (2015, Mnih et al.)

- Extension of Q-Learning to continuous state space
- Discrete control space
- Use neural network to represent  $Q(x, u, w)$
- $\epsilon$ -greedy policy with  $\epsilon$  decreasing over time
- Use experience replay:
  - Store transitions  $(x_t, u_t, l_t, x_{t+1})$  in replay buffer
  - Every few steps (e.g. 4) sampled randomly from replay buffer a mini-batch of transitions (size 32)
  - Use stochastic gradient descent to optimize  $w$

DQN uses fixed Q targets to stabilize training:

- compute Q targets using old fixed parameters  $w^-$
- optimize MSE betw Q-network and Q-learning Targets:

$$\underset{w}{\text{minimize}} \quad \mathbb{E} [l + \gamma \min_{u'} \phi(x', u', w^-) - \phi(x, u, w)]$$

- every C steps update target weights:

$$w^- := w$$

- Store  $w$  resulting in best performance as performance may get worse during training

## CONCLUSIONS

- Introduce function approximation to scale RL algorithms
- Use stochastic gradient descent for training neural networks
- Incremental alg. update w at every step with current transition
- Batch alg. introduce replay buffer to improve sample efficiency and break data dependency
- Action space is still assumed to be discrete

# **Policy Gradient Methods: Direct Policy Optimization**

---

Andrea Del Prete

December 5, 2025

Optimization and Learning for Robot Control—University of Trento

## Table of contents

---

1. From Value-Based to Policy-Based Methods
2. Policy Gradient Theorem
3. The REINFORCE algorithm
4. Reducing the variance
5. Actor-Critic Policy-Gradient Methods

## **From Value-Based to Policy-Based Methods**

---

## Context: From Value-Based to Policy-Based Methods

**Review:** So far, we have studied Value-Based methods (Q-Learning, SARSA, DQN).

- These methods learn action values ( $Q(x, u)$ ) and select actions based on them (e.g.,  $\epsilon$ -greedy). Their policies would not exist without the action-value estimates.

### Policy Gradient Methods (PGM):

- Learn a **parameterized policy**  $\pi(u|x, \theta)$  directly, which selects actions without consulting a value function.
- A value function may still be used to learn the policy parameter, but it is not required for action selection.
- PGM seek to **minimize cost**  $J(\theta)$  using gradient descent:  
$$\theta_{t+1} = \theta_t - \alpha \nabla J(\theta_t).$$

## Policy Parameterization for Continuous Actions

Policy-based methods naturally handle **continuous action spaces** by learning the statistics of the probability distribution, rather than individual action probabilities.

**Gaussian Policy Parameterization:** The policy  $\pi(u|x, \theta)$  can be defined as a normal (Gaussian) probability density function over a real-valued action  $u$ :

$$\pi(u|x, \theta) \propto \exp\left(-\frac{(u - \mu(x, \theta))^2}{2\sigma(x, \theta)^2}\right)$$

Where the policy parameter vector can be divided into two parts,  
 $\theta = [\theta_\mu, \theta_\sigma]$ .

- **Mean Approximation:**  $\mu(x, \theta_\mu)$ .
- **Standard Deviation Approximation:**  $\sigma(x, \theta_\sigma)$ .

# Why Optimize the Policy ( $\pi$ )? I: Continuous Action Spaces

For robotics systems, actions are usually continuous (e.g., joint torques, joint velocities).

- Value-based methods require efficiently solving  $\operatorname{argmin}_u Q^\pi(x, u)$ . This is a **challenge for continuous or high-dimensional action spaces**.
- Policy-based methods offer practical ways of dealing with large/continuous action spaces.

## Key Policy Advantages:

- naturally handle continuous action spaces.
- optimize what we care about (the performance measure  $J(\theta)$ ).

# Why Optimize the Policy ( $\pi$ )? II: Stochasticity and Simplicity

## 1. Intrinsically Stochastic Policies:

- PGM use a stochastic policy class  $\pi(u|x, \theta)$  that is differentiable with respect to  $\theta$ .
- The action probabilities change smoothly as a function of the learned parameter  $\theta$ .
- PGM enable the selection of actions with arbitrary probabilities and can find stochastic optimal policies (e.g., when bluffing in Poker).
- PGM can approach deterministic policies asymptotically, if optimal.

## 2. Simplicity of the Policy Function:

- The policy function  $\pi$  may be a **simpler function to approximate** than the value function  $Q$  or  $V$ .

## Policy Gradient Theorem

---

# Objective: Minimizing Global Cost $J(\theta)$

## Policy Search Formalization:

- A trajectory  $\tau$  is a sequence of states and actions (e.g., a robot's history).
- $C(\tau)$  is the corresponding cost-to-go.
- $\pi_\theta$  is the policy parameterized by  $\theta$ .

**Objective Function  $J(\theta)$**  : We aim to minimize the expected cost-to-go  $J(\theta)$ , defined as the sum of costs over all possible trajectories, weighted by their probability of occurrence under policy  $\pi_\theta$  :

$$\theta^* = \operatorname{argmin}_\theta J(\theta) = \operatorname{argmin}_\theta E_{\tau \sim \pi_\theta}[C(\tau)]$$

$$J(\theta) = \sum_{\tau} P(\tau|\theta) C(\tau)$$

**Solution:** Use Gradient Descent to find  $\theta^*$ :

$$\theta' = \theta - \alpha \nabla_\theta J(\theta)$$

## Step 1: Calculating the Gradient (Using the Gradient of Sum)

We start by taking the gradient of  $J(\theta)$  with respect to  $\theta$ :

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_{\tau} P(\tau|\theta) C(\tau)$$

We can swap the gradient and the summation (the gradient of a sum is the sum of gradients):

$$\nabla_{\theta} J(\theta) = \sum_{\tau} \nabla_{\theta} P(\tau|\theta) C(\tau)$$

**Challenge:** The term  $\nabla_{\theta} P(\tau|\theta)$  is hard to compute directly, as it depends on the unknown environment dynamics and the policy  $\pi_{\theta}$ .

## Step 1: Gradient Calculation

To transform the sum into a sampleable expectation, we use the Likelihood Ratio trick (or REINFORCE trick).

We multiply by 1, in the form  $\frac{P(\tau|\theta)}{P(\tau|\theta)}$  :

$$\nabla_{\theta} J(\theta) = \sum_{\tau} P(\tau|\theta) \frac{\nabla_{\theta} P(\tau|\theta)}{P(\tau|\theta)} C(\tau)$$

Recall the logarithm identity:  $\nabla_{\theta} \log x = \frac{\nabla_{\theta} x}{x}$  .

Substituting this identity, we get the gradient as an Expectation:

$$\nabla_{\theta} J(\theta) = \sum_{\tau} P(\tau|\theta) \nabla_{\theta} \log P(\tau|\theta) C(\tau)$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}} [\nabla_{\theta} \log P(\tau|\theta) C(\tau)]$$

**Intuition:** This gradient tries to increase the probability of paths ( $\tau$ ) with negative costs  $C(\tau)$  and decrease the probability of paths with positive costs.

## Step 2: Decomposing the Trajectory Probability $P(\tau|\theta)$

A trajectory  $\tau$  is defined by the sequence  $x_1, u_1, l_1, x_2, u_2, l_2, \dots$  over horizon  $H$ .

Under the Markov assumption, the probability of a trajectory  $P(\tau|\theta)$  is the product of the transition probabilities  $p(x_{t+1}|x_t, u_t)$  and the action probabilities defined by the policy  $\pi_\theta(u_t|x_t)$ :

$$P(\tau|\theta) = \prod_{t=1}^H p(x_{t+1}|x_t, u_t) \cdot \pi_\theta(u_t|x_t)$$

We now apply the logarithm and the gradient to this expression:

$$\nabla_\theta \log P(\tau|\theta) = \nabla_\theta \log \left[ \prod_{t=1}^H p(x_{t+1}|x_t, u_t) \cdot \pi_\theta(u_t|x_t) \right]$$

Using the property that the log of a product is the sum of logs:

$$\nabla_\theta \log P(\tau|\theta) = \nabla_\theta \left[ \sum_{t=1}^H \log p(x_{t+1}|x_t, u_t) + \sum_{t=1}^H \log \pi_\theta(u_t|x_t) \right]$$

## Step 2: Eliminating Unknown Dynamics

The crucial step relies on the fact that the gradient ( $\nabla_{\theta}$ ) operates only on the policy parameters  $\theta$ .

- The **transition probability**  $p(x_{t+1}|x_t, u_t)$  depends on the environment dynamics, but it **does not** depend on the policy parameters  $\theta$ .
- Therefore, the gradient of the first term is zero:  
$$\nabla_{\theta} \log p(x_{t+1}|x_t, u_t) = 0.$$

This simplifies the expression significantly :

$$\nabla_{\theta} \log P(\tau|\theta) = \nabla_{\theta} \sum_{t=1}^H \log \pi_{\theta}(u_t|x_t)$$

$$\nabla_{\theta} \log P(\tau|\theta) = \sum_{t=1}^H \nabla_{\theta} \log \pi_{\theta}(u_t|x_t)$$

**The Key Insight:** The environment dynamics model is **not required** to compute the gradient!

## Final Policy Gradient (Plain Policy Gradient)

Substituting the final expression for  $\nabla_{\theta} \log P(\tau|\theta)$  back into the expectation equation, we obtain the Policy Gradient Theorem (in its simple likelihood ratio form):

$$\nabla_{\theta} J(\theta) = E_{\tau \sim \pi_{\theta}} \left[ \left( \sum_{t=1}^H \nabla_{\theta} \log \pi_{\theta}(u_t | x_t) \right) C(\tau) \right]$$

**Monte Carlo Approximation** : The expectation can be approximated by sampling  $m$  trajectories  $(\tau^{(i)})$  from the current policy  $\pi_{\theta}$ :

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^H \nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} | x_t^{(i)}) C(\tau^{(i)})$$

This result forms the basis for the REINFORCE algorithm (Monte Carlo Policy Gradient).

## The REINFORCE algorithm

---

# REINFORCE: Derivation of the Update Rule

The final expression for the gradient expectation is:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m \sum_{t=1}^H \nabla_{\theta} \log \pi(u_t^{(i)} | x_t^{(i)}, \theta) C(\tau^{(i)})$$

**The REINFORCE Update Rule:**

$$\theta_{t+1} \leftarrow \theta_t - \alpha C_t \nabla \log \pi(U_t | X_t, \theta)$$

**Intuition of the Update:**

- The update moves  $\theta$  in the direction that **increases the log-probability** of the action  $U_t$  taken, proportional to the negative cost-to-go  $C_t$ .
- If  $C_t$  is negative, the probability of  $U_t$  is increased. If  $C_t$  is positive, the probability is decreased.
- This gradient tries to increase the probability of paths ( $\tau$ ) with negative cost-to-go and decrease the probability of paths with positive cost-to-go.

# REINFORCE: Key Characteristics

REINFORCE is a Monte Carlo (MC) Algorithm:

- It uses the **complete cost-to-go**  $C_t$  (the total discounted cost from time  $t$  until the end of the episode).
- Updates must be made in retrospect **after the episode is completed**.

Theoretical Property:

- REINFORCE is a **stochastic gradient method** whose expected update over an episode is in the opposite direction of the cost-to-go's gradient.
- It has good theoretical convergence properties (converging to a local optimum under standard conditions for decreasing  $\alpha$ ).

# REINFORCE (Monte Carlo Policy Gradient)

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic)

Input: A differentiable policy parameterization  $\pi(u|x, \theta)$

Parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta$  (e.g., to 0)

For each episode:

    Generate an episode  $X[0], U[0], L[0], \dots, L[T-1]$ , following  $\pi$

$C \leftarrow$  Compute cost-to-go

    For each step  $t = 0, 1, \dots, T - 1$ :

        # Policy Update

$\theta \leftarrow \theta + \alpha * C * \nabla \ln \pi(U_t | X_t, \theta)$

## Reducing the variance

---

# The Problem of High Variance

**The Trade-off:** REINFORCE is unbiased, but highly inefficient.

- Although the estimator is **unbiased** (its expectation is proportional to the true gradient), as a Monte Carlo method, it uses the complete cost-to-go  $G_t$  from a single rollout.
- The cost-to-go  $C_t$  is highly stochastic, meaning the sample gradient estimate is **very noisy** (high variance).
- This high variance leads to slow learning.

**Goal:** Find a way to reduce the variance of the gradient estimate without introducing bias.

# Exploiting Causality

The action at time  $t$  only affects the costs at times  $t' \geq t$ :

$$\begin{aligned}\nabla_{\theta} J(\theta) &\approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^H \nabla_{\theta} \log \pi(u_t^{(i)} | x_t^{(i)}, \theta) \right) C(\tau^{(i)}) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^H \nabla_{\theta} \log \pi(u_t^{(i)} | x_t^{(i)}, \theta) \right) \left( \sum_{k=0}^H I(x_k^{(i)}, u_k^{(i)}) \right) \\ &= \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^H \nabla_{\theta} \log \pi(u_t^{(i)} | x_t^{(i)}, \theta) \underbrace{\left( \sum_{k=0}^{t-1} I(x_k^{(i)}, u_k^{(i)}) + \sum_{k=t}^H I(x_k^{(i)}, u_k^{(i)}) \right)}_{\text{Independent of } u_t^{(i)}} \right)\end{aligned}$$

Removing the part of the cost that is independent of the current action can reduce the variance.

# Exploiting Causality

The improved policy gradient is:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{m} \sum_{i=1}^m \left( \sum_{t=1}^H \nabla_{\theta} \log \pi(u_t^{(i)} | x_t^{(i)}, \theta) \left( \sum_{k=t}^H I(x_k^{(i)}, u_k^{(i)}) \right) \right)$$

Still unbiased, however, still high variance.

## Variance Reduction: Introducing the Baseline

The primary technique to reduce variance without introducing bias is the introduction of an arbitrary **Baseline**  $b(x)$ .

**REINFORCE Update with Baseline:**

$$\theta_{t+1} \leftarrow \theta_t - \alpha (C_t - b(X_t)) \nabla \ln \pi(U_t | X_t, \theta)$$

The term  $(C_t - b(X_t))$  is now the "effective cost-to-go" used to scale the gradient step.

## Baseline Validity: The Unbiased Property

The generalized PGT remains valid **if and only if** the baseline  $b(x)$  does not vary with the action  $u$ .

The change in the gradient expectation due to the baseline is zero:

$$\begin{aligned} E_{\pi} [\nabla \ln P(\tau|\theta) b(x)] &= \sum_{\tau} P(\tau|\theta) \nabla \ln P(\tau|\theta) b(x) \\ &= \sum_{\tau} \nabla P(\tau|\theta) b(x) = \nabla \left( \sum_{\tau} P(\tau|\theta) b(x) \right) = \\ &= b(x) \nabla \left( \sum_{\tau} P(\tau|\theta) \right) = 0 \end{aligned}$$

Since the sum of probabilities of all trajectories must equal one ( $\sum_{\tau} P(\tau|\theta) = 1$ ), its gradient is null ( $\nabla(\sum_{\tau} P(\tau|\theta)) = 0$ ).

The expected value of the subtracted baseline term is zero. Thus, the gradient estimate remains **unbiased**.

## Baseline Choice: The State-Value Function

**Natural Choice:** The most effective choice for the baseline is an estimate of the **state-value function**,  $V(X_t, w)$ .

**The Advantage Function ( $A(x, u)$ ):** When  $b(X_t) = V(X_t, w)$ , the term scaling the gradient is an estimate of the Advantage function:

$$\hat{A}(X_t, U_t) = C_t - V(X_t, w)$$

**Intuition:** The policy probability  $\pi(U_t|X_t, \theta)$  is increased only if the action  $U_t$  yields a cost-to-go  $C_t$  that is **better** than the expected cost-to-go  $V(X_t)$ .

- This relative comparison allows the parameter  $\theta$  to differentiate between better and worse actions in a state where all returns might be generally high or generally low.
- Empirically, adding the state-value function as a baseline can make REINFORCE learn much faster.

# REINFORCE with Baseline (Pseudocode)

```
Input: Policy pi_theta (e.g., Neural Network)
Initialize policy parameters theta, baseline b

Loop forever (for each iteration/batch):
    1-Collect a batch of m trajectories {tau^(i)} by running
        current policy pi_theta
    2-For each time step calculate the total discounted
        cost-to-go C_t and the advantage A_t = C_t - b(X_t)
    3-Re-fit the baseline by minimizing ||b(X_t) - C_t||
    4-Update the policy using a policy-gradient estimate
        SUM_i=1^m SUM_t=1^H [log(pi(u_t^(i) | x_t^(i))) * A_t^(i)]
```

# **Actor-Critic Policy-Gradient Methods**

---

# From REINFORCE to Actor-Critic Methods

REINFORCE with Baseline still relies on the **full Monte Carlo cost-to-go** ( $C_t$ ), meaning updates are delayed until the end of the episode.

**Actor-Critic Goal:** Achieve online, incremental learning and further reduce variance by introducing **bootstrapping**.

- **Actor** ( $\pi_\theta$ ): The policy that selects actions.
- **Critic** ( $V_w(x)$ ): The learned value function that evaluates the policy's actions (assigns credit).

**The Trade-off:**

- Using bootstrapping (updating the value estimate from estimated values of subsequent states) reduces variance and accelerates learning.
- However, bootstrapping introduces **bias** and an asymptotic dependence on the quality of the function approximation.

# One-Step Actor-Critic Methods

We replace the full Monte Carlo cost-to-go  $C_t$  with the **one-step cost-to-go**  $C_{t:t+1}$  to enable bootstrapping and online updates.

The one-step Advantage estimate is the Temporal Difference (TD) error ( $\delta_t$ ):

$$\delta_t = L_{t+1} + \gamma V(X_{t+1}, w) - V(X_t, w)$$

**Policy Update (Actor):**

$$\theta_{t+1} \leftarrow \theta_t - \alpha_\theta \delta_t \nabla \ln \pi(U_t | X_t, \theta)$$

**Value Function Update (Critic):** The critic is typically updated using TD(0):

$$w_{t+1} \leftarrow w_t + \alpha_w \delta_t \nabla V(X_t, w)$$

**Benefit:** This provides a fully online and incremental algorithm, processing states, actions, and costs as they occur.

# One-step Actor-Critic (Pseudocode)

```
One-step Actor-Critic (episodic)
```

For each episode:

    Initialize X

    I <- 1

    Loop while X is not terminal:

        U <-  $\pi(\cdot | X, \theta)$

        Take action U, observe X', L

        # Calculate TD Error (Delta)

$\Delta \leftarrow L + \gamma * V_{\hat{w}}(X', w) - V_{\hat{w}}(X, w)$

        # Critic Update

$w \leftarrow w + \alpha_w * \Delta * \nabla V_{\hat{w}}(X, w)$

        # Actor Update

$\theta \leftarrow \theta - \alpha_t * I * \Delta * \nabla \ln \pi(U|X, \theta)$

$I \leftarrow \gamma * I$

$X \leftarrow X'$

\*Note: I is the eligibility trace factor that handles discounting.

# Advanced PG: Sophisticated Advantage Estimation

One-step TD error is low variance but high bias. MC cost-to-go ( $C_t$ ) is high variance but low bias. Advanced methods combine these:

## 1. Generalized Advantage Estimation (GAE):

- GAE uses an **exponentially weighted average** of  $n$ -step TD errors, similar to  $\text{TD}(\lambda)$ .
- GAE balances the bias-variance trade-off by combining returns of different lengths, providing a robust Advantage estimate.

## 2. Asynchronous Advantage Actor Critic (A3C):

Uses  $n$ -step Advantage estimation (bootstrapping).

## 3. TRPO and PPO:

Introduce constraints (Trust Regions) to stabilize large gradient steps, which is crucial for complex tasks like robotic locomotion.

## Summary of Key Concepts

---

1. **Policy Optimization:** Directly optimizes performance  $J(\theta)$  with respect to policy parameters  $\theta$ . Compatible with continuous actions.
2. **PGT:** Provides the theoretical foundation, guaranteeing that the gradient can be estimated from samples without needing to model state distribution derivatives.
3. **REINFORCE:** A Monte Carlo (MC) method derived from the PGT. It uses the full cost-to-go  $C_t$  and is **unbiased** but suffers from high variance.
4. **Baseline:** Subtracting the value function  $V^\pi(x)$  reduces variance without introducing bias (unbiased). This yields the Advantage  $C_t - V(X_t)$ .
5. **Actor-Critic:** Uses the learned value function (Critic) for **bootstrapping** (e.g., using TD error  $\delta_t$ ) to accelerate learning and enable online updates, at the cost of introducing bias.

## Application in Robotics

PG methods are highly relevant to robotics due to their ability to handle continuous, high-dimensional action spaces.

- **Continuous Control:** When controlling robotic manipulators, legged robots, or autonomous vehicles, the actions (torques, velocities) are continuous, where value-based methods struggle.
- **Locomotion:** Advanced PG methods like TRPO and PPO, often combined with GAE, have been essential in learning complex locomotion behaviors for simulated and real-world robots.
- **Historical Success:** Policy gradients were crucial for complex applications like acrobatic helicopter autopilots.

Policy Gradient methods provide a framework for learning sophisticated, continuous control policies that can often be simpler to represent than the corresponding value functions.

# Conclusion and Future Outlook

- Policy gradient methods offer a significantly different set of strengths and weaknesses compared to action-value methods (Q-Learning).
- While historically less studied than Q-Learning, they are now a subject of great excitement and ongoing research in Deep RL.

**Questions?**