

## ROBOT DYNAMIC REVIEW

Kinematics describe the motion without considering the forces involved.

HOMOGENEOUS TRANSFORMATION MATRIX represents the pose of a frame wrt another.

$$T_1^0 = \begin{bmatrix} R_1^0 & p \\ 0_{3 \times 3} & 1 \end{bmatrix} \quad p \text{ is the } 3 \times 1 \text{ position vector}$$

FORWARD KINEMATICS calculates the pose of the end-effector given the joint coordinates  $q$  by composing the transformation matrices of each link.

$$T_n^n(q) = T_1^0(q) T_2^1(q_2) + \dots + T_{n-1}^{n-1}(q_{n-1})$$

DIFFERENTIAL KINEMATICS relates joint velocities to end-effector linear and angular velocities. JACOBIAN MATRIX maps joint velocities ( $\dot{q}$ ) to the end-effector's Cartesian velocity ( $v_e$ ). The linear velocity  $v_e$  is composed of linear velocity  $\dot{p}$  and angular velocity  $\omega$ :

$$v_e = \begin{bmatrix} \dot{p} \\ \omega \end{bmatrix} = J(q) \dot{q}$$

STATICS relates forces and torques at the end-effector to the torques required at the joints. The joint torque ( $\tau$ ) required to balance an external wrench  $w = [f^T, m^T]^T$  applied at the end-effector are calculated using the transpose of the Jacobian:

$$\tau = J^T(q) w$$

DYNAMICS describes relationship between forces acting on the robot and its resulting motion.

$$H(q)\ddot{q} + C(q, \dot{q})\dot{q} + g(q) = \tau + J^T(q)w$$

## MOTION CONTROL

Strategies to compute required joint torques ( $\tau$ ) to make the robot follow a desired trajectory.

JOINT SPACE CONTROL have as objective to make the robot's joint configuration  $q(t)$  track a reference joint trajectory  $q^{ref}(t)$ . The tracking error is defined as  $e = q^{ref} - q$

- PD + GRAVITY COMPENSATION:  $\tau = K_p e + K_d \dot{e} + g(q)$  guarantees zero steady-error when desired velocity is zero.
- INVERSE DYNAMIC CONTROL:  $\tau = M(q)(\ddot{q}^{ref} + K_p e + K_d \dot{e}) + h(q, \dot{q})$  make the system linear  $\rightarrow \ddot{e} + K_d \dot{e} + K_p e = 0$

CARTESIAN SPACE CONTROL make end-effector pose  $x(t)$  track a reference Cartesian trajectory  $x^{ref}(t)$ .

- OPERATIONAL SPACE DYNAMICS:  $\tau = J^T F$  calculates a cartesian force  $F$  needed to drive the end-effector. The force is then mapped back to joint torques.
- CONTROL LAW WITH REDUNDANCY:  $\tau = J^T F + (I - J^T \bar{J}^T) \tau_0$  the null space of Jacobian can be used to perform secondary tasks. Control law is extended to null-space torque.

# 1. QP-BASED REACTIVE CONTROL

## JOINT SPACE CONTROL [standard PD control]

$$H(q)\ddot{q} \rightarrow h(q, \dot{q}) = \tau \Rightarrow \begin{cases} \tau^d = M\ddot{q}^d + h \\ \ddot{q}^d = \ddot{q}^r + K_d(\dot{q}^r - \dot{q}) + K_p(q^r - q) \end{cases}$$

↳ Dynamics equation

→ Trasformare il problema di come controllare un robot in un problema di ottimizzazione (risolvibile in real-time)

Optimization problem that allow us to come at the same solution:

$$(\tau^*, \ddot{q}^*) = \underset{\tau, \ddot{q}}{\operatorname{arg\,min}} \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2 = \frac{1}{2} \ddot{q}^T \ddot{q} - \ddot{q}^{d T} \ddot{q} + \frac{1}{2} \ddot{q}^{d T} \ddot{q} \quad \text{constant}$$

subject to  $M\ddot{q} + h = \tau$  L'accelerazione che troviamo deve rispettare la fisica del robot

$$\text{SOLUTION} \quad \begin{cases} \tau^* = \tau^d = M\ddot{q}^d + h \\ \ddot{q}^* = \ddot{q}^d \end{cases}$$

The advantage of this version of the problem is that we define:

- physical bound:  $\tau^{\min} \leq \tau \leq \tau^{\max}$  coppia massima
  - electric current limits:  $i^{\min} \leq K_i \tau \leq i^{\max}$
  - joint velocity limits:  $\dot{q}^{\min} \leq \dot{q} + \Delta t \ddot{q} \leq \dot{q}^{\max}$
- Andiamo a limitare l'accelerazione

At the end, we have to solve:

$$\underset{\tau, \ddot{q}}{\operatorname{min}} \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2 \text{ subject to } M\ddot{q} + h = \tau + \text{constraints defined above}$$

↳ Use a solver to solve the optimization problem efficiently

## VELOCITY CONSTRAINTS:

Since we can regulate only  $\ddot{q}$ , to manage velocity constraints we assume that  $\ddot{q} = \text{const}$  during time step  $\Delta t$ :  $\dot{q}(t + \Delta t) = \dot{q}(t) + \Delta t \ddot{q}$   
so I rely on the next value of the velocity. (Apisco sull'accelerazione)

## JOINT LIMITS:

→ Applicare limiti ai giunti: potrebbe richiedere accelerazioni eurmi

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

↳ In practice doesn't work. (In EXTRA material there's the proof)

↳ This problem is not solvable with reactive control methods

## TAXONOMY OF convex optimization problems

- LEAST SQUARE PROGRAM (LSP) →  $Ax \leq b$ ,  $Ax = b$

$$\rightarrow \|Ax - b\|^2$$

- QUADRATIC PROGRAMS (QP) → linear constraints

$$\rightarrow \text{convex quadratic cost } \frac{1}{2} x^T H x + g^T x, H \geq 0$$

can be solved efficiently ( $\leq 1 \mu s$ ) by solvers.

## LEAST SQUARE PROGRAMS

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

$$\nabla f = A^T A x - A^T b = \emptyset \rightarrow 2 \text{ cases to find the solution}$$

(1) IF  $A^T A$  is invertible  $x^* = \underline{(A^T A)^{-1} A^T b}$   
left-pseudoinverse of  $A$

(2) IF  $A^T A$  is NOT invertible, but  $A A^T$  is invertible:

$$x^* = \underbrace{A^T (A A^T)^{-1} b}_{\text{Right-pseudoinverse}} \leftarrow \begin{cases} Ax^* = A A^T (A A^T)^{-1} b \\ A^T A A^T (A A^T)^{-1} b - A^T b = A^T b - A^T b = \emptyset \end{cases} \rightarrow \begin{cases} f(x) = \emptyset \\ \nabla f = \emptyset \end{cases}$$

### SUMMARY

OBJECTIVE: minimize  $x^* = \arg \min_x \|Ax - b\|^2 \Leftrightarrow A^* = A^T b$

SOLUTION:  $A^*$  is computed using Singular Value Decomposition (SVD) (numerical method).

From Joint to cartesian space (task space)

We know that  $\ddot{x} = J\ddot{q} + \dot{J}\dot{q}$      $\ddot{x}^d = \ddot{x}^r + K_d(\dot{x}^r - \dot{x}) + K_p(x^r - x)$  [  $\ddot{x}$  → acceleration end-effector ]  
 We can compute  $\ddot{q}^d$  as:  $\begin{cases} \ddot{q}^d = J^+ (\ddot{x}^d - \dot{J}\dot{q}) \\ \tau^d = M\ddot{q}^d + h \end{cases} \Rightarrow \text{PID}$

If we treat the problem as an optimization problem:

$$(\tau^*, \ddot{q}^*) = \arg \min_{\ddot{q}, \tau} \|J\ddot{q} + \dot{J}\dot{q} - \ddot{x}^d\|^2 \equiv \|\ddot{x} - \ddot{x}^d\|^2 \quad \text{s.t. } M\ddot{q} + h = \tau$$

↳ It is still a least square program  $\Rightarrow \ddot{J}\ddot{q} + \dot{J}\dot{q} - \ddot{x}^d$

We can add same constraints as seen before.

### Task Models

- Generalize concept of e-e control
- Task = control objective
- Describe task with function  $e(\cdot)$  to minimize
- Assume that  $e(\cdot)$  describes error btw real and reference trajectory:  

$$e(x, u, t) = \underbrace{y(x, u)}_{\substack{\text{Actual state of task} \\ \text{Desired trajectory}}} - \underbrace{y^r(t)}_{\substack{\text{Robot state} \\ \text{const}}} \quad x \triangleq (q, \dot{q}) \quad u \triangleq \tau$$
- We cannot directly minimize  $e(\cdot)$  if it depends on  $q$  and  $\dot{q}$ , which are not decision variables of the LSP.
- IDEA: impose dynamics of  $e$  (e.g.  $\dot{e}$  or  $\ddot{e}$ ) such that:
  - (I)  $\lim_{t \rightarrow \infty} e(t) = 0$
  - (II) Dyn. of  $e$  is affine function of  $(\ddot{q}, \tau)$  ( $Ax + b$ )

## CASE #1 VELOCITY TASK FUNCTION

L'errore dipende dalla velocità dei giunti

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

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

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

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

$$\dot{e} = -K_e e^* \quad K \geq 0 \Rightarrow \lim_{t \rightarrow \infty} e(t) = 0 \quad * \text{Since } e = \text{const}, \forall K \quad -K_e = \text{const}$$

$$J\ddot{q} - \ddot{y}^r = -K_e$$

$$\underbrace{J\ddot{q}}_A - \underbrace{\ddot{y}^r}_B - K_e = 0 \Rightarrow A\ddot{q} - b = 0 \Rightarrow \|A\ddot{q} - b\|^2$$

X risolvere usiamo i minimi quadrati  
e risolviamo per  $\ddot{q}$

## CASE #2 : CONFIGURATION TASK FUNCTION

L'errore deriva dalla posizione dei giunti

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

$$\dot{e} = \dot{y} - \dot{y}^r = \frac{\partial y}{\partial q} \frac{dq}{dt} - \dot{y}^r = J\dot{q} - \dot{y}^r \text{ Not enough, I need it in function of } \ddot{q}$$

$$\ddot{e} = J\ddot{q} + J\dot{q} - \ddot{y}^r \text{ Linear function of } \ddot{q}$$

$$\text{Let's impose linear dynamics: } \ddot{e} = -K_p e - K_d \dot{e} \quad K_p \geq 0, K_d \geq 0$$

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

$$\underbrace{J\ddot{q}}_A - \underbrace{\ddot{y}^r - J\dot{q} - K_p e - K_d \dot{e}}_B \Rightarrow \text{Minimize } \|A\ddot{q} - b\|^2$$

The main difference between configuration and velocity task function is that in the first there are 2 feedback gains ( $K_d, K_p$ ), while in the second there is only 1.

## CASE #3 : CONTROL-INPUT TASK FUNCTION

Control inputs  $u \equiv \tau$  are part of the decision variables.

Therefore  $e(x, u, t) = e(u, t)$  must be affine:  $e(u, t) = A_u u(t) + b(t)$

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

Example:  $\|u\|^2$  or  $\|u - u^*\|^2$

## SUMMARY

### TASK FUNCTIONS

Affine function of  $u \equiv \tau$

Nonlinear function of  $x$ :

for function of  $\dot{q} \Rightarrow \dot{e}$

for function of  $q \Rightarrow \ddot{e}$

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

$$g(z) = \underbrace{[Ax \quad Au]}_A \underbrace{\begin{bmatrix} \ddot{q} \\ \tau \end{bmatrix}}_z - b \Rightarrow \min \|g(z)\|^2$$

## Task-Space Inverse Dynamic [TSID]

$$\underset{\substack{z = (\ddot{q}, \tau) \\ 2}}{\text{minimize}} \frac{1}{2} \|Az - b\|^2 \text{ subject to } [M \mid -I] z = -h$$

Modo compatto per scrivere  $M\ddot{q} - \tau = -h$

$$z = \begin{bmatrix} \ddot{q} \\ \tau \end{bmatrix}$$

For example:

- Joint-space control:  $A = [I, \theta]$ ,  $b = \ddot{q}^d \Rightarrow \min_{z=(\ddot{q}, \tau)} \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2$   
 $Az - b = \ddot{q} - \ddot{q}^d$
- End-effector control:  $A = [J, \theta]$ ,  $b = \ddot{x}^d - J\dot{q} \Rightarrow \min_{z=(\ddot{q}, \tau)} \frac{1}{2} \|\ddot{x} - \ddot{x}^d\|^2$   
 $Az - b = J\ddot{q} + \dot{J}\dot{q} - \ddot{x}^d$
- Under-actuated systems (less actuators than DoFs):  
 $\min_{z=(\ddot{q}, \tau)} \frac{1}{2} \|Az - b\|^2$  s.t.  $[M \mid -I] z = h$  vector that contain 0  
 $S^T$  in passive joints

## TSID FOR RIGID CONTACTS (constraints for motion)

- $C(q) = \text{const}$  contact points don't move
  - $J\dot{q} = \emptyset$  contact point velocities are null
  - $J\ddot{q} + \dot{J}\dot{q} = \emptyset$  contact point accelerations are null
- ↳ se il robot tocca un oggetto, il punto di contatto non può muoversi

**CONTACT FORCE**  $F$  is a decision variable in QP:

$$\min_{z=(\ddot{q}, \tau, F)} \|Az - b\|^2 \quad \text{s.t.} \quad \left\{ \begin{array}{l} M\ddot{q} + h = S^T \tau + J^T F \\ J\ddot{q} + \dot{J}\dot{q} = \emptyset \end{array} \right.$$

COPPIE DI MOTORI

$J$  is the Jacobian of  $C(q)$   $J \triangleq \partial C(q)/\partial q$   
 vincolo + importante, poiché il nostro controller (QP) comanda  $\dot{q}$  e  $\tau$ .

Quando il robot entra in contatto non solo deve eseguire il task ma deve anche "obbedire" a nuovi vincoli imposti dal contatto con un oggetto

Single constraint in matrix form →

$$\begin{bmatrix} M & -J^T & -S^T \\ J & \Theta & \emptyset \end{bmatrix} \begin{bmatrix} \ddot{q} \\ F \\ \tau \end{bmatrix} = \begin{bmatrix} -h \\ -J\ddot{q} \end{bmatrix}$$

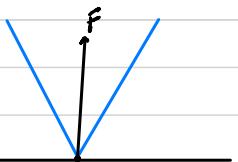
SOFT CONTACT  $\Rightarrow F = \text{cost}$   $F$  è sconosciuta, il sistema la calcola per mantenere i vincoli

HARD CONTACT  $\Rightarrow F$  depends on how much force I'm applying

$\hookrightarrow F$  fa parte di una variabile di costo. Ad esempio possiamo applicare una forza specifica minimizzando  $\|F - F^d\|^2$

## FRICITION CONE CONSTRAINTS

La forza deve essere fisicamente realistica



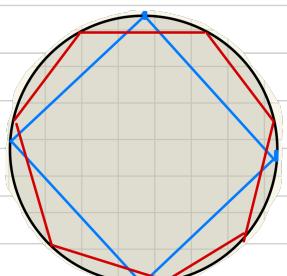
$$\|F_T\|^2 \leq \gamma F_N^2$$

tangential force  
(Atrito)

normal force

$\gamma$ : friction coefficient

friction force must be stay in a "friction cone"



Approximate quadratic cone with linear pyramidal cone:

$$BF \leq \emptyset$$

Include this in TSID's QP.

## Multy-Task Control

For complex redundant robots typically one has multiple control objectives, e.g. N task, each defined by a task function  $g_i(z) \triangleq \|Az - b\|^2$   $i = 1, \dots, N$

TSID QP:  $\min_{z=(\ddot{q}, \tau, F)} \sum_{i=1}^N w_i g_i(z) = \|Az - b\|^2$  s.t.  $\begin{bmatrix} M & -J^T & -S^T \\ J & \emptyset & \emptyset \end{bmatrix} \begin{bmatrix} \ddot{q} \\ F \\ \tau \end{bmatrix} = \begin{bmatrix} -h \\ -J\dot{q} \\ \text{Contact} \end{bmatrix}$

$\hookrightarrow$  Minimizzare la somma pesata di tutti i task

HARD CONSTRAINTS (Physics Laws)

$w_i \in \mathbb{R}^+$   $\triangleq$  user-defined weight. Large  $w_i \Rightarrow$  high importance to task  $i$

## Practical Implementation

Most industrial robots only allow for joint pos/vel commands. They have high joint friction  $\Rightarrow$  hard to achieve torque control.

hard to model and to compensate

Joint friction is mainly a consequence of high gear ratios.

What can I do if I cannot send torque commands?

(1) Solve QP as usual to compute  $z^* = (\ddot{q}^*, \tau^*)$

(2) Integrate  $\ddot{q}^* \Rightarrow \dot{q}^d = \dot{q}^* + \Delta t \cdot \ddot{q}^*$

(3) Send joint velocity commands  $\dot{q}^d$

OPTIMIZE  $\Rightarrow$  DISCRETE  
Continuous Time

DISCRETE  $\Rightarrow$  OPTIMIZE  
Discrete Time

(2)  
Hamilton-Jacobi-Bellman (HJB)

(1)  
Dynamic Programming (DP)

Global

Local

(3)  
Pontryagin Maximum Principle (PMP)  
Calculus of variations  
Indirect method

(3)  
Direct method

NOTEBOOK NUMERICAL OPTIMIZATION

## 2. OPTIMAL CONTROL

- No need for reference trajectory
- Objective clearly specified through cost function
- Can handle constraints

Try to find best state  $x(t)$  and control  $u(t)$  trajectories in a time interval  $t \in [0, T]$

### O.C. PROBLEM

$$\left\{ \begin{array}{l} \text{minimize}_{x(\cdot), u(\cdot)} \int_0^T e(x(t), u(t), t) dt + e_f(x(T)) \\ \text{subject To } \dot{x}(t) = f(x(t), u(t), t) + g(x(t), u(t), t) \leq 0 \quad \forall t \in [0, T] \\ x(0) = x_0 \quad \text{Initial condition (not mandatory)} \\ g(x(t), u(t), t) \leq 0 \quad \text{Path constraints} \end{array} \right.$$

Fixed length of the Problem

RUNNING COST

TERMINAL COST (depends only on the terminal state)

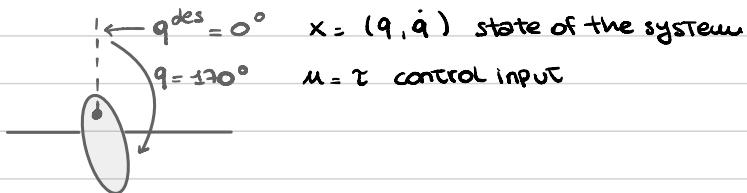
IL robot deve seguire le leggi della fisica

DYNAMICS

cost function

- $x(\cdot), u(\cdot)$  are trajectories  $\Rightarrow$  infinite-dimensional objects
  - constraints are infinitely many
  - $\neq$  optimization problem (at time instant between  $[0, T]$ )
- Stiamo cercando intere funzioni, non numeri

### EXAMPLE Pendulum



$$\min_{x(\cdot), u(\cdot)} \int_0^T (w_q q(t)^2 + w_{\dot{q}} \dot{q}(t)^2 + w_{\tau} \tau(t)^2) dt + 10^2 q(T)^2$$

How much fast pendulum reach desired  $q=0$

Torque and acceleration are linear dependant so we can optimize one

subject to

$$\underbrace{I \ddot{q} = \tau + mg \sin(q)}_{\text{System dynamic of a pendulum}} \Rightarrow \begin{bmatrix} \dot{q} \\ \ddot{q} \end{bmatrix} = \begin{bmatrix} \dot{q} \\ I^{-1} (\tau + mg \sin(q)) \end{bmatrix} \quad \forall t$$

$$\begin{aligned} |\tau(t)| &\leq \tau^{\max} & \forall t \in [0, T] \\ q^{\min} &\leq q(t) \leq q^{\max} & \forall t \in [0, T] \\ q(T) &= \emptyset \end{aligned}$$

- NOTE**
- $10^2 q(T)$  terminal cost must be in function only of the terminal state
  - $q(T) = \emptyset$  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. Same for joint velocities, we can define  $\dot{q}(T)=0$ .
  - $w_q q(t)^2$  It's possible to use the absolute value instead of square. With square you are more tolerant with small errors and penalize more big errors.

# Non Linear Programming [2.2 Notebook]

Non linear programming has as objective to minimize a non linear cost function  $f(x)$ , that have equality and disequality constraints.

For an UNCONSTRAINED problem, to verify that  $x^*$  is a LOCAL MINIMA, we have to check:

$$(1) \nabla f(x^*) = 0$$

$$(2) \nabla^2 f(x^*) > 0 \quad \text{Hessian matrix must be positive. Verify that } x^* \text{ is a minimum.}$$

→ To find a solution we use iterative numerical methods:

- Gradient descent  $x_{k+1} = x_k + \gamma \nabla f(x_k)$  converge slowly

- NEWTON METHOD  $x_{k+1} = x_k - (\nabla^2 f)^{-1} \nabla f$  use 2° order informations, faster convergence.

This method is unstable, so it's necessary two improvements

- Hessian regularization:  $H + \lambda I$  guarantee that direction is always descent

- Line search: we looking for an optimal step length  $\alpha \in [0, 1]$ , to prevent divergences.

If we have a CONSTRAINED problem we have to find  $x^*$  that satisfy the constraints.

KKT CONDITIONS generalize gradient concept and represent necessary conditions for optimality.

- LAGRANGIAN FUNCTION: combines original cost function with constraints, weighted by coefficients called Lagrange multipliers ( $\lambda$  for eq.,  $\mu$  for dis.).

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

$g(x) = \text{Equality constraints}$   
 $h(x) = \text{Inequality constraints}$

- KKT CONDITIONS: for a point  $x^*$ , to be a minimum, must be exist a set of  $\lambda, \mu$  such that:

1) STATIONARITY:  $\nabla_x L(x, \lambda, \mu) = 0 \rightarrow$  we can't move in any direction without violate any of the constraints.

2) PRIMAL ADMISSIBILITY:  $x^*$  must respect all constraints  $\begin{cases} g(x^*) = 0 & \text{equality constraints} \\ h(x^*) \leq 0 & \text{inequality constraints} \end{cases}$

3) DUAL ADMISSIBILITY:  $\mu \geq 0$  this is the "cost" of the constraint. If  $\mu = 0$  that constraint doesn't limit us.

4) COMPLEMENTARY SLACKNESS:  $\mu_i h_i(x^*) = 0 \forall i \rightarrow x^*$  can be only on the frontier of the constraint or inside it. If  $x^*$  is inside of admissible region of the constraint, it is not affected by it.

## EQUALITY CONSTRAINTS $g(x^*) = 0$

Define feasible cone:  $F(x) = \{d \mid \nabla g(x)^T d = 0\}$  Insieme di tutte le direzioni  $d$  in cui posso muovermi partendo da  $x$ .  $\nabla g(x)^T d = 0 \Rightarrow d$  deve essere perpendicolare al gradiente del vincolo  $\nabla g(x)$

Set of descent directions:  $\nabla_x(x) = \{d \mid d^T \nabla f(x) < 0\}$  Insieme delle direzioni che riducono  $f(x)$

Local minimum  $x \rightarrow \nabla f(x) \wedge F(x) = 0$  In un minimo non esiste nessuna direzione che sia contrapposizionalmente ammessa e di discesa. Se così non fosse,  $d$  mi porterebbe in una direzione migliore, perciò  $x^*$  non sarebbe un minimo.

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

↑ I 2 gradienti sono paralleli

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

↓  
KKT conditions:  $\nabla_L(x, \lambda) = 0 \Leftrightarrow \begin{cases} \nabla f(x) + \nabla g(x) \lambda = 0 & \text{Derivation wrt } x \\ g(x) = 0 & \text{Derivation wrt } \lambda \end{cases}$

## INEQUALITY CONSTRAINTS

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

→ Active means that we are on boundaries, so that constraint limits us. Otherwise the constraint is inactive and we can move in any direction inside the region.

Feasible cone  $F(x) = \{d \mid \nabla g(x)^T d = 0, \nabla h_i(x)^T d \leq 0, i \in A(x)\}$  Considering only active constraints, by moving in any direction of they can only decrease or remain the same

$$\text{KKT conditions} \left\{ \begin{array}{l} g(x) = 0, h(x) \leq 0 \\ d^T \nabla f(x) \geq 0 \quad \forall d \in F(x) \end{array} \right. \Leftrightarrow \nabla f(x) + \nabla g(x) \lambda + \nabla h_A(x) \gamma_A = 0 \\ \gamma_A \geq 0$$

$$\Leftrightarrow \left\{ \begin{array}{l} \nabla f(x) + \nabla g(x) \lambda + \nabla h(x) \gamma = 0 \\ \gamma \geq 0 \\ \gamma h(x) = 0 \end{array} \right. \text{complementary condition} \Rightarrow \begin{array}{l} \text{if } h(x) = 0 \rightarrow \gamma_i \text{ can be positive} \\ \text{if } h(x) < 0 \rightarrow \gamma_i = 0 \end{array}$$

## SOLVING KKT CONDITIONS

$$\min f(x)$$

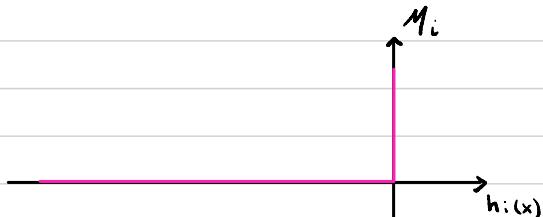
$$\text{s.t. } g(x) = 0 \quad \Rightarrow \text{KKT system} \\ h(x) \leq 0$$

$$\left\{ \begin{array}{l} \nabla_x L(x, \lambda, \gamma) = 0 \\ g(x) = 0 \\ h(x) \leq 0 \\ \gamma \geq 0 \\ \gamma^T h(x) = 0 \end{array} \right. \quad \begin{array}{l} \text{for equality constraints} \\ \text{extra for inequality constraints} \end{array}$$

$$L(x, \lambda, \gamma) = f(x) + \lambda^T g(x) + \gamma^T h(x)$$

## Direct Methods

Newton method cannot be directly used to solve inequality constraints because solution space is not smooth.



KEY IDEA: discretize OCP  $\Rightarrow$  use NLP solver to find local optimum

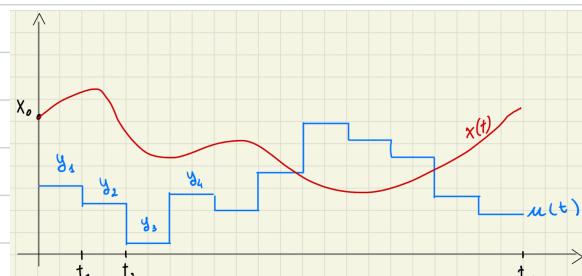
DISCRETIZATION: 1) parametrize state/control trajectories  
2) enforce constraints on a time grid to  $t_1 < \dots < t_N$

### ① SINGLE SHOOTING (sequential):

- discretize only  $u(t)$
  - compute  $x(t)$  integrating dynamics
  - No need to have dynamic constraints
- Discretizzo solo il controllo  $u(t)$  e ne derivo  $x(t)$  di conseguenza

### ② Discretize $u(t) = y_i \quad \forall t \in [t_i, t_{i+1}]$

optimization variable



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

$$\begin{cases} \dot{x}(t) = f(x, u, t) \\ x(0) = x_0 \text{ fixed and known} \end{cases}$$

c. NLP

$$\min_y \int_0^{t_f} \begin{array}{l} \text{Running cost} \\ l(x(t, y), u(t, y), t) dt \end{array} + \begin{array}{l} \text{terminal cost} \\ l_f(x(t_f, y)) \end{array}$$

La funzione di costo viene calcolata integrando i costi lungo la traiettoria  $x(t)$  ottenuta.

s.t.  $\begin{array}{l} g(x(t_i, y), u(t_i, y), t_i) \leq 0 \\ i = 0, \dots, N \end{array}$

Tutti i vincoli vengono imposti solo nei punti della griglia temporale

Discretize path constraints

Running cost integral can be computed when integrating dynamics:

$$\begin{cases} c(t) = \int_0^t l(\cdot, \cdot) dt \\ \dot{c}(t) = l(\cdot, \cdot) \\ c(0) = 0 \end{cases} \Rightarrow \begin{array}{l} \bar{x} = (x, c) \\ \dot{\bar{x}} = \begin{bmatrix} f(x, u, t) \\ l(x, u) \end{bmatrix} \end{array}$$

### NUMERICAL INTEGRATION

given ode:  $\dot{x}(t) = f(x(t), t)$ ,  $x(0) = x_0$

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

If first derivatives of  $f$  are bounded  $\Rightarrow$  unique solution

### EXPLICIT EULER

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

$$\hookrightarrow h f(x, t) \approx x(t+h) - x(t) \Rightarrow x(t+h) = x(t) + h f(x, t) \quad \text{small } h \Rightarrow \uparrow \text{acc.} \uparrow \text{comp. cost}$$

### RUNGE-KUTTA METHODS

$x_{n+1} = x_n + h \sum_i b_i k_i$  approximate average  $\dot{x}$  with weighted average of  $k_i$ , with weights  $b_i$ , with  $\sum b_i = 1$ ,  $b_i \geq 0$

$k_i = f\left(x_n + h \sum_j a_{ij} k_j, t_n + c_i h\right)$  approximate value of  $\dot{x}$  at time  $t_n + c_i h$ , with  $c_i \in [0, 1]$

- $a_{ij}$  and  $c_i$  say when and where evaluate dynamics
- $q$  is the order of the method  $\rightarrow$  how many evaluations of dynamics to compute.
- for  $q > s$  exists many versions of same order  $P$  (consistency).

### SUMMARY

- 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

Discretizzo sia  $x(t)$  che  $u(t)$

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

$x(t) = s_i \quad \forall t \in [t_i; t_{i+1}] \rightarrow$  one polynomial for each time interval

Sia lo stato che il controllo vengono parametrizzati. Vengono rappresentati da polinomi a tratti, le variabili di ottimizzazione diventano i coefficienti di questi polinomi (o i loro valori  $s_i, y_i$  nei punti di griglia)

b. Replace Dynamics  $\dot{x} = f(x, u)$  with: APPROXIMATION

$$\frac{s_{i+1} - s_i}{t_{i+1} - t_i} - f(s_i, y_i) = 0 \quad \text{rewriting} \rightarrow c_i(s_i, s_{i+1}, y_i) = 0$$

Approximation of dynamics

↳ la dinamica diventa un VINCOLO

Approximation of dynamics has to be equal to real dynamics to have best approximation possible

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

d. NLP (sparse)

$$\min_{s, y} \sum_{i=0}^{N-1} l_i(s_i, y_i) + l_F(s_N) \quad \text{Minimize sum of approximated costs at each interval + terminal cost}$$

s.t.  $s_0 - x_0 = 0$  (initial conditions)

↳ sequenza dei punti deve partire dalla condizione iniziale

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

$g_i(s_i, y_i) \leq 0 \quad i=0, \dots, N-1 \leftarrow$  discretized path constraints

Check constraints only at discrete time inst.

KKT matrix is SPARSE  $\rightarrow$  most of the entry of the matrix are zero. This is useful because solvers can check sparsity and the presence of zeros lead to a more efficient computation.

$$H_{ij} = \frac{\partial}{\partial s_i} \left( \frac{\partial L}{\partial s_j} \right) = 0 \quad \forall i \neq \{j, j+1, j-1\} \quad H = \begin{bmatrix} & & & & \text{zeros} \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \ddots \end{bmatrix}$$

## 3 MULTIPLE SHOOTING

→ hybrid between single shooting and collocation. Tries to solve instability problems in single shooting.

a. Discretize control  $u(t)$  on coarse grid to  $t_1 < \dots < t_N$   $u(t) = y_i \quad \forall t \in [t_i, t_{i+1}]$

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

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

$x_i(t_i) = s_i$  Instead of one single integration, compute  $N$  integration starting from  $s_i$  (start of interval)

Obtain  $x_i(t)$ , then numerically compute integrals:  $l_i(s_i, y_i) = \int_{t_i}^{t_{i+1}} l(x_i(t), y_i) dt$

c. NLP

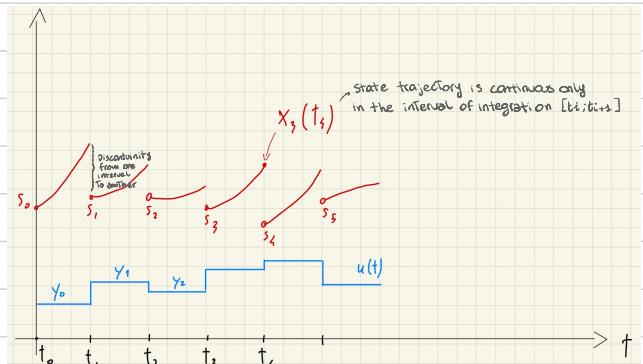
$$\min_{s, y} \sum_{i=0}^{N-1} l_i(s_i, y_i) + l_F(x_N)$$

s.t.  $s_0 - x_0 = 0$  initial conditions

Resulting of numerical integration from  $s_i$   $s_{i+1} - x_i(t_{i+1}, s_i, y_i) = 0$  continuity constraints

$g_i(s_i, y_i) \leq 0$  path constraints

Starting point of next segment must be exactly the point in which is finished previous segment



## DIRECT METHODS - SUMMARY

### ① SINGLE SHOOTING

- ⊕ Small problem
- ⊖ Need initial guess only for  $u(t)$
- ⊕ 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 adopt time grid

### ② COLLOCATION

- ⊖ Large problem
- ⊖ Cannot adapt time grid
- ⊕ Sparse problem
- ⊕ Work well with unstable systems
- ⊕ Can initialize  $x(t)$

## Model Predictive Control (MPC)

- Use optimal control for controlling system (robot). Non calcola la traiettoria di controllo una sola volta, ma riprende un processo in 3 fasi: a ogni passo temporale

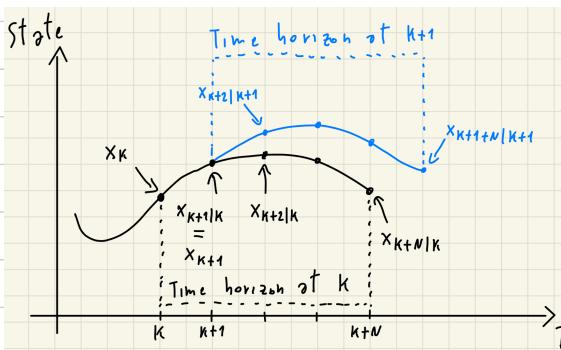
1. Solve a finite-horizon OCP using current state as initial state.

$$x^*, u^* = \arg \min_{x, u} \sum_{i=0}^{N-1} l(x_i, u_i) \quad \begin{matrix} \text{collocation} \\ \text{cost function is the sum of "transition} \end{matrix}$$

$$\begin{aligned} \text{s.t.} \quad & x_{i+1} = f(x_i, u_i) \quad i=0, \dots, N-1 \rightarrow \text{Finite-Horizon} \\ & x_{i+1} \in X \quad u_i \in U \quad i=0, \dots, N-1 \\ & x_0 = x^{\text{mean}} \end{aligned}$$

2. Apply computed control only for this step:  $u = u_0^*$
3. Repeat the process for the next step. In this way the control remain optimal despite changes in environment. For example, the presence of an obstacle can change set  $X$ .

AL passo temporale successivo, il sistema si troverà in un nuovo stato (dovuto al comando  $u_0^*$  ed eventuali disturbi). L'intero processo ricomincia dal punto 1.



### RECEDING HORIZON ORIZZONTE MOBILE

OCP at time  $K$  and  $K+1$  optimize over different horizons so can result in different Trajectories.

Even without disturbances, predicted and actual trajectory can be different. This happens because model  $f(x, u)$  is just an approximation. HPC try to fix iteratively this error.

CHALLENGES → FEASIBILITY: how to ensure feasibility and deal with other cases? ASSICURARE CHE CI SIA SEMPRE UNA SOLUZIONE  
 → STABILITY: can MPC stabilize the system? MPC è in grado di far convergere il sistema senza oscillazioni incontrollate  
 → COMPUTATION TIME: is it sufficiently fast? Deve essere risolto tra un passo temporale e l'altro

## INFINITE HORIZON MPC

$N = \infty \Rightarrow$  fornisce "garanzie teoriche", nella pratica si usa l'orizzonte finito

In this case pred.-act. trajectories are the same. This follows from Bellman's optimality principle:

Given optimal Traj.  $x^*(0), \dots, x^*(N)$ , the subtraj. is optimal for OCP in horizon  $[k, N]$  starting from  $x(k)$ .  $\Rightarrow$  optimal problem doesn't change at each step

With  $N = \infty$ , if  $\text{cost } l(x,u) \geq \alpha \|x\|_1 \forall x, u$  for some  $\alpha > 0$ , then having a finite cost implies stability.

Feasibility of the first OCP implies feasibility of all the OCP's.

IDEA: add terminal cost and constraints to finite-horizon OCP to mimic an infinite horizon.

$$V_N(\bar{x}) = \min_{\bar{u}} \sum_{k=0}^{N-1} l(x_k, u_k) + l_F(x_N) \quad \begin{matrix} \text{Estimation of total cost that would accumulate from } N \text{ to } \infty \\ \text{terminal cost} \end{matrix}$$

$$\text{s.t. } x_{k+1} = f(x_k, u_k) \quad k=0, \dots, N-1$$

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

$$x_0 = \bar{x}$$

$x_N \in X_F$  terminal constraint force terminal state to finish in a Terminal region  $X_F$

How to choose  $l_F(\cdot)$ ,  $x_N$ ?

### 1 FEASIBILITY

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

• INPUT CONSTRAINTS ONLY  $\rightarrow$  always feasible

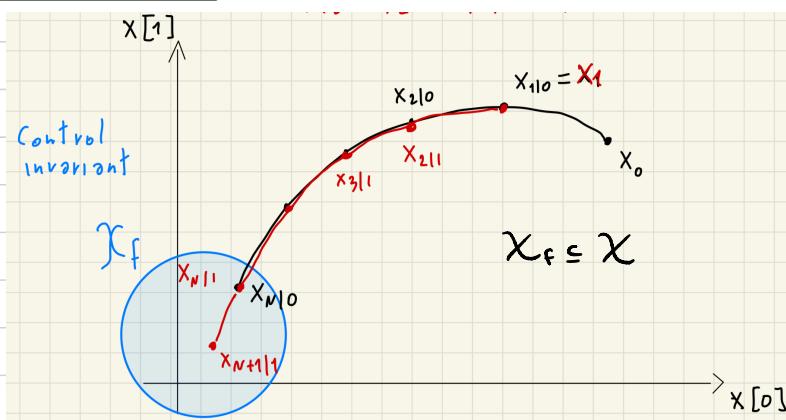
• HARD STATE CONSTRAINTS  $\rightarrow$  if  $N < \infty$  there is no guarantee that OCP remains feasible, even in nominal case

MAXIMUM OUTPUT ADMISSIBLE SET THEORY:  $\begin{cases} N < \infty \text{ is enough} \\ \text{hard to know which value of } N \text{ is enough} \\ \text{computation time may be too high to that } N \end{cases}$

**THEOREM** If  $X_F$  is control invariant  $\Rightarrow$  MPC is recursively feasible

**DEFINITION** A set  $S$  is control invariant if  $\forall x \in S, \exists u \in U \text{ s.t. } f(x, u) \in S$

#### VISUAL PROOF



$X_F$  is CONTROL INVARIANT SET: a system is of ci if, once you're in it, exists always a control that allows to remain in the set.

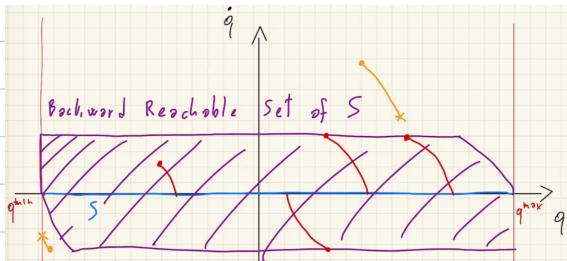
SUMMARY: an invariant set is like a safe region in which we can always find a solution

## CONTROL - INVARIANT SETS

Set of equilibria  $S$  is always C.I., but small:  $S = \{x \in \mathcal{X} \mid \exists u \in \mathcal{U} : x = F(x, u)\}$

Can use backward reachable set of  $S$ , i.e. states starting from which you can reach  $S$ .

↳ visto che  $S$  è troppo piccolo da raggiungere per MPC,  
cerchiamo di farlo arrivare in  
una regione da cui è possibile  
raggiungere  $S$



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

$$\min_{x, u} \quad 1$$

$$\text{s.t. } \begin{aligned} x_{i+1} &= f(x_i, u_i) & i=0, \dots, N-1 \\ x_{i+1} &\in \mathcal{X}, u_i \in \mathcal{U} & i=0, \dots, N-1 \\ x_N &= x_{N+1} \\ x_0 &= \bar{x} \end{aligned}$$

$$\begin{cases} \bar{x} \in B_N(S) & \text{if 3 solution} \\ \bar{x} \notin B_N(S) & \text{otherwise} \end{cases}$$

- Use FUNCTION APPROXIMATION to learn an approximate representation of  $B_N(S)$  from examples.
- To BUILD A DATASET, repeat the following steps:
  - sample state  $\bar{x}$
  - if  $\bar{x} \notin \mathcal{X} \Rightarrow$  add data point  $(\bar{x}, 0)$  to dataset
  - else try to solve OCP
  - if solution found  $\Rightarrow$  add  $(\bar{x}, 1)$  to dataset
  - else add  $(\bar{x}, 0)$  to dataset
- train classifier to predict label 0/1 given state  $\bar{x}$ .
- use the classifier as terminal constraint in MPC.

## 2 STABILITY

LYAPUNOV FUNCTION:  $V: \mathbb{R}^n \rightarrow \mathbb{R}$  is exponential Lyapunov function if  $\exists d_1, d_2, d_3 > 0$

↳ "misura di energia"  $V(x)$  del sistema.  
GOAL  $\rightarrow$  stabilizzare all'origine  $x=0$ , il punto con energia minima

$$\text{s.t. } \forall x \in \mathcal{X} :$$

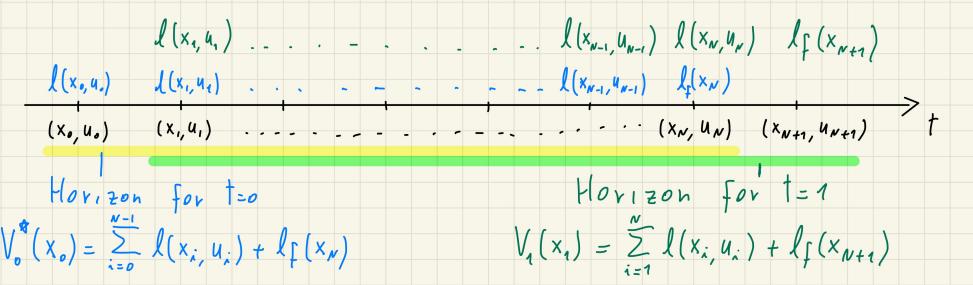
$$d_2 \|x\| \leq V(x) \leq d_1 \|x\| \Rightarrow V(x) > 0 \quad \forall x \neq 0$$

↳  $V(x)$  è una funzione convessa  $\Downarrow V(0) = 0$

$$V(F(x)) - V(x) \leq -d_3 \|x\| \quad \text{Ad ogni passo del sistema l'energia diminuisce}$$

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

$$\|x_{k+1}\| \leq c \gamma^k \|x_k\| \text{ for some } c > 0, \gamma \in [0, 1]$$



Value function as a Lyapunov function

Is  $V^*(x_1)$  a Lyapunov function?

$\text{so ASS. #2}$

$$V(x_1) - V^*(x_0) = \underline{l(x_N, u_N)} + \underline{l_f(x_{N+1})} - \underline{l_f(x_N)} - \underline{l(x_0, u_0)} \leq d_3|x_0| \quad \forall x_0 \in \mathcal{X}_F$$

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

$\leq -d_1|x_0| \text{ ASS. #3}$

### STABILITY ASSUMPTIONS

ASS #1:  $f(0,0)=0, \ell(0,0)=0, \ell_f(0)=0$

ASS #2:  $\forall x \in \mathcal{X}_F, \exists u \in \mathcal{U} \text{ s.t. } \ell_f(f(x,u)) - \ell_f(x) \leq -\ell(x,u)$  terminal cost decrease

ASS. #3:  $\exists d_1 > 0 \text{ s.t. } \forall x \in \mathcal{X}_N, \forall u \in \mathcal{U} \quad \mathcal{X}_N \text{ is set of states from which ocp has solution}$

$\ell(x,u) \geq d_1|x| \text{ cost lower bound}$

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

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 in equilibrium  $\Rightarrow$  Define new state  $\bar{x} = x - x^*$ , control  $\bar{u} = u - u^*$ , dynamics:

$$\begin{cases} \dot{x}^+ = f(x, u) \\ x^* = f(x^*, u^*) \end{cases} \quad \begin{cases} \dot{x}^+ = \bar{f}(\bar{x}, \bar{u}) \\ \bar{x}^+ = \bar{f}(\bar{x}, 0) \end{cases}$$

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

### STABILITY OPTIONS:

(1) Pick  $\mathcal{X}_F$  and  $\ell_f$  such that:  $-\mathcal{X}_F$  is control invariant

$-\forall x \in \mathcal{X}_F \quad \exists u \in \mathcal{U} \text{ s.t. } \ell_f(f(x,u)) - \ell_f(x) \leq -\ell(x,u)$

(2) Pick  $\mathcal{X}_F = \{0\} \Rightarrow$  small basin of attraction (this is a special case of (1))

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

$\hookrightarrow$  most common, especially for non linear systems. But large  $N \Rightarrow$  large computation time

### STABILITY EXTENSIONS

- Trajectory tracking can be cast as regulation of time-varying systems.

- Stability + recursive feasibility can still be ensured with time varying  $\mathcal{X}_F, \ell_f(\cdot)$ .

- If cost measures some kind of energy consumption  $\Rightarrow$  does not satisfy  $\ell(x,u) \geq d_1|x|$   
 $\Rightarrow$  Economic MPC

### ③ COMPUTATION TIME

KEY IDEA: warm start, for example using the solution of previous ocp as initial guess for current ocp.

Shift trajectory back by 1 step:  $u_n^{\text{guess}} \leftarrow u_{n-1}^*$

Use zero as initial guess for last time step  $u_{N-1}^{\text{guess}} \leftarrow 0$

Don't iterate until convergence, just do 1 Newton step. Possibly skip line search.

Use hardware accelerators (e.g. GPU).

### 3. REINFORCEMENT LEARNING

- Tries to find global optimal policy
- Assumes Dynamic is unknown
- Finite size state/control spaces
- Assumes Dynamic is stochastic (consider friction change during motion)
- Infinite Horizon

} Differences with optimal control

### Markov Decision Processes

Describe environment for RL problems.

**Markov Property:** FUTURE IS INDEPENDENT FROM THE PAST GIVEN THE PRESENT.  $P(x_{t+1}|x_t) = P(x_{t+1}|x_1, \dots, x_t)$

State transition probability  $P_{xx'} = P(x_{t+1} = x' | x_t = x) \Rightarrow$  matrix  $P = \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix}$  Each row sum = 1 due to Total probability law

Markov Process is defined as the tuple  $\langle X, P \rangle$  where  $X$  finite set of states

↳ Markov Chain

Since  $P$  has all nonnegative entries, and it's irreducible, by Perron - Frobenius theorem we know that:

PROP 1. The largest eigenvalue  $r$  of  $P$  satisfies:  $\min_i \sum_j P_{ij} \leq r \leq \max_j \sum_i P_{ij} \Rightarrow r=1$

PROP 2. All eigenvalues of  $P$  are smaller than 1:  $\lambda_i(P) < 1 \forall i$

↳ !! IMPORTANT PROPERTIES

In MP dynamics is encoded as pdf:  $P(x^+|x)$  Probability of transitioning from state  $x$  to next state  $x^+$ .

In OCP, for det. systems, dynamics was encoded as  $x^+ = f(x)$  or  $x^+ = f(x, w)$   $\xrightarrow{\text{uncertainty}}$  TRY TO USE THIS ONE

### REWARD PROCESS

Tuple  $\langle X, P, C, \gamma \rangle \rightarrow C = \text{cost function } C_x = \mathbb{E}[l_{t+1} | x_t = x]$  we work with expected values  $\Rightarrow$  stochastic

↳  $\gamma \in [0, 1]$  = discount factor, defined by user. It prevents cost to go to infinity

Negative cost  $\Rightarrow$  reward

### COST TO GO

Total discounted cost from time  $t$  to  $\infty$ :  $J_t = l_t + \gamma l_{t+1} + \dots = \sum_{k=0}^{\infty} \gamma^k l_{t+k}$  somma di tutti i costi futuri da tempo  $t$  a  $\infty$

↳  $\gamma$  represents preference for later costs over immediate cost

$\left\{ \begin{array}{l} \gamma \rightarrow 0 \text{ myopic evaluation} \\ \gamma \rightarrow 1 \text{ far sighted evaluation} \end{array} \right.$	Found a trade-off to prevent $J_t = \infty$ (since we sum till infinity)
--	--

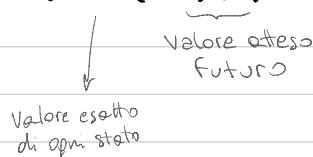
### VALUE FUNCTION

Cost starting from state  $x$ :  $V(x) = J_t(x_t = x) = l_t + \gamma(l_{t+1} + \gamma l_{t+2} + \dots)$

↳ can be decomposed in:  $V(x) = l(x) + \gamma V(f(x)) \Rightarrow$  Bellman Equation

↳ matrix form  $V = C + \gamma PV \Rightarrow \begin{bmatrix} V(s) \\ \vdots \\ V(n) \end{bmatrix} = \begin{bmatrix} l(s) \\ \vdots \\ l(n) \end{bmatrix} + \gamma \begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & \ddots & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix} \begin{bmatrix} V(1) \\ \vdots \\ V(n) \end{bmatrix}$  IL valore di stare in uno stato è il costo  
due paghi adesso + il valore scontato di stare finito dopo  
Expected value to go from state s to other n states

Linear equation in  $V$ :  $V = (I - \gamma P)^{-1} C$  the solution is found with iterative methods



## HDP

Tuple  $\langle X, U, P, C, \gamma \rangle \rightarrow$  Introduced control input  $u$  (action  $a$  in RL terminology)

$$\begin{aligned} P_{xx'}^u &= P(X_{t+1} = x' \mid X_t = x, U_t = u) \\ C_x^u &= f_t(X_t = x, U_t = u) \end{aligned} \quad \left\{ \text{cost and prob. depend also on control input}$$

Policy  $\Rightarrow$  acts like a controller, its a distribution over control inputs given states:  $\pi(u|x) = P(U_t = u \mid X_t = x)$

We will assume deterministic policies  $u = \pi(x)$ . We can use stochastic policies during policy search process.

### ACTION-VALUE FUNCTION

It's the cost starting from state  $x$ , applying input  $u$ , and then following policy  $\pi$ . You don't know the future cost, it depends on the next action you decide to take.

$$Q^\pi(x, u) = J_t(X_t = x, U_t = u, M_{t+1} \sim \pi)$$

Can be decomposed as:  $Q^\pi(x, u) = \ell(x, u) + \gamma Q^\pi(f(x, u), \pi(f(x, u)) = \ell(x, u) + \gamma V^\pi(f(x, u))$

where:  $V^\pi(x) = Q^\pi(x, \pi(x))$  Relation between value function and action-value function  
first action it's the only that doesn't come from action set  $\pi$

### OPTIMAL VALUE FUNCTION

- 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:  $\pi^* \leq \pi \wedge \pi^* \quad [\text{where } \pi \leq \pi^* \text{ if } V^\pi(x) \leq V^*(x) \forall x]$

### OPTIMAL POLICY

Se conosci la funzione Q ottimale, non devi pianificare nulla. Ti basta essere avido (greedy) e scegliere l'azione che minimizza il costo istantaneo combinato col futuro. Ecco perché in Reinforcement Learning (es. Q-Learning) cerchiamo disperatamente di stimare Q: se hai Q, hai il controllo perfetto.

↳ If we know  $Q^*(x, u)$  we immediately have optimal policy  $\pi^*(x) = \min_u \ell(x, u) + \gamma V(f(x, u))$

### BELLMAN OPTIMAL EQUATION

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

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

$$\begin{aligned} V^*(x) &= \min_u \ell(x, u) + \gamma V^*(f(x, u)) \\ Q^*(x, u) &= \ell(x, u) + \gamma \min_u Q^*(f(x, u), u) \end{aligned} \quad \left\{ \begin{array}{l} \cdot \text{Non linear} \\ \cdot \text{No closed-form solution} \\ \cdot \text{Iterative Algorithms} \end{array} \right.$$

$$\text{LINEAR: } V^* = \ell + \gamma P V^* \quad | \quad \text{NON LINEAR: } V^* = \min_u \ell + \gamma P V^*$$

# Dynamic Programming

Assume full knowledge of MDP, we have 2 classes of problems:

1. PREDICTION { Input: MDP  $\langle X, U, P, C, \gamma \rangle$  Policy  $\pi^*$   
Output: Value function  $V^*$

2. CONTROL { Input: MDP  
Output: Optimal Value Function  $V^*$   
Optimal Policy  $\pi^*$

Prediction (Evaluation): Ti do una policy fissa  $\pi$ . Tu dimmi quanto vale ( $V\pi$ ). È un lavoro da contabile: non cambi nulla, calcoli solo il disastro (o il successo) futuro.

Control (Optimization): Non mi interessa la policy attuale. Voglio la migliore policy possibile ( $\pi^*$ ) e il miglior valore possibile ( $V^*$ )

## FINITE HORIZON

- theory is simpler
- used in optimal control
- policy and value depends on time

## INFINITE HORIZON

- complex theory
- used in RL
- Policy and value are stationary
- good approximation of problems with finite but long horizon

All algorithm of DP rely on Bellman operators.

BELLMAN EXPECTATION OPERATOR  $T^\pi(V) = C^\pi + \gamma P^\pi V$  Modifies a given value function

BELLMAN OPTIMALITY OPERATOR  $T(V) = \min_{\pi \in \Pi} C^\pi + \gamma P^\pi V$  Approssima il valore assumendo che d'ora in poi farò sempre la mossa migliore

$$(TV)(x) = \min_{u \in U} l(x, u) + \gamma V(F(x, u)) \leftarrow \text{Different Notation}$$

Approssima il valore basandosi sulle probabilità restanti delle probabilità, seguendo la policy attuale.

## PREDICTION PROBLEM

PROBLEM: evaluating a given policy  $\pi$

SOLUTION: Start from arbitrary estimate of value function  $V_0$

Apply iteratively Bellman expectation operator:  $V_{k+1}(x) = l(x, \pi(x)) + \gamma V_k(F(x, \pi(x))) \quad \forall x \in X$   
 $\downarrow V_{k+1}(x) = T^\pi(V_k)$

Guaranteed to converge to  $V^\pi$ :  $\lim_{k \rightarrow \infty} V_k = V^\pi$

Stop when  $\|V_{k+1} - V_k\| < \text{threshold}$  where threshold  $\rightarrow 0$

## 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 wrt  $V^{\pi_k} \Rightarrow \pi_{k+1}(x) = \arg \min_u l(x, u) + \gamma V^{\pi_k}(F(x, u))$

Always converges to  $\pi^*$ :  $\lim_{k \rightarrow \infty} \pi_k = \pi^*$

Una volta che hai i valori veri, approssima la policy scegliendo l'azione che minimizza il costo immediato + costo futuro

Evaluate exactly a policy can take many iterations. If we compute an approximation of  $V^{\pi_k}$  we use MN evaluation iteration to compute  $V^{\pi_k}$ .

Under some mild assumptions it converges to  $V^*$ .

for  $MN = \infty$  we recover policy iteration

for  $MN = 0$  we get value iteration

Every time you update  $V$  you use a policy that is greedy wrt  $V$ :

$$V^k(x) = l(x, \pi^k) + \gamma V^{k-1}(F(x, \pi^k)) \quad \text{with } \pi^k(x) = \arg \min_u l(x, u) + \gamma V^{k-1}(F(x, u))$$

$$\Rightarrow V^k(x) = \min_u l(x, u) + \gamma V^{k-1}(F(x, u))$$

## CONTROL: VALUE ITERATION

Algorithm:

- Start with arbitrary value function  $V_0$

For  $k=0 \dots \infty$

For any  $x \in \mathcal{X}$

$$V_{k+1}(x) = \min_{u \in U} l(x, u) + \gamma V_k(f(x, u))$$

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

- Optimal policy  $\pi^*$  is computed at the end from  $V^* \Rightarrow \pi^* = \arg \min_u l + \gamma V^*$

## Model Free Prediction

PROBLEM: finding value function of unknown MDP by acting and observing.

SOLUTION: use MONTE CARLO METHODS.

↳ IDEA: value = mean return. Can only be applied to episodic MDP, and all episodes must end.

### MC POLICY EVALUATION

Goal: Learn  $V^\pi$  from episodes of experience under policy  $\pi$ :  $x_1, m_1, l_1, x_2, m_2, l_2, \dots, x_n \sim \pi$

Cost-to-go is total discounted cost:  $J_t = l_t + \gamma l_{t+1} + \dots + \gamma^{T-1} l_{t+T-1}$

Value function is expected cost-to-go under policy  $\pi$ :  $V^\pi(x) = E[J_t | x_t = x]$

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

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  $N(x) \rightarrow \infty$

### INCREMENTAL MC UPDATES

↳ IDEA: compute mean incrementally:  $V_N = \frac{1}{N} \sum_{j=1}^N J_j = \frac{J_N + \sum_{i=1}^{N-1} J_i}{N} = V_{N-1} + \frac{1}{N} (J_N - V_{N-1})$

In a non-stationary problems it can be useful to track a running mean, i.e. forget all episodes:  $V(x) = V(x) + \alpha (J - V(x))$

## TEMPORAL DIFFERENCE LEARNING

It's an alternative to MC.

Estimate cost-to-go by 1-step look ahead (TD):  $V(x_+) = V(x_+) + \alpha_t \left( l_t + \gamma V(x_{t+1}) - V(x_+) \right)$

Policy evaluation is TD with  $\alpha_t = 1$  and sampling all states uniformly.

PROPERTIES:

- stochastic approximation algorithm

- If TD converges  $\Rightarrow$  it must converge to  $V^\pi$

- Convergence is guaranteed if step-size sequence satisfies RM conditions:

$$\sum_{t=0}^{\infty} \alpha_t = \infty \quad \sum_{t=0}^{\infty} \alpha_t' < \infty$$

TD ERROR

TD TARGET

MC

wait end of episode to know cost

episodic environments

cost-to-go  $J_t$  is unbiased estimate of  $V_{\pi}(x_t)$

High variance, zero bias

Not very sensitive to initialization

TD

Learn after every step

non-terminating environment

target is biased estimate of  $V_{\pi}(x_t)$

target has lower variance than  $J_t$

More sensitive to initial value

### EXAMPLE SLIDE 149, 150, 151

BOOTSTRAPPING: update involves an estimate

SAMPLING: updates sample an expectation

		Root Samp.	NO	YES
		Exhaustive Search	DP	
P	Known		NO	TD
	Unknown	HC		

## Model free control

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

Since we don't know MDP, we can't compute  $V(x')$ .

We have to learn  $Q$  instead of  $V$ .

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

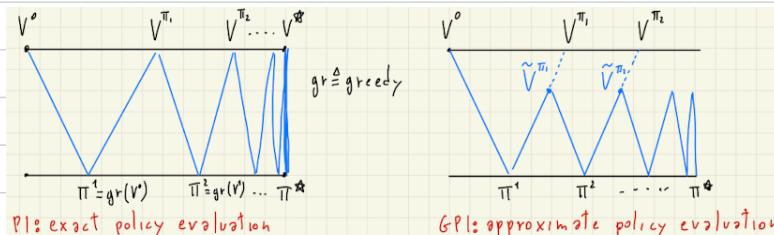
### ACTOR-CRITIC METHODS

Extension of policy iteration to unknown MDP. Policy Iteration alternates

"policy evaluation" and "policy improvement".

Since we don't know MDP, policy evaluation can take infinitely many samples.

The solution is to improve policy based on approximate value function.



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 transitions is typically not the same that is evaluated and improved. The behavior policy mixes small amount of exploration into target policy.

## SARSA

Use TD to learn  $Q$  from on-policy samples. After each transition  $(x, u, r, x', u')$  update  $Q$  as follows:

$$S(Q) = r(x, u) + \gamma Q(x', u') - Q(x, u)$$

$$Q_{k+1}(x, u) = Q_k(x, u) + \alpha_k S(Q)$$

If  $\pi$  was fixed  $\Rightarrow$  SARSA = TD

SARSA can diverge in off-policy situations.

To choose control we take policy that is greedy wrt  $Q$ . Greedy policy can be computed on the fly  $\Rightarrow$  faster

- GREEDY STRATEGY: always choose control with best estimated cost. It can fail to find best action. A good learner must choose suboptimal control to explore.
- $\epsilon$  GREEDY STRATEGY: choose a random control with probability  $\epsilon$ . Choose greedy control with probability  $1-\epsilon$ . Simple way to ensure exploration. Probability  $\epsilon$  can change over time (typically decreasing to 0).

Using  $\epsilon$ -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 pairs are visited infinitely many times.

- policy converges on greedy policy

- 2) Robbins-Monro sequence of step sizes:  $\sum_{k=1}^{\infty} \alpha_k = \infty$     $\sum_{k=1}^{\infty} \alpha_k^2 < \infty$

### ON-POLICY

e.g. SARSA

Learn about policy  $\pi$  from experience sampled from  $\pi$

Behavior policy  $\approx$  Target policy

Limited exploration  $\Rightarrow$  low sample efficiency

### OFF-POLICY

e.g. Q-learning

Learn about policy  $\pi$  from experience sampled from  $\gamma$

Behavior policy  $\gamma$  can be different from target policy  $\pi$

Higher sample efficiency

## DIRECT METHODS: Q-LEARNING

Iteratively update an estimate  $Q_\pi(x, u)$  of  $Q^*(x, u)$ . After observing transition  $(x, u, r, x')$  update with:

$$S(Q) := r + \gamma \min_{w \in V} Q(x', w) - Q(x, u) \quad \rightarrow \text{Use } \min_w Q(x', w) \text{ instead of } Q(x', u)$$

$$Q_{k+1}(x, u) := Q_k(x, u) + \alpha_k S(Q)$$

Error  $\delta(Q(x, u))$  is independent from behavior policy, which is the one that generates  $u \Rightarrow$  off-policy

Target policy is:  $\pi(x) = \arg \min_u Q(x, u)$

At stochastic equilibrium,  $V(x, u)$  visited infinitely often:  $E[\delta(Q)] = 0 = TQ - Q \Rightarrow Q = Q^*$

$Q_k$  converges when appropriate learning rates are used and sufficient exploration is ensured.

## PSEUDOCODE

### SARSA

```

INPUTS: Q, ε
For k = 0 ... K
    x ← env.reset()
    u ← EpsGreedy(Q, x, ε)
    For t = 0 ... N
        x', l ← env.step(u)
        u' ← EpsGreedy(Q, x', ε)
        Qtarget ← l + γ Q[x', u']
        Q[x, u] += α (Qtarget - Q[x, u])
        x, u ← x', u'
    ε ← updateEpsilon(ε)

```

### Q-LEARNING

```

INPUTS: Q, behavior policy πB
For k = 0 ... K
    x ← env.reset()
    For t = 0 ... N
        u ← πB(x)
        x', l ← env.step(u)
        Qtarget ← l + γ minu' (Q[x', u'])
        Q[x, u] += α (Qtarget - Q[x, u])
        x ← x'

```

## Summary

DP	TD
Policy Evaluation $V(x) = l + \gamma V(x')$	TD-Learning $V(x) \leftarrow l + \gamma V(x')$
Policy Iteration $V(x) = l + \gamma V(x')$ $\pi(x) = \arg \min_u l + \gamma V(x'(u))$	SARSA $Q(x, u) \stackrel{d}{\leftarrow} l + \gamma Q(x', u')$ $\pi(x) = \begin{cases} \arg \min_u Q(x, u) & \text{with prob. } 1-\epsilon \\ \text{random} & \text{with prob. } \epsilon \end{cases}$
Value Iteration $V(x) = \min_u l + \gamma V(x'(u))$	Q-Learning $Q(x, u) \stackrel{d}{\leftarrow} l + \gamma \min_{u'} Q(x', u')$

## Value function Approximation [continuous time]

Use function approximation to represent  $V(x)/Q(x, u)$

$$\begin{array}{l} \text{High dimensional functions} \\ \left. \begin{array}{l} \hat{V}(x, w) \approx V(x) \\ \hat{Q}(x, u, w) \approx Q(x, u) \end{array} \right\} w: \text{parameters defining the function} \end{array}$$

$$x \rightarrow \boxed{w} \rightarrow \hat{V}(x, w) \quad \text{Neural networks}$$

$$x \rightarrow \boxed{w} \rightarrow \hat{Q}(x, u, w) \quad \text{Differentiable Approximator: linear combination of features, Fourier basis, Polynomials, ...}$$

$$x \rightarrow \boxed{w} \rightarrow \hat{Q}(x, u_*, w) \quad \text{Non-stationary target } V$$

$$\hat{Q}(x, u_{\text{m}, w}) \quad \text{Ahn i.i.d. data}$$

### VALUE APPROXIMATION BY STOCHASTIC GRADIENT DESCENT

$$J_w = E_{\pi} [(V_{\pi}(x) - \hat{V}(x, w))^2]$$

$$\Delta w = -\frac{\alpha}{2} \nabla_w J = \alpha E_{\pi} [(V_{\pi} - \hat{V}) \nabla_w \hat{V}]$$

SGD samples the gradient:  $\Delta w = \alpha (V_{\pi}(x) - \hat{V}(x, w)) \nabla_w \hat{V}(x, w)$

Expected update is equal to full gradient update.

How we get target  $y = V(x)$ ?  $\left\{ \begin{array}{l} \text{For MC is cost-to-go } J_t(x) : \Delta w = \alpha (J_t - \phi(x_t, w)) \nabla_w \phi(x_t, w) \\ \text{For TD is TD target: } \Delta w = \alpha (r_t + \delta \phi(x_{t+1}, w) - \phi(x_t, w)) \nabla_w \phi(x_t, w) \end{array} \right.$

### INCREMENTAL CONTROL WITH FUNCTION APPROXIMATION

Generalized Policy Iteration with approximate Action-Value Function:  $\min_w E_{\pi} [(\phi(x, u, w) - Q_{\pi}(x, u))^2]$

Use SGD to find local minimum.

Use  $\epsilon$ -greedy policy to ensure exploration.

### CONVERGENCE OF CONTROL ALGORITHMS

	VALUE ENCODING		
	Table Look up	Linear	Nonlinear
MC control	✓	(V)	X
Sarsa	✓	(V)	X
Q-Learning	✓	X	X

(V) = chatter around optimum

LIMITATIONS OF INCREMENTAL CONTROL  $\left\{ \begin{array}{l} \text{Poor data efficiency} \\ \text{Data is not i.i.d.} \end{array} \right.$

## BATCH LEARNING

Instead of learning incrementally, learn after having collected a batch of training data.

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

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

- 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.
- Use fixed  $Q$  targets to stabilizing training  
↳ Store  $w$  resulting in best performance as performance may get worse during training.