



Notes of the Course

Advanced Optimization-Based Robot Control

prof. Del Prete Andrea

andrea.delprete@unitn.it

February 26, 2023

Author:

Matteo Dalle Vedove (*matteodv99tn@gmail.com*)

Contents

1	Reactive Controls	4
1.1	Joint space motion control	4
1.2	Cartesian space motion control	6
1.3	Impedance control	7
1.4	QP-based reactive control	8
1.4.1	Joint space control	8
1.4.2	Review of convex optimization	10
1.4.3	Cartesian space control	11
1.4.4	Task space control	12
2	Optimal Control	14
2.1	Dynamic programming	15
2.2	Hamilton-Jacobi-Bellman	18
2.3	Pontryagin minimum principle	18
2.4	Direct methods	19
2.4.1	Single shooting	19
2.4.2	Numerical integration	20
2.4.3	Computation of sensitivities	22
2.4.4	Collocation	23
2.4.5	Multiple shooting	24
2.5	Linear quadratic regulator	24
2.5.1	Steady-state regulation	27
2.5.2	Inhomogeneous systems and tracking problems	27
2.6	Differential dynamic programming	28
2.6.1	Iterative LQR	30
2.7	Model predictive control	30
2.7.1	Feasibility	32
2.7.2	Stability	33
2.7.3	Computation time	34
3	Reinforcement Learning	35
3.1	Markov decision process	36
3.2	Dynamic programming	38
3.2.1	Prediction: iterative policy evaluation	39
3.2.2	Control: policy and value iteration	40
3.3	Model-free prediction	41
3.3.1	Monte Carlo policy evaluation	41
3.3.2	Temporal difference learning	42
3.4	Control learning methods	44
3.4.1	Q learning	44

3.4.2	Actor-critic	46
-------	------------------------	----

Chapter 1

Reactive Controls

Usually while dealing with manipulators (or robot in generals) we are interested in controlling it's *dynamic*, so affecting the system's behavior. Let call $q \in \mathbb{R}^n$ the *Lagrangian coordinate* that are used to describe the system itself and denoted with $\dot{q} \in \mathbb{R}^n$ their velocity, the dynamic of the system is described by the following differential equation:

$$M(q)\ddot{q} + C(q, \dot{q})\dot{q} + g(q) = \tau + J^\top(q)w \quad (1.1)$$

where $M \in \mathbb{R}^{n \times n}$ is positive-definite *mass-matrix* (and usually depends just on the system's configuration), $C \in \mathbb{R}^{n \times n}$ is a matrix that takes into account both centrifugal and Coriolis force, while $g \in \mathbb{R}^n$ is the elements describing the gravity; $\tau \in \mathbb{R}^n$ is the set of joint actions, while w is the *wrench* (the external actions) that are weighted by the end effector's Jacobian J in order to be projected in the joint space. Usually both C and g are condensed in the so called *bias forces* $h \in \mathbb{R}^n$, rewriting (1.1) as

$$M(q)\ddot{q} + h(q, \dot{q}) = \tau + J^\top(q)w \quad (1.2)$$

It must be observed that M, h, J are all non-linear expression with respect to the system's states q, \dot{q} , so linear theory can't be properly applied. Still (1.2) is linear with respect to \ddot{q} and τ (once q, \dot{q} and w are fixed).

1.1 Joint space motion control

The *control* problem addresses the question of finding the proper actuator actions τ that makes the system behaves according to a desired trajectory; mathematically, denoted with $q^{ref}(t)$ the *reference trajectory* desired for all the joints, the control problem wants to find the proper time-history $\tau(t)$ that make such motion feasible, so

$$\text{find } \tau(t) \quad \text{such that } q(t) \simeq q^{ref}(t)$$

Direct dynamic controller Assumed that the system is known and deterministic (hypothesis carried out throughout the whole course), then it means that both M and h in (1.2) are known; the simplest strategy that we can use to solve this problem is by means of the *Proportional Integrative Derivative (PID)* controlled for which

$$\tau(t) = K_p e(t) + K_d \dot{e}(t) + K_i \int_0^t e(\zeta) d\zeta \quad (1.3)$$

where $e(t) = q^{ref}(t) - q(t)$ is the *error* between the reference trajectory and the real measured coordinates and $K_p, K_d, K_i \in \mathbb{R}^{n \times n}$ are respectively the *proportional*, *derivative* and *integral gain matrices* that must be positive-definite.

As a rule of thumb such controller parameters K_x are tuned by firstly setting the proportional gain, following the derivative one and lastly the integral part. In general the best tracking performance is achieved by increasing the proportional gain K_p : this however leads to a *stiff* system that's usually less safer for human interaction (since high forces are exerted by the robot itself). For this reason more advanced techniques are used to solve the control problem in joint coordinate.

The controller in (1.3) exploits the *direct kinematic*, since we directly compute a value of the torque that must be applied to the different joints.

Inverse dynamic control In inverse dynamic the approach followed is a bit different, where the controller computes a desired acceleration \ddot{q}^d for the joint coordinates and the torques are instead computed by exploiting the system's dynamic (1.2):

$$\begin{cases} \tau &= M\ddot{q}^d + h \\ \ddot{q}^d &= \ddot{q}^{ref} + K_d\dot{e} + K_p e \end{cases} \quad (1.4)$$

In this case we observe that τ linearly depends \ddot{q}^d : for this reason this technique is also known as *feedback linearization* or *computed torques*.

Using the control (1.4) for the closed-loop dynamic provides

$$\begin{aligned} M\ddot{q} + h &= \tau \\ &= M\ddot{q}^d + h \\ \ddot{q} &= \ddot{q}^d = \ddot{q}^{ref} + K_d\dot{e} + K_p e \\ 0 &= \ddot{e} + K_d\dot{e} + K_p e \end{aligned}$$

At this point the stability of the system can be analyzed with linear techniques by reducing the previous differential equation to first order; defined in fact the vector $y = (e, \dot{e}) \in \mathbb{R}^{2n}$, the system reduces to the following linear dynamics:

$$\dot{y} = Ay \quad \text{with } A = \begin{bmatrix} 0 & I \\ -K_p & -K_d \end{bmatrix}$$

Linear theory provides us ways to tune K_p, K_d in such a way that the overall matrix turns stable (by enforcing that all eigenvalues of A have a negative real part).

Controller (1.4) usually leads to less stiffer system (that's a desirable feature) but a very good a-priori knowledge of the systems dynamics is required, i.e. (1.2) must be known upfront.

Integral gain In (1.3) it has been necessary to embed an integral term in the control law to compensate the effect of gravity (since no system's knowledge is required): this effect is non-linear in the robots configuration and the integral gain aims at compensating such contribution asymptotically.

In (1.4) the integral term is instead missing since all gravitational actions are modeled by the bias force h , thus is automatically compensated by the computed torques law (\ddot{q}^d must depends just on the error in tracking the reference trajectory and shouldn't compensate the gravity).

1.2 Cartesian space motion control

Natural environment can't be easily described in joint coordinate, but the easiest and intuitive choice is in cartesian space. A solution to control systems directly in cartesian coordinate is by solving at each control's sampling period the inverse kinematic problem and then compute the desired joint torque; this solution theoretically works, however we have to deal with the high non-linearity of the kinematic chain that makes the analytical solution of the problem impossible even for the simplest examples.

The approach here suggested aims at directly compute the desired joint torques $\tau(t)$ for which the end-effectors trajectory $x(t) \in \mathbb{R}^6$ (in cartesian coordinate) follows a desired motion $x^{ref}(t)$:

$$x(t) \simeq x^{ref}(t)$$

Operational space control Pre-multiplying the systems dynamic (1.2) by JM^{-1} allows us to map the equation from joint to cartesian space, reaching the following result:

$$\begin{aligned} JM^{-1}M\ddot{q} + JM^{-1}h &= JM^{-1}\tau \\ \ddot{x} - \dot{J}\dot{q} + JM^{-1}h &= \end{aligned} \quad (1.5)$$

where $J\ddot{q}$ has been substituted by the derivative of $\dot{x} = J\dot{q}$ that evaluated to $\ddot{x} = J\ddot{q} + \dot{J}\dot{q}$.

Assuming that the Jacobian J is full-row rank (i.e. the end effector is allowed to move in all directions) then $JM^{-1}J^\top \in \mathbb{R}^{6 \times 6}$ is invertible; pre-multiplying (1.5) by the matrix Λ defined as $(JM^{-1}J^\top)^{-1}$ gives

$$\Lambda\ddot{x} + \underbrace{\Lambda(JM^{-1}h - \dot{J}\dot{q})}_{\mu} = \Lambda JM^{-1}\tau$$

Regarding the joint torques τ as a linear combination of a vector $f_\tau \in \mathbb{R}^6$ describing the joint action in cartesian space, then we have that $\tau = J^\top f_\tau$: the previous equation can so be rewritten as

$$\begin{aligned} \Lambda\ddot{x} + \mu &= \Lambda JM^{-1}J^\top f_\tau = (JM^{-1}J^\top)^{-1} JM^{-1}J^\top f_\tau \\ &= f_\tau \end{aligned} \quad (1.6)$$

The controller that we can use to compute the joint torques given a reference trajectory in cartesian space can be simply regarded as

$$\begin{cases} \tau &= J^\top f_\tau \\ f_\tau &= \Lambda\ddot{x}^d + \mu \\ \ddot{x}^d &= \ddot{x}^{ref} + K_d(\dot{x}^{ref} - \dot{x}) + K_p(x^{ref} - x) \end{cases} \quad (1.7)$$

To prove the effectiveness of the controller (1.7) it's simply necessary to compute the close-loop dynamics

$$\begin{aligned} JM^{-1}(M\ddot{q} + h) &= \tau = J^\top \Lambda\ddot{x}^d + h \\ J\ddot{q} + JM^{-1}h &= JM^{-1}J^\top \Lambda\ddot{x}^d + JM^{-1}h \\ J\ddot{q} &= \ddot{x}^d \\ \ddot{x} - \dot{J}\dot{q} &= \end{aligned}$$

As long as the systems doesn't move too fast, the term $\dot{J}\dot{q}$ is negligible and the end effectors acceleration \ddot{x} matches the desired one \ddot{x}^d .

Redundancy resolution Controller (1.7) aims just at moving the end-effector in cartesian coordinate; we observe that the number of actuated joint exceeds the degrees of freedom required, then there are theoretically infinite configuration that allows to achieve the same final pose. The suggested control tracks the desired configuration with no problem, but it might happen over time that the system might approach *singular configurations* for which J is no longer full-row rank (losing so some direction of motion): in this case $JM^{-1}J^\top$ is no longer invertible and the control do not exists.

Redundancy resolution aims so at solving this problem: since there are infinitely many configurations that allows to achieve a desired x , the idea is now to a torque τ that do not impact the end-effector motion but helps avoiding singular configuration.

Observed that $\Lambda JM^{-1} = (JM^{-1}J^\top)^{-1}JM^{-1}$ is the *left-inverse* of J^\top , i.e.

$$\text{let } (JM^{-1}J^\top)^{-1}JM^{-1} = J^{\top\dagger} \quad \text{then } J^{\top\dagger}J^\top = I$$

Note that in general $J^\top J^{\top\dagger} \neq I$.

If we now consider any torque τ_0 in the kernel of $J^{\top\dagger}$ it happens that the dynamic of the system in cartesian space is not affected at all; recalling (1.6) it happens that

$$\Lambda \ddot{x} + \mu = J^{\top\dagger}(\tau + \tau_0) = J^{\top\dagger}\tau \quad \forall \tau_0 \in \ker\{J^{\top\dagger}\}$$

Computing τ_0 can be hard since there's no upfront guarantee that the motion will avoid singular configuration, however we can compute such torque in a smart way as

$$\tau_0 = (I - J^\top J^{\top\dagger})\tau_1 \quad \forall \tau_1 \in \mathbb{R}^n$$

Such choice always lies in $\ker\{J^{\top\dagger}\}$, in fact

$$\begin{aligned} J^{\top\dagger}\tau_0 &= J^{\top\dagger}(I - J^\top J^{\top\dagger})\tau_1 \\ &= (J^{\top\dagger} - J^{\top\dagger}J^\top J^{\top\dagger})\tau_1 \\ &= (J^{\top\dagger} - J^{\top\dagger})\tau_1 = 0 \quad \forall \tau_1 \in \mathbb{R}^n \end{aligned}$$

With this premise we can improve controller (1.7) by adding an extra term τ_1 that's proportional to a reference joint configuration that the system should tends to in order to avoid singular configurations:

$$\begin{cases} \tau &= J^\top f_\tau (I - J^\top J^{\top\dagger})\tau_1 \\ f_\tau &= \Lambda \ddot{x}^d + \mu \\ \ddot{x}^d &= \ddot{x}^{ref} + K_d(\dot{x}^{ref} - \dot{x}) + K_p(x^{ref} - x) \\ \tau_1 &= M(K_{p1}(q^{ref1} - q) - K_{d1}\dot{q}) + h \end{cases} \quad (1.8)$$

1.3 Impedance control

Controllers (1.7) and (1.8) allow to command the system's torque directly in cartesian coordinate; such technique to properly work requires a perfect knowledge of the system of interest.

Lack of information in the model can still be compensated to a certain extent by increasing the proportional gains in the controller, improving the tracking capability of the reference trajectory; still this solution leads to an unsafe stiff system, since in case of contact the controller might require force high enough to brake either the robot or the environment.

To overcome this limitation we can use *interaction control*, particularly useful to model contact of the end-effector with the surrounding environment. The underlying idea is that if we have access to an estimate of the contact force, then we can regulate the motion accordingly. Based on this simple idea there are *direct force controllers*, mainly focusing on controlling the motion based just on the applied contact forces (not covered in this course), and *indirect force controls*. In the latter case it's not specified a force that the end-effector needs to apply, but the contact is modeled by a mathematical relation relating force and position: for this reason the method is also known as *impedance control*.

Example of impedance control (in cartesian space) is the one where the contact is modeled by a mass-spring-damper dynamics of the form

$$M_d \ddot{x} + B \dot{x} + Kx = f \quad (1.9)$$

where f is the contact force and M_d, B, K are matrices that can be tuned to achieve a desired contact behavior.

Once the contact model (1.9) is defined, the controller needs to compute at each time instant the joint torques τ to follow the desired behavior. Starting from (1.6), the dynamic projected in cartesian space is simply $\Lambda \ddot{x} + \mu = J^{\top} \tau + f$; solving explicitly (1.9) for the acceleration provides the desired value $\ddot{x}^d = M_d^{-1}(f - B\dot{x} - Kx)$ that substituted in the dynamics

$$\Lambda M_d^{-1}(f - B\dot{x} - Kx) + \mu = J^{\top} \tau + f$$

Regarding the torques τ with the respective cartesian projection $J^{\top} f_{\tau}$, then the previous equation can be solved for f_{τ} as

$$\begin{aligned} f_{\tau} &= \Lambda M_d^{-1}(f - B\dot{x} - Kx) + \mu - f \\ &= -\Lambda M_d^{-1}(B\dot{x} + Kx) + \mu + (\Lambda M_d^{-1} - I)f \end{aligned}$$

At this point in order to compute τ (from f_{τ}) it's required an estimate of the contact force f that, unfortunately, usually comes from very noisy sensors and with time delays that might lead to unstable systems. A trick that we can use to overcome this limitation is to build the contact model (1.9) with $M_d = \Lambda$: with such choice f_{τ} no longer depends from the contact force f (since ΛM_d^{-1} will evaluate simply to I). At the end the resulting torque simplifies to the expression

$$\tau = J^{\top} f_{\tau} = J^{\top} (-Kx - B\dot{x} + \mu) \quad (1.10)$$

1.4 QP-based reactive control

Up to now only simplest controller have been presented for which optimality condition have not been discussed. In this section we will focus more on *optimization-based approaches* that exploits minimization to determine the optimal control solution with the addition of constraints (such actuator effort, joint position and velocity limits...) that until now it was not possible to consider.

1.4.1 Joint space control

To control a desired trajectory in joint space (so given q^{ref}), controller (1.4) has been proposed for which

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

The same result can be achieved by solving the following *optimal control problem* (OCP):

$$\begin{aligned} \min_{\tau, \ddot{q}} \quad & \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2 = \frac{1}{2} (\ddot{q} - \ddot{q}^d)^\top (\ddot{q} - \ddot{q}^d) \\ \text{sub. to:} \quad & M\ddot{q} + h = \tau \end{aligned} \quad (1.11)$$

where $\frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2$ is the *cost function* that needs to be minimized with respect to the *minimization variables* \ddot{q}, τ and that must be subjected to the *constraint* $M\ddot{q} + h = \tau$ (the dynamic law).

Since in (1.11) the cost is a quadratic function, the minimum is obtained by setting the argument of the norm to zero, i.e. the optimal acceleration is $\ddot{q}^* = \ddot{q}^d$; this implies that the optimal commanded torque is simply $\tau^* = M\ddot{q}^d + h$, as in (1.4).

In this case the solution of the minimization problem was trivial, thus applying (1.4) would have been easier, however the formulation (1.11) leaves room for adding other constraints that until were impossible to consider.

Problem (1.11) is usually referred as a *quadratic program* (QP) since the optimization problem has a quadratic cost.

Actuator effort bounds Real actuators can't exert infinitely high forces, but their actions present a saturation limit that acts as a bound. Considering a symmetric behavior on the forces that can be generated, the actuator effort can be bounded as

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

Furthermore while dealing with electrical actuators, the motor drivers can handle just a limited amount of current, thus limiting the overall torque that can be generated; this limitation can be described by the following inequality constraint:

$$-i^{max} \leq K_i \tau \leq i^{max}$$

Joint velocity bounds Actuators also present saturated joint velocities \dot{q}^{max} . We observe that in (1.11) the joint velocity \dot{q} is not a minimization argument, but is a value that's estimated by the robot itself. To bound the velocity we can use an integral constraint; since we can't act on the current velocity, we might minimize in a way that at the next time-step the velocity doesn't exceed the saturation limit. Exploiting the Euler approximation, it holds that $\dot{q}(t + \Delta t) = \dot{q}(t) + \Delta t \ddot{q}$, thus a constraint that we might add to take into account for velocities is

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

Joint position bounds Actuators also have bounded joint position boundaries, i.e. q is constrained to lie in a interval $[q^{min}, q^{max}]$. To satisfy such constraint we might be tempted to exploit the same trick used for the joint velocity saturation, by so computing q at the next time-step by integrating \dot{q} , reaching a formulation of the form

$$q^{min} \leq q + \Delta t \dot{q} + \frac{1}{2} \Delta t^2 \ddot{q} \leq q^{max} \quad (\circ) \quad (1.12)$$

Even if this solution theoretically works, in practice it leads to *unfeasible QPs*: in order not to exceed the joint position saturation at high velocity, the controller is required to set a "high" acceleration \ddot{q} , conflicting with other boundaries set on the actuator efforts.

A heuristic that can be used to overcome this limitation is by exploiting the same constraint (1.12), but using a wider time-step Δt in order to "predict further in the future" the behavior of the system.

QP problem With the discussions said so far, the optimal control problem in (1.11) can be extended in order to embed all bounds described so far, reaching the formulation

$$\begin{aligned}
 \min_{\tau, \ddot{q}} \quad & \frac{1}{2} \|\ddot{q} - \ddot{q}^d\|^2 = \frac{1}{2} (\ddot{q} - \ddot{q}^d)^\top (\ddot{q} - \ddot{q}^d) \\
 \text{sub. to:} \quad & M\ddot{q} + h = \tau \\
 & -\tau^{\max} \leq \tau \leq \tau^{\max} \\
 & -i^{\max} \leq K_i \tau \leq i^{\max} \\
 & -\dot{q}^{\max} \leq \dot{q} + \Delta t \ddot{q} \leq \dot{q}^{\max} \\
 & q^{\min} \leq q + \Delta T \dot{q} + \frac{1}{2} \Delta T^2 \ddot{q} \leq q^{\max}
 \end{aligned} \tag{1.12}$$

Such OCP is harder than (1.11) and the solution is usually non-trivial, for this reason QP solvers numerical methods have been developed. Since the cost is quadratic we can exploit *convex optimization* that will be revised next.

1.4.2 Review of convex optimization

In general not all optimization problems have a solution; moreover for control purposes the execution time of numerical algorithms is fundamental since we want to control the system with less delay as possible.

Taxonomy of convex optimization problems Optimization problems (1.11) and (1.12) are particular minimization referred usually as *least square problems (LSP)* that are characterized by having linear constraint (of the form $Ax \leq b$ or $Ax = b$) and a cost function that is a squared norm of a linear function, i.e. in the form $\|Ax + b\|^2$.

LSPs are a subclass of *quadratic programs (QPs)* optimization problems that are still characterized by linear constraints and a *convex quadratic cost* in the form

$$\frac{1}{2} x^\top H x + g^\top x$$

with H that's a semi-positive definite matrix.

Numerical solvers are usually developed for QP problems, but they can be used also to solve least square programs (still "ad hoc" solvers might find the solution more efficiently). In general LSPs and QPs are *convex optimization problems* that can be solved in low time with off-the-shelf softwares.

Least square programs LSPs are characterized by a cost function of the form

$$f(x) = \frac{1}{2} \|Ax - b\|^2 = \frac{1}{2} (Ax - b)^\top (Ax - b) = \frac{1}{2} x^\top A^\top A x - b^\top A x + \frac{1}{2} b^\top b \tag{1.13}$$

Calling the product $A^\top A$ as H that by definition is positive definite, then the similarity of LSPs with QPs is clear.

Since the problem is convex then there's only one stationary point x^* of f that's also the global minimum; x^* can be computed by setting to zero the gradient of the cost f , i.e. solving

$$\nabla f = A^\top A x - A^\top b = 0$$

Assuming that $A^\top A$ is invertible the explicit solution of the minimizing argument is

$$x^* = (A^\top A)^{-1} A^\top b \tag{†}$$

where $(A^\top A)^{-1}A^\top$ is the *left pseudo-inverse* of A . If $A^\top A$ is not invertible but conversely AA^\top is, then we can use the *right pseudo-inverse* to obtain the solution

$$x^* = A^\top (AA^\top)^{-1}b \quad (\ddagger)$$

Plugging such optimal solution in (1.13) leads in fact to a zero (minimum) cost:

$$f(x^*) = \|AA^\top (AA^\top)^{-1}b\|^2 = \|Ib - b\|^2 = 0$$

In general the solution of least square minimization problem is determined by means of the *Moore-Penrose pseudo-inverse* A^+ as

$$x^* = \min_x \|Ax - b\|^2 = A^+b \quad (1.14)$$

where

$$A^+ = \begin{cases} A^\top (AA^\top)^{-1} & \text{if } A \text{ is a full-row rank matrix} \\ (A^\top A)^{-1}A^\top & \text{if } A \text{ is a full-column rank matrix} \end{cases} \quad (1.15)$$

We observe now that both solutions (†) and (‡) have been obtained considering the respective definition of the Moore-Penrose pseudo-inverse based on the condition of the matrix A (leading to the invertibility of $A^\top A$). Moreover, if A is neither full-row and columns rank, still A^+ can be computed by means of the singular value decomposition, thus an optimal solution (1.14) always exists for (1.13).

1.4.3 Cartesian space control

We can use quadratic programs also to solve control problem in cartesian space. Recalling the desired cartesian acceleration $\ddot{x}^d = \ddot{x}^{ref} + K_p(x^{ref} - x) + K_d(\dot{x}^{ref} - \dot{x})$ defined in (1.7), we can formulate the control problem as the following optimization:

$$\begin{aligned} \min_{\tau, \ddot{q}} \quad & \frac{1}{2} \|\ddot{x} - \ddot{x}^d\|^2 \\ \text{sub. to:} \quad & M\ddot{q} + h = \tau \end{aligned}$$

The main issue with this formulation is that the cost function lacks the dependence of the *decision variables* that are just τ and \ddot{q} ; in this case we have simply to add a new constraint relating \ddot{x} with \ddot{q} , i.e.

$$\begin{aligned} \min_{\tau, \ddot{q}} \quad & \frac{1}{2} \|\ddot{x} - \ddot{x}^d\|^2 \\ \text{sub. to:} \quad & M\ddot{q} + h = \tau \\ & \ddot{x} = J\ddot{q} + \dot{J}\dot{q} \end{aligned} \quad (1.16)$$

To improve the numerical performance of the algorithm we can remove the lastly added constraint and explicitly substitute it in the cost function, reaching the following OCP:

$$\begin{aligned} \min_{\tau, \ddot{q}} \quad & \frac{1}{2} \|J\ddot{q} + \dot{J}\dot{q} - \ddot{x}^d\|^2 \\ \text{sub. to:} \quad & M\ddot{q} + h = \tau \end{aligned} \quad (1.17)$$

Exploiting (1.14) the optimal joint accelerations and torques are computed simply as

$$\ddot{q}^* = J^+ (\ddot{x}^d - \dot{J}\dot{q}) \quad \tau^* = M\ddot{q}^* + h \quad (\Delta)$$

where in this case the coefficient b is made by $\ddot{x}^d - \dot{J}\dot{q}$ (since all this quantities are known beforehand).

TODO: confrontare risultato con quello ottenuto in precedenza (pseudoinversa pesata)

1.4.4 Task space control

In this section we'll further generalize the concept of control in such a way that the target is not just a joint configuration or the end effector in cartesian coordinate, but a *task*, a *control objective* that the robot should achieve.

Each task is usually described by an *error function* $e(\cdot)$ that needs to be minimized; in all previous seen cases the error was just the difference between the current configuration and a reference one, i.e.

$$e(x, u, t) = y(x, u) - y^{ref}(t) \quad (1.18)$$

where $x = (q, \dot{q})$ are in this case the *states* of the system and $u = \tau$ the inputs (using a more general notation).

Since QP problems requires an error function e that must be affine in the decision variable \ddot{q} , a transformation is required (such the one to solve the controls in cartesian space); minimizing directly e is not possible as in general it depends only on q and \dot{q} , not \ddot{q} explicitly.

To reach this goal we must impose the following two conditions on e (or on it's derivative \dot{e}, \ddot{e}) for which:

1. $\lim_{t \rightarrow \infty} e(t) = 0$;
2. the dynamics of e is affine in the decision variables \ddot{q}, u .

Velocity task function Let's consider a case for which the error function depends just on the joint velocity (and eventually on time), so in the form $e(x, u, y) = e(\dot{q}, t) = y(\dot{q}) - y^{ref}(t)$; differentiating in time the error function provides

$$\dot{e}(\dot{q}, t) = \dot{y} - \dot{y}^{ref} = \frac{\partial y}{\partial \dot{q}} \frac{d\dot{q}}{dt} - \dot{y}^{ref} = J\ddot{q} - \dot{y}^{ref}$$

This expression clearly satisfies the second requirement of being affine with respect to \ddot{q} ; to ensure the second condition we can set \dot{e} as proportional to the current error:

$$\dot{e} = -K_p e$$

If K_p is positive definite the linear dynamics is asymptotically stable, satisfying the second condition. The final cost function that can be used to solve this problem is so

$$\begin{aligned} J\ddot{q} - \dot{y}^{ref} &= -K_p e & \text{with } K_p > 0 \\ J\ddot{q} - \dot{y}^{ref} + K_p e &= 0 \end{aligned} \quad (1.19)$$

Position task function Considering now a task depending just on the joint's position in the form $e(x, u, t) = e(q, t) = y(q) - y^{ref}(t)$, then an affine function in \ddot{q} can be obtained considering 2 differentiation in time as follows:

$$\xrightarrow{d/dt} \quad \dot{e} = J\dot{q} - \dot{y}^{ref} \quad \xrightarrow{d/dt} \quad \ddot{e} = J\ddot{q} + \dot{J}\dot{q} - \ddot{y}^{ref} \quad (\circ)$$

In this way the second requirement for e is required. To study the first condition we need to reduce the dynamic to a first order exploiting the variable $z = (e, \dot{e})$, and considering for this case a proportional-derivative controller the linear system boils down to

$$\dot{z} = \begin{pmatrix} \dot{e} \\ \ddot{e} \end{pmatrix} = \begin{bmatrix} 0 & I \\ -K_p & -K_d \end{bmatrix} \begin{pmatrix} e \\ \dot{e} \end{pmatrix} = Az$$

In this case matrices $K_p, K_d > 0$ must be chosen in such a way that $\ddot{e} = -K_p e - K_d \dot{e}$ is asymptotically stable (A must be Hurwitz). Since now also the first condition is met, we can rewrite (o) as

$$\begin{aligned} J\ddot{q} + \dot{J}\dot{q} - \ddot{y}^{ref} &= -K_p e - K_d \dot{e} \\ J\ddot{q} + \dot{J}\dot{q} + K_p e + K_d \dot{e} &= 0 \end{aligned}$$

Linear theory suggests that such dynamic is exponentially stable, however in this case no constraints have been considered (but they can limit the actual stabilization performance).

Input task function The only way to build a task function $e(u, t)$ that depends just on the input torques $\tau = u$ is by ensuring that e is actually affine in u (that's a decision variable), i.e. must be of the form

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

Task summary As just presented we might have 3 main kind of task function, respectively:

- one that's affine in the input $\tau = u$;
- two that are non-linear functions of the states x but for which an affine function of \ddot{q} can be achieved by means of one or two differentiation in time for velocity and position tasks respectively.

In general a combination of multiple tasks can be collapsed to a single task exploiting an "augmented" state vector $z = (\ddot{q}, u)$ that determines the following affine function:

$$g(z) = \begin{bmatrix} A_x & A_u \end{bmatrix} \begin{pmatrix} \ddot{q} \\ u \end{pmatrix} - b = Az - b$$

QP-based control The *task-space inverse dynamics (TSID)*, also referred as *QP-based control*, solves optimization problems in the for

$$\begin{aligned} \min_{z=(\ddot{q}, \tau)} \quad & \|Az - b\|^2 \\ \text{sub. to:} \quad & \begin{bmatrix} M & -I \end{bmatrix} z = -h \\ & \text{other constraints} \end{aligned} \tag{1.20}$$

This is a generalization of the QP joint space control in (1.12), choosing $A = \begin{bmatrix} I & 0 \end{bmatrix}$ and $b = \ddot{q}^d$, or the cartesian space control in (1.17), choosing $A = \begin{bmatrix} J & 0 \end{bmatrix}$ and $b = \dot{x}^d - \dot{J}\dot{q}$.

TODO: underactuated systems, rigid contacts, multi-task control

Chapter 2

Optimal Control

Optimal control problems (OCPs) are constrained minimization problems in the form

$$\begin{aligned}
 & \min_{x(\cdot), u(\cdot)} \int_0^T \ell(x(t), u(t), t) dt + \ell_f(x(T)) \\
 \text{sub. to: } & \dot{x}(t) = f(x(t), u(t), t) & \forall t \\
 & g(x(t), u(t), t) \leq 0 & \forall t \\
 & x(0) = x_0
 \end{aligned} \tag{2.1}$$

where $\ell(\cdot)$ is the *running cost*, $\ell_f(\cdot)$ is the *terminal (final) cost*, f is the *system's dynamic*, g are the *path constraints* and x_0 is the *initial condition*. Moreover x are the so called *states* (that for manipulator are simply joint coordinates and velocities) while u are the *inputs* (joint torques).

For ease of notation the time-dependence of both x and u will now be dropped.

Optimal control problems from a first look are very similar to quadratic programs like (1.12) at 10, however the main difference in this case is that both $x(t)$ and $u(t)$ are *trajectories* that needs to be identified. In QPs we assumed in fact that a desired motion $y^{ref}(t)$ was given, while OCP solves a problem with theoretically infinitely many degrees of freedom that build the optimal trajectory y^{ref} itself.

Example: simple pendulum Let's consider the classical example of a simple pendulum that we want to "swing up"; choosing q in such a way that for $q = 0$ the pendulum is in the upright configuration, then the simplest OCP that we can formulate for this problem is

$$\begin{aligned}
 & \min_{x, \tau} q(T) = 0 \\
 \text{sub. to: } & I\ddot{q} = \tau + mg \sin(q)
 \end{aligned}$$

where the lonely constraint is just the system's dynamic. In this case we used just a terminal cost to enforce to reach the desired motion at the end of the time horizon. Even if this problem lead to the desired configuration, no "rules" on how the system's behavior are set, implying that there's no penalization for a fast convergence or high velocities.

One way to reach the desired configuration in less time possible is to consider an integral cost that adds when we are out of target, considering the following OCP:

$$\begin{aligned}
 & \min_{x, \tau} \int_0^T q^2 dt \\
 \text{sub. to: } & \begin{pmatrix} \dot{q} \\ \ddot{q} \end{pmatrix} = \begin{pmatrix} \frac{dq}{dt} \\ I^{-1}(\tau + mg \sin q) \end{pmatrix}
 \end{aligned}$$

where in this case the dynamic has been explicitly written in state-space after a transformation to first order given the state $x = (q, \dot{q})$.

Solving this problem would lead to the generation of infinitely high torques and joint velocities in order to reach the target configuration in little time as possible to reduce the integral value.

To avoid this problem we might add as integral costs also the square of the joint velocities and the torque, so the solver needs to find the perfect trade-off to reach the desired configuration in low time but without generating high velocities or torques:

$$\begin{aligned} \min_{x, \tau} \quad & \int_0^T (q^2 + \dot{q}^2 + \tau^2) dt + 10^3 q(T) \\ \text{sub. to:} \quad & I\ddot{q} = \tau + mg \sin(q) \end{aligned}$$

In this case it has been added also a heavy-weighted terminal cost to ensure that the system actually reach the desired configuration (since the trade-off in the integral might lead to a non-requested solution).

Using integral cost to bound velocities and torques is an easy way to solve the problem, but similar results can be obtained by setting hard constraints on some variables, i.e.

$$\begin{aligned} \min_{x, \tau} \quad & \int_0^T q^2 dt \\ \text{sub. to:} \quad & I\ddot{q} = \tau + mg \sin(q) \\ & |\tau| \leq \tau^{max} \\ & |\dot{q}| \leq \dot{q}^{max} \end{aligned}$$

Optimal control method families Tables/ocpfamilies.tex Once an optimal control problem (2.1) has been states, different families of solvers, described in table ??, can be used to determine the optimal trajectories. In particular there are two "orthogonal" categorizations:

- *continuous-time* solvers are searching for the continuous-time solution of the problem itself, while *discrete-time* methods rely on a discretized (approximated) version of the same problem;
- *global* solvers are searching the global optimum for the problem, while *local* are starting from a guess and reach the closes minima point; the latter of course is faster at runtime.

For robotic applications mainly we use *direct method* solvers; still we start introducing the other ones as they present critical aspect that are useful in understanding direct methods and the next chapter on *Reinforcement Learning*.

2.1 Dynamic programming

Dynamic programming (DP) methods are based on the *Bellman's optimality principle* for which

« given the optimal control starting from A at time $t = 0$ and reaching B at $t = T$,
then given any condition \bar{x} in the optimal trajectory for $\bar{t} \in (0, T)$ then the the optimal
solution from \bar{x} to B is still the same. »

This is a simple, yet powerful, concept after which all minimization algorithms have been developed on. In particular for what concerns dynamic programming problem (2.1) is discretized in N time-steps and rewritten as

$$\begin{aligned} \min_{X,U} \quad & \sum_{i=0}^{N-1} \ell(x_i, u_i) \\ \text{sub. to:} \quad & x_{i+1} = f(x_i, u_i) \quad \forall i = 0, \dots, N-1 \end{aligned} \quad (2.2)$$

where¹ $X \in \mathbb{R}^{n \times N}$ and $U \in \mathbb{R}^{m \times N}$ are the matrices containing all the states and inputs at the discrete-time steps, i.e.

$$X = \begin{bmatrix} x_1 & \dots & x_N \end{bmatrix} \quad U = \begin{bmatrix} u_0 & \dots & u_{N-1} \end{bmatrix} \quad (2.3)$$

We point out that in (2.2) the terminal cost has been removed just to simplify the calculation of the following proof, but can be added at any time.

Furthermore we state that dynamic programming can be used only to solve *unconstrained minimization problems* (or, better, the lonely constraint allowed is the one of the systems dynamic).

Solution with QP After discretization of (2.1), (2.2) can be seen as a "simple" optimization problem in the variables X, U , similarly to (1.11) at page 9.

QP solvers can so be used to obtain the global optimum of the problem, however (2.2) suffers of dimensionality: increasing the number of states/inputs or improving the discretization N heavily affects the computational load. Problems might become too large to be handled just with very simple scenarios, so solving (2.2) with QPs is not a feasible solution.

After discretization OCP (2.2) reduces to a simple optimization problem such (1.11) at page 9; the main issue of using QP solvers for this problem is the high dimensionality of the problem itself that scales badly with both the number of states/inputs and the number of discretization steps.

Solution with the Bellman's principle To numerically solve (2.2), DP solvers exploit the Bellman's optimality principle; considering in fact $M \in \mathbb{N}$ such that $0 < M < N$, then the overall problem can be casted to the following unconstrained minimization:

$$\min_{X,U} \left(\underbrace{\sum_{i=0}^{M-1} \left(\ell(x_i, u_i) + \mathcal{I}(x_{i+1} - f(x_i, u_i)) \right)}_{=C_0(X_{1:M}, U_{0:M-1})} + \underbrace{\sum_{i=M}^{N-1} \left(\ell(x_i, u_i) + \mathcal{I}(x_{i+1} - f(x_i, u_i)) \right)}_{=C_M(x_M, X_{M+1:N}, U_{M:N-1})} \right) \quad (+)$$

In this expression it has been introduced the *indicator function* \mathcal{I} , a mathematical "trick" used to convert the dynamic constraint into a cost; to ensure the condition, the system is infinitely penalized when such constraint isn't satisfied, so by considering

$$\mathcal{I}(x) = \begin{cases} 0 & \text{if } x = 0 \\ \infty & \text{otherwise} \end{cases} \quad (2.4)$$

¹in this case we assume that the state x has n dimensions while inputs u have dimension m .

Algorithm 1 pseudo-code for solving a discretized optimal control problem using dynamic programming.

```

initialize  $V_n(x_N) = \ell_f(x_N)$  ▷ look-up table that needs to be defined  $\forall x_N$ 

for  $i$  from  $N - 1$  to  $0$  do
     $V_i(z) = \min_u (\ell(z, u) + V_{i+1}(f(z, u)))$ 
end for

for  $i$  from  $0$  to  $N - 1$  do ▷ at runtime
     $u_i^*(x) = \min_u (\ell(x, u) + V_{i+1}(f(z, u)))$ 
end for

```

At this point we can rewrite (†) as the sum of two independent minimization problems as follows:

$$\min_{X_{1:M}, U_{0:M-1}} C_0(X_{1:M}, U_{0:M-1}) + \underbrace{\min_{X_{M+1:N}, U_{M:N-1}} C_M(x_M, X_{M+1:N}, U_{M:N-1})}_{=V_M(x_M)} \quad (2.5)$$

Here we observe that the first is a "canonical" minimization, while the second one is parametric in x_M , i.e. the value of the second minimization is affected by the state x_M reached in the optimization of C_0 .

Such second term is usually referred as the *value function* $V_M(x_M)$ and evaluates the optimal cost that needs to be paid from time $t = M$ starting from a (generic) initial condition x_M , or, in other words, the "optimal cost-to-go given x_M as initial condition. Exploiting the definition of the value function we have been able to split the initial problem (†) into two (almost) independent minimization.

In this way (2.5) can be concisely rewritten as

$$V_0(x_0) = \min_{X_{1:M}, U_{0:M-1}} C_0(X_{1:M}, U_{0:M-1}) + V_M(x_m) \quad (2.6)$$

i.e. "the cost that we pay starting from $t = 0$ with x_0 initial condition and behaving optimally".

Choosing now $M = 1$ allow us to rewrite the value function with a *recursive formulation* as follows:

$$V_i(x_i) = \min_{u_i} (\ell(x_i, u_i) + V_{i+1}(f(x_i, u_i))) \quad (2.7)$$

Here we see that the minimization is independent from the current state x_i that is the parameter of the value function, but what needs to be determined is the value of u_i that allows to reduce the overall cost of the current time-step combined with the cost-to-go.

DP algorithm Dynamic programming exploits the recursive definition of the value function (2.7) to solve the optimization problem; as presented in algorithm 1, the main idea is to regard V as a look-up table where for each current state, the optimal cost-to-go is given. Such table is built backward in time, i.e. we start from the terminal cost (that's known given the parametric final state x_N) and, going backward in time, we compute each value in the table in order to minimize the overall cost-to-go given the current configuration. Finally at run-time the optimal controls u^* are chosen in such a way that they minimize the cost-to-go toward the final target.

The main advantage of this algorithm is that it doesn't provide an optimal trajectory, but rather a *feedback policy*: once the look-up table is filled, then for each state the system

reach at run-time the optimal behavior can be automatically determined.

This comes however with a big issue: the parametric minimization itself. This method to work requires the generation of a look-up table that suffers of dimensionality based on the discretization used to describe the states and the inputs as well as on time.

In the simplified case where the systems dynamic is linear and the cost is simply a quadratic function in x and u , then the value function V is convex and a simpler closed-form solution can be obtained, leading to the sub-class problems of the so called (*discrete-time*) *linear quadratic regulators* (LQR), described in more depth in section 2.5 (page 24).

2.2 Hamilton-Jacobi-Bellman

As the name suggest, the *Hamilton-Jacobi-Bellman* (HJB) method is based on the Bellman's optimality principle and solves the optimal trajectory directly in continuous-time (while in dynamic programming a discretization was required). The underlying ides of HJB is to take DP and push the discretization step toward an infinitesimal value.

Considering the value function (2.7) as been obtained from a δ time-step discretization, then $V(\cdot)$ depends just on the current time t_i and the parametric state

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

where $\ell_d = \ell(x, u)\delta$ is the discretized cost. Pushing δ towards 0, then the value function can be approximated by it's first Taylor series expansion as

$$\begin{aligned} V(t_i, x) &= \min_u \ell_d(x, u) + V(t_i, x) + \frac{\partial V}{\partial t} \delta + \frac{\partial V}{\partial x} \frac{dx}{dt} \delta \\ 0 &= \min_u \ell(x, u)\delta + \dot{V}\delta + \nabla_x V(t_i, x)f(x, u)\delta \end{aligned}$$

that divided by the common term δ allows us to reach the *Hamilton-Jacobi-Bellman equation*

$$-\dot{V} = \min_u \ell(x, u) + \nabla_x V f(x, u) \quad (2.8)$$

This expression models the Bellman's optimality principle in a continuous-time setting; the main issue is that (2.8) is not a ordinary differential equation (ODE), but a partial differential equation (PDE) since it depends on the differentiation of the unknown function V with respect to both time t and states x .

Still if the dynamic is linear and the cost is quadratic, the problem simplifies to the solution of a continuous-time LQR (see sec. 2.5).

HJB solvers still exploits the backward integration from the final cost $V(T, x_T) = \ell_f(x_T)$ to determine the optimal feedback policy, but the problem is now harder to numerically integrate.

2.3 Pontryagin minimum principle

Methods based on the *Pontryagin minimum principle* starts from the HJB equation (2.8) and is based on the observation that the optimal control u^* depends just on the gradient $\nabla_x V$ of the value function, and not by $V(\cdot)$ itself. Introducing the *adjoint variable* $\lambda(t)$ defined as $\left(\nabla_x V(t, x(t))\right)^\top$, then the optimal control can be computed as

$$u^*(t, x, \lambda) = \min_u \underbrace{\ell(x, u) + \lambda^\top f(x, u)} \quad (2.9)$$

TODO: Rivedere bene il PMP e capire come funziona e come l'equazione è derivata

2.4 Direct methods

Direct methods (DM) is a family of optimal control solvers highly used in robotic applications; they mainly search for local optima trajectories of a discretized optimal control problem from the continuous-time one.

As seen so far, the main issue associated to the solution of optimal control problems is related to the high dimensionality (due to the time-continuity of the dynamics). The main idea in direct methods is to discretize the time axis and define a parametrization (such a piecewise constant function) for the inputs u inside such interval.

After discretization the OCP becomes an optimization problem that can be solved by means of *non-linear programming* (NLP) solvers. Direct methods works as follows:

- once the time-axis is discretized, we parametrize the trajectories of the decision variables of the problem (inputs and/or states); for this purpose we can use piece-wise constant functions or polynomials (or any other function);
- we enforce the constraints (due to both dynamic and path constraints) on the defined time-grid of times $t_0 < t_1 < \dots < t_N$.

Intuitively this is the simplest algorithm that we might come up with (compared also with the solution previously discussed), however we have to deal with issue of *over-discretization* (leading to a numerical cost growth) or *under-discretization* (smaller dimensionality but at runtime some constraints might be violated since they are enforced only in few discrete time-steps).

With respect to all other presented solution, direct methods are overall faster and for this very reason are used in robotic also in feedback loop; for this last point it's crucial to develop algorithms that are able to solve OCPs fast enough to keep up with the robot's dynamic.

There are mainly 3 algorithms for solving OCP problems with direct methods:

- *single shooting* (or *sequential*): the discretization is performed on the controls trajectory $u(t)$ that's so is the lonely decision variable; the state trajectory $x(t)$ is instead obtained by integration of the dynamics (that so can be removed from the constraints since the condition is always satisfied by construction);
- *collocation* (or *simultaneous*): the discretization is performed on both inputs u and states x ; in this case the dynamic is used to enforce the constraint between such trajectories;
- *multiple shooting*: this can be regarded as a combination of both single shooting and collocation since the states x are discretized on a coarser time-grid (with respect to the inputs u). Intermediate values for the states are compute by integration. Such method can be regarded as a sequence of single shootings.

The choice of the method that should be used heavily depends on the NLP solver chosen; collocation leads to a problem with bigger matrices that are however sparser, so some software might exploit such aspect; conversely, single shooting have a lower system's dimensionality but require a good integrator to accurately compute the state trajectories.

2.4.1 Single shooting

In *single shooting* we discretize only the control $u(t)$ on a fixed time-grid $0 = t_0 < t_1 < \dots < t_N = t_f$; for sake of simplicity u is usually parametrized as piece-wise within the

integration time-step. Called y the decision variable (that describes the parametrization of u) and once the system can be integrated given such parameters, then the NLP problem that needs to be solved is of the form

$$\begin{aligned} \min_y \quad & \int_0^{t_f} \ell(x(t_i, y), u(t_i, y), t) dt + \ell_f(x(t_f, y)) \\ \text{sub. to:} \quad & g(x(t_i, y), u(t_i, y), t_i) \quad \forall i \in [0, N-1] \end{aligned} \quad (2.10)$$

where in this case g is the *discretized path constraints*. The dynamic is not included since it's knowledge has been used to compute the state trajectory $x(t_i, y)$. Regarding the running cost as a time-dependent function of the form

$$c(t) = \int_0^T \ell(\cdot) dt$$

then the related computation can be reduced to the integration of the following ordinary differential equation:

$$\begin{cases} \dot{c}(t) = \ell(\cdot) \\ c(0) = 0 \end{cases}$$

Calling $\bar{x} = (x, c)$ the *augmented state vector*, then the overall ordinary differential equation that must be integrated to solve the OCP (2.10) is

$$\dot{\bar{x}} = \begin{pmatrix} f(x, u, t) \\ \ell(x, u) \end{pmatrix} \quad (2.11)$$

Once this expression is integrated by means of an integrator and its derivative with respect to the decision variable u is computed, then NLP solvers might be able to converge to the local minima by exploiting a gradient descent.

2.4.2 Numerical integration

Given a ordinary differential equation in the form $\dot{x} = f(x, t)$ subjected to an initial condition $x(0) = x_0$, then numerical integrators deals with the computation of $x(t)$ for all times $t \in [0, T]$; f in this case is not written as dependent from the input since when we evaluate the dynamic u has a fixed and known value, i.e. integration is not performed parametrically.

As long as f is *Lipshitz continuous* (a stronger condition then requiring a continuous first derivative of f), then it is proven that the ODE has a unique solution in the neighborhood of the initial point. Based on this assumption we can numerically integrate by using different algorithms, such the Euler method.

Euler method Recalling the formal definition of the incremental ratio

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

the Euler method uses the approximation for a "sufficiently small" h to compute the next state $x(t+h)$ given the initial condition $x(t)$ simply by reverting the equation:

$$x(t+h) \approx x(t) + h f(x, t) \quad (2.12)$$

This method to be accurate requires h to be "very small": accuracy is so inversely proportional to the numerical time complexity (to be more accurate we need lower h , increasing the number of required evaluation to cover the same time span). As we'll see better later, this is called a " 1^{st} order method", so with a minimum consistency; in general higher the order, more favorable is the trade-off between time complexity and accuracy.

Midpoint method The midpoint method is another integration algorithm that falls in the category of the 2^{nd} order. Once we compute

$$K_1 = f(x(t), t)$$

and

$$K_2 = f\left(x(t) + \frac{h}{2}K_1, t + \frac{h}{2}\right)$$

then the state at the next time-step is approximated by the formula

$$x(t+h) \approx x(t) + hK_2 \quad (2.13)$$

The idea behind this formula is that first we compute the slope K_1 and use half of a Euler step to approximate a "mean" slope K_2 inside the discretization time-step. Given the same value of h , the midpoint method is usually more accurate than Euler, but it also requires a double value of function f evaluation (2 vs 1).

Properties of numerical integration Given an ODE with exact trajectory solution $x(t)$ and an integrator whose output is described by $\hat{x}(t, t_0, x(t_0))$ where $t_0, x(t_0)$ are respectively the time and initial condition for starting the integration, then we define the *local error* $e(\cdot)$ as

$$e(t) = x(t) - \hat{x}(t, t-h, x(t-h)) \quad (2.14)$$

while the *global error* $E(\cdot)$ is simply

$$E(t) = x(t) - \hat{x}(t, t_0, x(t_0)) \quad (2.15)$$

As the name suggests, the global error evaluates the "distance" between the truth value and the integrated given the same initial condition, while the local error measures the distance considering as initial true condition just one time-step ahead.

With this definitions any "good" numerical integration scheme must satisfy the following 3 conditions:

- i) *convergence*, i.e. $\lim_{h \rightarrow 0} E(t) = 0$;
- ii) *consistency order*: $\lim_{h \rightarrow 0} e(t) = \mathcal{O}(h^{p+1})$ where $p > 0$ is the *order* of the numerical integration; such parameter is key since it rules the asymptotic convergence of the system;
- iii) *stability*. This concept is now well defined, but intuitively we can say that the global error E must be bounded for $t \rightarrow \infty$; this requirement is important especially for so called "stiff" ODEs.

Consistency order of the explicit Euler method Rewriting the local error (2.14) in a more readable index notation $e(t) = x_{n+1} - \hat{x}_{n+1}(x_n) = e_{n+1}$, the definition of \hat{x}_{n+1} given by the Euler method is simply $x_n - hf(x_n, t_n)$. Considering now the Taylor series expansion of the true solution around x_n that's

$$x_{n+1} = x_n + h\dot{x}_n + \mathcal{O}(h^2)$$

then we can show that

$$\begin{aligned} e_{n+1} &= x_{n+1} - \hat{x}_{n+1} = x_n + h\dot{x}_n + \mathcal{O}(h^2) - x_n - hf(x_n, t_n) \\ &= \mathcal{O}(h^2) \end{aligned}$$

proving that the consistency order is actually $p = 1$.

Runge-Kutta methods Both Euler and midpoint methods can be regarded as "instantiations" of the more general definition of the *Runge-Kutta (RK) methods*, one-step integration algorithms of the form

$$x_{n+1} = x_n + h \sum_{i=0}^q b_i K_i \quad (2.16)$$

with

$$K_i = f \left(x_n + h \sum_{j=0}^q a_{ij} K_j, t_n + c_i h \right)$$

where a_{ij}, b_i, c_i are all tabled coefficients; q is the so called *order* of the Runge-Kutta method that in general might differ from the consistency order p . Usually we call *RK4* a Runge-Kutta method of order $q = 4$.

It's true that $p \leq q$, and in particular for $q \geq 5$ the consistency order is always lower than the method's one.

As example, the explicit Euler method is a *RK1* integration scheme with $b_1 = 1, a_{11} = 1$ and $c_1 = 0$; the midpoint integrator is instead a *RK2* scheme.

Every-time it happens that $a_{ij} = 0$ for all $j \geq i$, then the *method* is said *explicit* and the integration can be performed by simple evaluations of the function f (while for *implicit methods* it is required to solve an implicit non-linear problem in x_{t+1} , increasing the cost complexity).

Runge-Kutta methods are unequivocally described by the so called *Butcher tableaux* storing the parameters for the integration in a form:

$$\begin{array}{c|ccc} c_1 & a_{11} & \dots & a_{1q} \\ \vdots & \vdots & \ddots & \\ c_q & a_{q1} & & a_{qq} \\ \hline & b_1 & \dots & b_q \end{array} \quad (2.17)$$

The best trade-off between accuracy and numerical complexity is favorable up to order $q = 4$ (to get a consistency order $p = 5$ it is required at a *RK6* method); one of the most widely used Butcher tableau is the *RK4* explicit method described as

$$\begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6} \end{array}$$

2.4.3 Computation of sensitivities

Optimal control problems, as already mentioned, are usually solved exploiting gradient-based search thus the computation of the derivative is a key concept. Any OCP (2.1) has a cost that needs to be minimized: this value usually depends on the state trajectory $x(t)$ that is computed by integration of the dynamics considering the input $u(t)$.

In this section we will deal with the computation of the *sensitivities*, so the derivative of the result of an integration.

Given the sequence of piecewise constant inputs $u(t) = y_i$ (for any $t \in [t_i, t_{i+1}]$), the states are computed by integration of $\dot{x} = f(x, y_i, t)$ for all $t \in [t_i, t_{i+1}]$. Exploiting the Euler discretization, the cost function can be approximated as

$$c(y) = \int_0^T \ell(x(t), u(t), t) dt + \ell_f(x(T)) \approx \sum_{i=0}^{N-1} \ell(x_i, y_i) h + \ell_f(x_N)$$

In order to apply numerical optimization we can compute the gradient of $c(\cdot)$ with respect to the optimization variable y as

$$\begin{aligned}\frac{dc(y)}{dy} &= \sum_{i=0}^{N-1} \frac{d\ell(x_i, y_i)}{dy} h + \frac{d\ell_f(x_N)}{dy} \\ &= \sum_{i=0}^{N-1} \left(\frac{\partial \ell}{\partial x_i} \frac{dx_i}{dy_i} + \frac{\partial \ell}{\partial y_i} \right) h + \frac{\partial \ell_f}{\partial x_N} \frac{dx_N}{dy_i}\end{aligned}\quad (2.18)$$

Typically the cost function ℓ is relatively simple and its partial derivative can be easily obtained analytically (since we can design ℓ); the main issue in (2.18) is that the derivative dx_i/dy_i depends on the numerical integration of the state.

Since Runge-Kutta methods are single-step integration methods, they can be described using a so called *integration function* ϕ that determines the next step given just the current state and the applied input:

$$x_{i+1} = \phi_i(x_i, y_i)$$

whose derivative is simply

$$\frac{dx_{i+1}}{dy} = \frac{\partial \phi_i}{\partial x_i} \frac{dx_i}{dy} + \frac{\partial \phi_i}{\partial y} \quad (2.19)$$

In this equation we might observe that the derivative of x_{i+1} depends on its previous state x_i ; in this case if we firstly fix $dx_0/dy = 0$ (since the initial state is known and do not depend on the applied input), then each next derivative can be computed integrating forward in time. Furthermore we can see that the matrix

$$\frac{\partial \phi_i}{\partial y} = \begin{bmatrix} \frac{\partial \phi_i}{\partial y_0} & \frac{\partial \phi_i}{\partial y_1} & \cdots & \frac{\partial \phi_i}{\partial y_{N-1}} \end{bmatrix}$$

is sparse, since the next state depends just on the currently applied input (not the others), implying

$$\frac{\partial \phi_i}{\partial y_j} = 0 \quad \forall i \neq j$$

The only element that still needs to be covered in (2.19) is the derivative $\partial \phi_i / \partial x_i$ whose value depends just on the integration scheme used.

Sensitivities of the explicit Euler method Let's consider the Euler method (2.12) defined by the integrator function

$$\phi(x, u) = x + h f(x, u)$$

In this case the sensitivities can be easily computed from the analytical definition as

$$\frac{\partial \phi}{\partial x} = I + h \frac{\partial f}{\partial x} \quad \frac{\partial \phi}{\partial u} = h \frac{\partial f}{\partial u}$$

2.4.4 Collocation

Collocation is another category of direct methods that discretizes also the states $x(t)$ (in this case still with piecewise constant functions for simplicity). Calling

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

the discretized values for the states, then the dynamic $\dot{x} = f(x, u)$ can be approximated using the Euler method as

$$\underbrace{\frac{\overbrace{s_{i+1} - s_i}^{\approx \dot{x}}}{\overbrace{t_{i+1} - t_i}^{\approx \Delta t}} - f\left(\frac{\overbrace{s_{i+1} + s_i}^{\approx x}}{2}, y_i\right)}_{=c_i(s_i, s_{i+1}, y_i)} = 0 \quad \forall i \quad (2.20)$$

Considering this discretized problem, the cost integral for a given time-step can be regarded as

$$\int_{t_i}^{t_{i+1}} \ell(x(t), u(t)) dt \approx \ell\left(\frac{s_i + s_{i+1}}{2}, y_i\right) (t_{i+1} - t_i) = \ell_i(s_i, s_{i+1}, y_i)$$

Based on the presented discretization of the OCP (note that this is just one of the infinitely many that we might come up with), then an approximation of the problem that can be solved using NLP software is

$$\begin{aligned} \min_{s, y} \quad & \sum_{i=0}^{N-1} \ell_i(s_i, s_{i+1}, y_i) + \ell_f(s_N) \\ \text{sub. to:} \quad & s_0 - x_0 = 0 && : \text{initial condition} \\ & c_i(s_i, s_{i+1}, y_i) = 0 && : \text{discretized dynamics} \\ & g_i(s_i, y_i, t_i) = 0 && : \text{discretized constraints} \end{aligned} \quad (2.21)$$

This problem has an higher number of decision variables with respect to single shooting (since we have to add the cardinality of y with the one of s), but comes with the advantage of having a *sparser problem*: each constraint in fact depends just con the current position and, at worst, a neighbor state, i.e.

$$\frac{\partial^2 c_i}{\partial s_i \partial s_{i \pm j}} = 0 \quad \forall j > 1$$

2.4.5 Multiple shooting

Multiple shooting is a direct method that can be regarded as a combination of single shooting and collocation; from the latter it inherits the fact that the states are discretized, but in this case on a coarser time-grid with respect to the one of the controls. Inside the "bigger" state time-chunk, the dynamic is integrated to have a smooth trajectory.

By doing so no dynamic constraint should be considered in the OCP problem, but we have to ensure the state continuity $x_i(t_i) = s_i$ at each boundary of the state time-axis.

TODO: migliorare spiegazione e aggiungere

2.5 Linear quadratic regulator

Given a linear discrete-time system described as

$$x_{t+1} = Ax_t + Bu_t \quad (2.22)$$

we want to design a set of controls u_i for which:

- the resulting state trajectory $x(t)$ is "small", staying close to zero leading to a good regulation (subtly implying that $x(t)$ is modeling the error from a reference target);

- the control sequence u_i is also "small" in order to reduce the required actuator effort.

Intuitively such conditions are in contrast since one can't have a good regulation performance without asking a minimum amount of effort, so a trade-off between the two requirements is necessary. *Linear quadratic regulators (LQR)* deals with this by solving the following unconstrained minimization problem with costs in quadratic form:

$$\min_{U, X} \mathcal{J}(U, X) = \sum_{i=0}^{N-1} \left(x_i^\top Q x_i + u_i^\top R u_i \right) + x_N^\top Q_f x_N \quad (2.23)$$

$$\text{sub. to: } x_{t+1} = Ax_t + Bu_t$$

where $Q, Q_f \geq 0$ and $R > 0$ are all symmetric matrices weighting respectively the cost of the states and the controls. Choosing Q "big" (with respect to R) tells that is more important having a good system regulation, while if R is dominant we want to rather bound the actuator effort instead of the regulation performance.

We can observe that (2.23) has a quadratic cost and is subjected to a linear constraint, so it's a *convex QP problem* as seen in sec. 1.4.2 (page 10). Describing the whole state dynamic in a vector $X = (x_0, \dots, x_N)$ it can be shown that

$$\underbrace{\begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_N \end{pmatrix}}_X = \underbrace{\begin{bmatrix} 0 & & & \\ B & 0 & & \\ AB & B & 0 & \\ \vdots & \vdots & \ddots & \\ A^{N-1}B & A^{N-2}B & B & 0 \end{bmatrix}}_G \underbrace{\begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ \vdots \\ u_{N-1} \end{pmatrix}}_U + \underbrace{\begin{bmatrix} I \\ A \\ A^2 \\ \vdots \\ A^N \end{bmatrix}}_H x_0 \quad (\dagger)$$

where the matrix G has been obtained recursively applying (2.22) and collecting all elements in a matrix form. With this notation (2.23) can be compactly rewritten as

$$\min_{U, X} \mathcal{J}(U, X) = \|\bar{Q}X\|^2 + \|\bar{R}U\|^2 \quad (2.24)$$

$$\text{sub. to: } x_{t+1} = Ax_t + Bu_t$$

where

$$\bar{Q} = \text{diag}\{Q^{1/2}, \dots, Q^{1/2}, Q_f^{1/2}\} \quad \bar{R} = \text{diag}\{R^{1/2}, \dots, R^{1/2}\}$$

Moreover we can substitute the definition of the dynamics (\dagger) and restate the problem as the following unconstrained minimization:

$$\min_U \mathcal{J}(U) = \|\bar{Q}(GU + Hx_0)\|^2 + \|\bar{R}U\|^2 = \left\| \begin{bmatrix} \bar{Q}G \\ \bar{R} \end{bmatrix} U + \begin{bmatrix} \bar{Q}Hx_0 \\ 0 \end{bmatrix} \right\|^2 \quad (2.25)$$

Recalling (1.14), page 11, the solution of this problem can be obtained by means of the Moore-Penrose pseudo-inverse; since the matrix that needs to be inverted has size $N(n+m) \times Nm$, the overall numerical complexity for QP solvers is $\mathcal{O}(N^3nm^3)$.

Solving such problem with QP solvers is highly inefficient, so for this reason dynamic programming (see sec. 2.1, page 15) is usually used to reduce the computational cost to $\mathcal{O}(Nn^3)$. The idea is that if we are able, given the value function $V_{t+1}(z)$, to compute $V_t(z)$, then we can start from V_N and go backward in time to compute the optimal trajectory.

We start off defining the value function as

$$V_t(z) = \min_{w, U_{t+1}} z^\top Qz + w^\top R w + \sum_{k=t+1}^{N-1} \left(x_k^\top Q x_k + u_k^\top R u_k \right) + x_N^\top Q_f x_N$$

$$\text{sub. to: } x_{k+1} = Ax_k + Bu_k$$

$$u_t = w$$

and we observe the recursive definition

$$V_t(z) = \min_w \underbrace{z^\top Qz + w^\top R w}_{\ell(z,u)} + V_{t+1}(\underbrace{Az + Bw}_{f(z,w)}) \quad (2.26)$$

Observing the similarity with (2.5), ℓ is the cost that we pay at the current time considering the input $u_t = w$, while the second term is the cost-to-go given the state that we land in. The main advantage now in LQR (with respect to the general case of dynamic programming) is the linear dynamic that will help us reach an easier solution.

By construction the value function V_{t+1} can always be regarded as quadratic and simplified to the notation

$$V_{t+1}(z) = z^\top P_{t+1} z \quad \text{with } P_{t+1} = P_{t+1}^\top \geq 0, \quad P_N = Q_f$$

Observed that this definition surely holds for the last time-step, we can elaborate (2.26) as

$$\begin{aligned} V_t(z) &= \min_w (z^\top Qz + w^\top R w + V_{t+1}(Az + Bw)) \\ &= z^\top Qz + \min_w (w^\top R w + (Az + Bw)^\top P_{t+1} (Az + Bw)) \\ &= z^\top Qz + \min_w (w^\top (R + B^\top P_{t+1} B) w + 2z^\top A^\top P_{t+1} B w) + z^\top A^\top P_{t+1} A z \\ &= z^\top (Q + A^\top P_{t+1} A) z + \min_w \underbrace{(w^\top H w + 2g^\top w)}_{f(w)} \end{aligned}$$

where to simplify the notation it has been used $H = R + B^\top P_{t+1} B$ and $g^\top = z^\top A^\top P_{t+1} B$. From the last equality we clearly see that the value function always have a fixed cost based on the parameter z , while the second term can actually be minimized in the control input w . Being f a convex quadratic function we can exploit 1.14 to determine the optimal control input as

$$w^* = -H^{-1}g = -(R + B^\top P_{t+1} B)^{-1} B^\top P_{t+1} A z = K_t z \quad (2.27)$$

Plugging back the optimal solution inside the definition of the value function provides

$$\begin{aligned} V_t(z) &= z^\top (Q + A^\top P_{t+1} A) z + g^\top H^{-1} H H^{-1} - 2g^\top H^{-1} g \\ &= z^\top \underbrace{(Q A^\top P_{t+1} A - A^\top P_{t+1} B H^{-1} B^\top P_{t+1} A)}_{P_t = P_{t+1}^\top \geq 0} z \end{aligned} \quad (2.28)$$

With this we proved that if V_{t+1} is a quadratic function, then also V_t must be so; furthermore it happens that P_t is always symmetric and semi-positive definite.

The main advantage of solving an LQR problem with dynamic programming is that the *optimal control input* u turns out to be not an *open-loop trajectory*, but a *linear feedback policy* that depends on a time-varying matrix K_t as shown in (2.27).

Algorithm 2 shows the pseudo-code to solve the LQR problem, summarizing the procedure yet described.

Time-varying system The definition of the LQR algorithm works fine also with time-varying systems $x_{t+1} = A_t x_t + B_t u_t$; the lonely drawback in this case is that no steady-state regulation exists.

Algorithm 2 pseudo-code for solving the linear quadratic regulator problem.

```

initialize  $P_N = Q_f$ 

for  $t$  from  $N$  to  $1$  do ▷ backward pass
     $P_{t-1} = Q + A^\top P_t A - A^\top P_t B (RB^\top P_t B)^{-1} B^\top P_t A$ 
end for

for  $t$  from  $0$  to  $N - 1$  do ▷ forward pass
     $K_t = -(R + B^\top P_{t+1} B)^{-1} B^\top P_{t+1} A$ 
     $u_t^* = K_t x_t$  ▷ at runtime
end for
  
```

2.5.1 Steady-state regulation

While running the algorithm we usually observe that for N "large", the majority of the time-steps shows a constant feedback matrix K_t that converges to zero only toward the end of the OCP time horizon. This initial constant solution is called *steady-state regulator* that's associated to the solution of the *discrete-time algebraic Riccati equation* defined as

$$P_{ss} = Q + A^\top P_{ss} A - A^\top P_{ss} B (RB^\top P_{ss} B)^{-1} B^\top P_{ss} A \quad (2.29)$$

The solution of P_{ss} can be obtained by means of direct or recursive method and allows to compute the constant feedback

$$K_{ss} = -(R + B^\top P_{ss} B)^{-1} B^\top P_{ss} A$$

that can be use as long as we are not close to the end of the time horizon.

2.5.2 Inhomogeneous systems and tracking problems

The linear quadratic regulator can be solved also to the extend problem of *inhomogeneous systems*, so when at both cost and dynamic we add a fixed term, so considering the following NLP:

$$\begin{aligned}
 \min_{u, X} \quad & \sum_{t=0}^{N-1} (x_t^\top \quad u_t^\top \quad 1) \begin{bmatrix} Q_t & S_t & q_t \\ S_t^\top & R_t & s_t \\ q_t^\top & s_t^\top & 0 \end{bmatrix} \begin{pmatrix} x_t \\ u_t \\ 1 \end{pmatrix} + (x_N^\top \quad 1) \begin{bmatrix} Q_f & q_N \\ q_f^\top & 0 \end{bmatrix} \begin{pmatrix} x_N \\ 1 \end{pmatrix} \\
 \text{sub. to:} \quad & x_{i+1} = Ax_i + Bu_i + c_i \\
 & x_0 = x^{\text{init}}
 \end{aligned} \quad (2.30)$$

The solution of such regulator is achieved by applying a proper transformation that turns the problem into the yet studied formulation in (2.23). Defining the "augmented states" $\bar{x} = (x, 1)$, then the discrete-time dynamic can be rewritten as

$$\bar{x}_{t+1} = \begin{pmatrix} x_{t+1} \\ 1 \end{pmatrix} = \begin{bmatrix} A_t & c_t \\ 0 & 1 \end{bmatrix} \begin{pmatrix} x_t \\ 1 \end{pmatrix} + \begin{bmatrix} B_t \\ 0 \end{bmatrix} u_t = \bar{A}_t \bar{x}_t + \bar{B}_t u_t$$

while the cost becomes

$$\mathcal{J} = \sum_{t=0}^{N-1} (\bar{x}_t^\top \quad u_t^\top) \begin{bmatrix} \bar{Q}_t & \bar{S}_t \\ \bar{S}_t^\top & R_t \end{bmatrix} \begin{pmatrix} \bar{x}_t \\ u_t \end{pmatrix} + \bar{x}_N^\top \bar{Q}_f \bar{x}_N$$

with

$$\bar{Q}_t = \begin{bmatrix} Q_t & q_t \\ q_t^\top & 0 \end{bmatrix} \quad \bar{S}_t = \begin{bmatrix} S_t \\ s_t^\top \end{bmatrix}$$

The lonely difference with respect to the "standard" LQR case is that now in the cost there's also a mixed cost term $x^\top S u$ (previously not present).

Tracking problem Let's assume that there are two reference trajectories $x^{\text{ref}}, u^{\text{ref}}$ for both states and inputs; the inhomogeneous formulation (2.30) of the LQR problem can be also used to optimize the tracking of such references. Considering the cost

$$\mathcal{J} = \sum_{t=0}^{N-1} (x_t - x_t^{\text{ref}})^\top Q_t (x_t - x_t^{\text{ref}}) + (u_t - u_t^{\text{ref}})^\top R_t (u_t - u_t^{\text{ref}})$$

we can reduce to the inhomogeneous formulation by simply expanding the quadratic terms:

$$\begin{aligned} \mathcal{J} &= \sum_{t=0}^{N-1} x_t^\top Q_t x_t + u^\top R u - 2x_t^{\text{ref}} Q x - 2u_t^{\text{ref}} R u \\ &= \sum_{t=0}^{N-1} \begin{pmatrix} x^\top & u^\top & 1 \end{pmatrix} \begin{bmatrix} Q & 0 & Qx_t^{\text{ref}} \\ 0 & R & Ru_t^{\text{ref}} \\ x_t^{\text{ref}\top} Q & u_t^{\text{ref}\top} R & 0 \end{bmatrix} \begin{pmatrix} x \\ u \\ 1 \end{pmatrix} \end{aligned}$$

2.6 Differential dynamic programming

Differential dynamic programming (DDP) can be regarded as an extension of the linear quadratic regulator to non-linear dynamics and cost. Solutions are obtained by linearization of the problem, and for this reason convergence is not ensured, but the underlying step ideas of the algorithms are as follows:

1. linearize the cost around the current trajectory;
2. solve the LQR problem to get the variation of the controls u ;
3. perform a line-search to ensure convergence.

Denoting with U_i the set of control inputs (u_i, \dots, U_{N-1}) , we generally define the value function simply as

$$V(z, i) = \min_{U_i} \mathcal{J}_i(z, U_i)$$

and exploiting the Bellman's optimality principle it leads to the following recursive definition:

$$V(z, i) = \min_w \underbrace{\left(\ell_i(z, w) + V(f(z, w), i+1) \right)}_{Q(z, w)}$$

The dependency on w and z of Q is usually non-linear, however given a reference trajectory (\bar{x}, \bar{u}) it's possible to perform a linearization up to the second order (in order to get a quadratic cost) exploiting the Taylor series expansion:

$$Q(\bar{x} + z, \bar{u} + w) \approx Q(\bar{x}, \bar{u}) + \begin{bmatrix} Q_x^\top & Q_u^\top \end{bmatrix} \begin{pmatrix} z \\ w \end{pmatrix} + \frac{1}{2} \begin{pmatrix} z^\top & w^\top \end{pmatrix} \begin{bmatrix} Q_{xx} & Q_{xu} \\ Q_{xu}^\top & Q_{uu} \end{bmatrix} \begin{pmatrix} z \\ w \end{pmatrix} \quad (2.31)$$

where with the subscripts² we intend the differentiation operation defined as

$$Q_i = \frac{\partial Q}{\partial i} \quad Q_{ij} = \frac{\partial^2 Q}{\partial i \partial j}$$

with Q_i vector and Q_{ij} a matrix. Further denoting with V' the value function $V(\cdot, i+1)$ evaluated at the next time step, we can explicitly compute the different terms of expansion of Q in terms of the cost ℓ and dynamic f :

$$\begin{aligned} Q_x &= \ell_x + f_x^\top V'_x \\ Q_u &= \ell_u + f_u^\top V'_x \\ Q_{xx} &= \ell_{xx} + f_x^\top V'_{xx} f_x + \cancel{V'_x f_{xx}} \\ Q_{uu} &= \ell_{uu} + f_u^\top V'_{xx} f_u + \cancel{V'_x f_{uu}} \\ Q_{xu} &= \ell_{xu} + f_x^\top V'_{xx} f_u + \cancel{V'_x f_{xu}} \end{aligned} \quad (2.32)$$

In this expression some terms are neglected since the evaluation of the derivatives f_{ij} will lead to *tensors* (3D matrices) that are hard to deal with; furthermore it has been observed the contribution of such elements is in practice negligible.

The optimal cost of the linearized quadratic expression (2.31) is found as

$$w^* = \min_w Q(z, w) = -Q_{uu}^{-1} (Q_u + Q_{ux}z) = \bar{w} + Kz \quad (2.33)$$

Substituting back this optimal result in 2.31 provides us

$$\begin{aligned} V(z, i) &= \bar{Q} + Q_x^\top z + Q_u w^* + \frac{1}{2} z^\top Q_{xx} z + \frac{1}{2} w^{*\top} Q_{uu} w^* + \frac{1}{2} z^\top Q_{xu} w^* = \dots \\ &= \Delta V + V_x z + \frac{1}{2} z^\top V_{xx} z \end{aligned} \quad (2.34)$$

with

$$\begin{aligned} \Delta V &= \bar{Q} - \frac{1}{2} Q_u^\top Q_{uu}^{-1} Q_u \\ V_x &= Q_x - Q_{xu} Q_{uu}^{-1} Q_u \\ V_{xx} &= Q_{xx} - Q_{xu} Q_{uu}^{-1} Q_{ux} \\ \bar{Q} &= Q(\bar{x}, \bar{u}) + V'(\bar{x}, \bar{u}) = \bar{\ell} + \bar{V}' \end{aligned} \quad (2.35)$$

Regularization Looking at the terms in (2.35) it can be easily observe that their computation relies on the inversion of the matrix Q_{uu} that, unluckily, can't be guaranteed a-priori. In general to avoid the possibility of having singular matrices that needs to be inverted, we can perform a *regularization* by adding a rescaled identity matrix:

$$\bar{Q}_{uu} = Q_{uu} + \mu I \quad (2.36)$$

with μ "small". This do not provide any mathematical guarantee that \bar{Q}_{uu} is actually invertible, but numerically it becomes more likely. Intuitively one can regard such regularization as adding a new term in the linearized cost proportional to the input itself, i.e.

$$\bar{\ell}(\bar{x} + z, \bar{u} + w) = \ell(\bar{x} + z, \bar{u} + w) + \frac{1}{2} \mu \|w\|^2$$

²notation that will be also used later while defining the value function V_i , V_{ij} and the costs ℓ_{ij} .

Based on the regularization, terms in (2.35) are reconsidered as

$$\begin{aligned}\Delta V &= \bar{Q} + Q_u^\top \bar{w} + \frac{1}{2} \bar{w}^\top Q_{uu} \bar{w} \\ V_x &= Q_x^\top + Q_u^\top K + \bar{w}^\top Q_{uu} K + \bar{w}^\top Q_{xu}^\top \\ V_{xx} &= Q_{xx} + K^\top Q_{uu} K + Q_{xu} K + K^\top Q_{xu}^\top\end{aligned}\tag{2.37}$$

with

$$\bar{w} = -\bar{Q}_{uu}^{-1} Q_u \quad K = -\bar{Q}_{uu}^{-1} Q_{ux}\tag{2.38}$$

that are based on the optimal control solution (2.33) exploiting the regularized matrix.

2.6.1 Iterative LQR

One way to solve differential dynamic programming is by means of the *iterative linear quadratic regulator (iLQR)* algorithm, summarized as follows:

1. given an initial set of control inputs U and the dynamic f of the system, we simulate the sequence of states X ;
2. we perform a *backward pass*: exploiting the definition of the value function, we start from the end and compute the updated values of the inputs and the corresponding gains;
3. we perform a *forward pass* with a line-search to find the optimal update step.

Steps 2 and 3 are iterated until convergence. Algorithm 3 presents the whole procedure of the iLQR.

Line-search Once the backward pass is done we have a set of variational inputs \bar{w} that can be applied to the yet applied inputs \bar{U} ; the goal now is to estimate "how much times" (α) of such quantity we should add to obtain the best cost reduction.

To do so we firstly need to compute the cost given by the value function at time 0 considering a state variation $x - \bar{x} = z = 0$ that evaluates to

$$V_0(0) = \Delta V_0 =$$

TODO: finire di spiegare l'algoritmo e aggiornare lo pseudo-codice per il costo e regolarizzazione dinamica

2.7 Model predictive control

In *model predictive control (MPC)* we use the resulting optimal control for controlling a system/robot inside a control loop. The main idea is to use the current state as initial condition for solving a finite horizon OCP and apply just the first torque computed (not the whole trajectory); at the next time-step the controller will solve the same problem but with a new initial condition.

Algorithm 3 pseudo-code for solving the differential dynamic programming problem using the iterative LQR algorithm.

Require: U initial set of inputs, f dynamic of the system, x_0 initial state

$X \leftarrow f(U, x_0)$

$\mu \leftarrow \mu_0$

while is not converged **do**

$\bar{U} \leftarrow U$

▷ store reference input trajectory

$\bar{X} \leftarrow X$

▷ store reference state trajectory

$c_0 \leftarrow \ell(\bar{X}, \bar{U})$

▷ cost with current trajectories

$V_x(N) \leftarrow \nabla_x \ell_N$

▷ backward pass steps

$V_{xx}(N) \leftarrow \nabla_{xx} \ell_N$

for i from $N - 1$ to 0 **do**

compute $Q_x, Q_u, Q_{xx}, Q_{uu}, Q_{xu}$ as in (2.32)

$\bar{Q}_{uu} \leftarrow Q_{uu} + \mu I$

▷ regularization

$\bar{w} \leftarrow -\bar{Q}_{uu}^{-1} Q_u$

▷ eq. (2.38)

$K \leftarrow -\bar{Q}_{uu}^{-1} Q_{ux}$

end for

$\alpha \leftarrow 1$

▷ forward pass steps

for k from 1 to k_{max} **do**

▷ perform a line-search

$X, U \leftarrow \text{simulate system with input } \bar{U} + \alpha \bar{w} + K(x - \bar{x})$

$c \leftarrow \ell(X, U)$

$\Delta \ell(\alpha) \leftarrow \Delta V_0 - \ell(\bar{U})$

if $c < \gamma c_0$ **then**

exit the *for* loop

else

$\alpha \leftarrow k_\alpha \alpha$

▷ $k_\alpha \in (0, 1)$

end if

end for

if $k = k_{max}$ **then**

algorithm has converged

return U, K

▷ open loop control and feedback gains

end if

end while

Considering an already-discretized OCP problem, MPC deals with the solution of the following problem at each time-step:

$$\begin{aligned} X^*, U^* = \arg \min_{X, U} \quad & \sum_{k=0}^{N-1} \ell(x_k, u_k) \\ \text{sub. to:} \quad & x_{k+1} = f(x_k, u_k, k) \quad k = 0, \dots, N-1 \\ & x_{k+1} \in \mathcal{X}, u_k \in \mathcal{U} \quad k = 0, \dots, N-1 \\ & x_0 = x^{\text{meas}} \end{aligned} \quad (2.39)$$

At each time-step the applied control is simply $\tau = u_0^*$, while the rest of both U^*, X^* is unused (numerically to improve convergence speed, we might use the yest computed U^* as initial guess for the next OCP).

Challenges in model predictive control This control technique seems promising as it solves at each iteration the optimal control possible, however we have to tackle few challenges that are not that irrelevant, in particular

- i) *feasibility* deals with the fact that OCPs might be not always feasible, so that we might reach states where not all constraints can be satisfied. To avoid this issue we have to ensure *recursive feasibility*;
- ii) *stability* deals with the stabilization of the system itself. In fact at each iteration the OCP looks on a different and "increased" time-horizon after which the problem is solved. It might happen in fact that in order to "stabilize in the future", the system tends to initially diverge from the desired configuration: if this process is iterated we might go toward a system instability (and not stability);
- iii) *computation time*: solving OCP problem requires time, and if we think to re-compute everything at each iteration we need to do it fast.

2.7.1 Feasibility

Infinite horizon MPC Let's consider a case where at each iteration the time horizon is set to infinite ($N = \infty$): this means that each time the OCP sees always the "same amount of time". Assuming that there are no disturbances, predicted and actual trajectories computed each time are the same, since they follow from the Bellman's optimality principle.

If we now consider a cost $\ell(x, u) \geq \alpha \|x\|$ for any combination x, u with $\alpha > 0$, then having a finite cost implies that the system is stable (since eventually x will set to zero, thus no cost is added).

Moreover, since predicted and actual trajectories are equal, it turns out that we have *recursive feasibility* along a closed-loop trajectory.

The main idea now is to add a terminal cost and a constraint to a finite horizon OCP in order to "mimic" an infinite horizon problem, reminiscing the concept of the value function.

Terminal cost and constraint With the idea just said, the problem that we want to solve in order to improve feasibility becomes

$$\begin{aligned}
 V_N^*(x_0) = \arg \min_{X, U} \quad & \sum_{k=0}^{N-1} \ell(x_k, u_k) + \ell_f(x_N) \\
 \text{sub. to:} \quad & x_{k+1} = f(x_k, u_k, k) \quad k = 0, \dots, N-1 \\
 & x_{k+1} \in \mathcal{X}, u_k \in \mathcal{U} \quad k = 0, \dots, N-1 \\
 & x_0 = x^{\text{meas}} \\
 & x_N \in \mathcal{X}_f \quad \text{terminal constraint}
 \end{aligned} \tag{2.40}$$

The main objective now is concerning the generation of both $\ell_f(\cdot)$ and \mathcal{X}_f .

Feasibility If an OCP is subjected to only input constraints, then it happens that it's always feasible; on the other hand, if the problem is subjected to hard state constraint, if $N < \infty$ there's no guarantee that the OCP remains feasible, even while staying in the nominal case.

The *maximum output admissible theory* tells us that in general a limited amount of horizon time-steps $N < \infty$ is enough to guarantee the problem feasibility; in particular it's necessary to use a *maximal control invariant set* as a terminal constraint to ensure that the closed-loop convergence. Such theory is based on the following theorem:

« If \mathcal{X}_f is a control invariant set, then the model predictive control problem is recursively feasible. »

In particular

« a set \mathcal{S} is control invariant if $\forall x \in \mathcal{S}$ there exists an input $u \in \mathcal{U}$ such that $f(x, u) \in \mathcal{S}$. »

i.e. when we start in the set \mathcal{S} it's always possible to stay there.

In practice computing control invariant sets is generally hard, specially for non-linear systems. For linear system with dynamic $x^+ = Ax$ and output $y = Cx$ with $y \in \mathcal{Y} = \{y \in \mathbb{R}^p : h_i(y) \leq 0, i = 1, \dots, s\}$ then the maximal output admissible set is found as

$$\mathcal{O}_\infty = \{x \in \mathbb{R}^n : h_i(CA^+x) \leq 0, i = 1, \dots, s, t = 0, \dots, \infty\}$$

TODO: finire questa parte

2.7.2 Stability

In order to address the stability issue of MPC, we have to firstly define what's a positive invariant set and an exponential Lyapunov function.

« A set \mathcal{S} is positive invariant set if $\forall x \in \mathcal{S}$ it holds that $f(x) \in \mathcal{S}$ »

« Suppose that \mathcal{X} is a positive invariant set, a function $V : \mathbb{R}^n \rightarrow \mathbb{R}^+$ is an exponential Lyapunov function if $\exists \alpha_1, \alpha_2, \alpha_3 > 0$ such that

$$\begin{aligned}
 V(x) &\geq \alpha_1 \|x\| \\
 V(x) &\leq \alpha_2 \|x\| \\
 V(f(x)) - V(x) &\leq -\alpha_3 \|x\| \\
 &\text{for all } x \in \mathcal{X}. \text{ »}
 \end{aligned}$$

With this premise

« If there exists a Lyapunov function in the set \mathcal{X} , then the origin is exponentially stable in \mathcal{X} ; furthermore, if $\mathcal{X} = \mathbb{R}^n$ the origin is globally exponentially stable. »

Since the value function $V(\cdot)$ represents the cost-to-go that we expect to decrease over time, then we can regard it as a Lyapunov function. Given the optimal value function $V_0^*(x_0)$ and the value function(not necessarily optimal) at the sequent time-step $V_1(x_1)$ as

$$V_0^*(x_0) = \sum_{i=0}^{N-1} \ell(x_i, u_i) + \ell_f(x_N) \quad V_1(x_1) = \sum_{i=1}^N \ell(x_i, u_i) + \ell_f(x_{N+1})$$

then their difference evaluates to

$$V_1(x_1) - V_0^*(x_0) = \ell(x_N, u_N) + \ell_f(x_{N+1}) - \ell_f(x_N) - \ell(x_0, u_0) \leq -\alpha_3 \|x_0\| \quad \forall x_N \in \mathcal{X}_f$$

To have $V(\cdot)$ as a Lyapunov function we need to assume that

- $f(0,0) = 0, \ell(0,0) = 0, \ell_f(0) = 0$;
- \mathcal{Z} is closed and \mathcal{X}_f is compact (i.e. closed and bounded);
- for any state $x \in \mathcal{X}_f$ exists a control $u \in \mathcal{U}$ such that

$$f(x, u) \in \mathcal{X}_f$$

so \mathcal{X}_f is a control invariant set, and that the terminal cost decreases, i.e.

$$\ell_f(f(x, u)) - \ell_f(x) \leq -\ell(x, u)$$

Furthermore must exists two constants $\alpha_1, \alpha_f > 0$ that allows to lower bound the running cost and upper bound the terminal cost, i.e.

$$\begin{aligned} \ell(x, u) &\geq \alpha_1 \|x\| & \forall x \in \mathcal{X}_N, \forall u \in \mathcal{U} \\ \ell_f(x) &\leq \alpha_f \|x\| & \forall x \in \mathcal{X}_f \end{aligned}$$

where \mathcal{X}_N is the set of states which the OCP has a solution.

If all this axioms are satisfied, then $V_N^*(\cdot)$ is a Lyapunov function.

TODO: aggiungere considerazioni finali con sistema lineare

2.7.3 Computation time

To improve the convergence time of the numerical algorithms to solve OCPs, we can use a *warm start* that exploits the previous solution of the same problem as initial guess for the current OCP. This is done by shifting back by 1 time-step the optimal trajectory, i.e. $u_k^{\text{guess}} = u_{k+1}^*$; the last element is instead initialized to zero ($u_{N-1}^{\text{guess}} = 0$). To bound also the computation time, we don't iterate until a full convergence, but we just do a fixed amount of Newton iteration (1 is usually fine).

Chapter 3

Reinforcement Learning

Reinforcement learning (RL) is a method originally developed within the data-science community that, in the last years, has also been extended to robotic and control applications. It has a lot of similarities with optimal control, however key differences are that:

- reinforcement learning tries to find the *global optimal policy* for arbitrarily complex problems;
- assumes that the system's *dynamic* is *unknown* and assumed to be *stochastic*¹;
- initially it was developed to solve problem with discrete sets for both states and control, however in the recent days it has also been extended to continuous-time systems²;
- typically solves an infinite horizon problem.

Terminology Even if concepts are very similar, reinforcement learning uses slightly different terminology and notation, as can be seen in table 3.1. The main difference is in the attitude toward the problem: in optimal control we want to minimize a cost, while in reinforcement learning we want to maximize the reward.

For sake of simplicity, in this chapter we will mainly stick with the optimal control notation and consider the problem as a minimization (not a maximization of the reward); only exception is the subtle difference between *value function* and *optimal value function* that will be taken from the reinforcement learning vocabulary.

¹in this chapter we will mainly consider deterministic dynamic as subset of the more general formulation.

²in general the discretization of the continuous-time dynamics scales badly with the growth of the problem size

Table 3.1: terminology and notation comparison between optimal control and reinforcement learning.

optimal control	reinforcement learning
state x	state s
control u	action a
plant/dynamic	environment
cost	reward
cost-to-go	return
cost-to-go of a policy	value function
value function	optimal value function

3.1 Markov decision process

A *Markov decision process* (MDP) is used to describe the environment of a reinforcement learning problem; a MDP is characterized by being fully observable and satisfying the *Markov property* for which "the future is independent from the past, given the present", mathematically speaking

$$\Pr \{x_{t+1}|x_t\} = \Pr \{x_{t+1}|x_t, x_{t-1}, \dots, x_1\} \quad (3.1)$$

where $\Pr \{X|Y\}$ is the *conditional probability* of X happening given that Y has already happened. The probability of switching from one state to another is described by the *state transition matrix* $\mathcal{P}_{xx'} = \Pr \{x_{t+1} = x'|x_t\}$ defined as

$$\mathcal{P}_{xx'} = \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \\ P_{n1} & & P_{nn} \end{bmatrix} \quad P_{ij} \in [0, 1] \quad (3.2)$$

Each state transition matrix, to be valid, must have that the sum of all the elements of each row adds up to 1 (since each state must evolve surely at each time-step). Furthermore in case of deterministic dynamic each element is either 0 (transition not possible) or 1 (only transition admissible).

Markov process

A *Markov process* (MP) is described by the tuple

$$\langle \mathcal{X}, \mathcal{P} \rangle \quad (3.3)$$

where \mathcal{X} is the (finite) set of the allowed states and \mathcal{P} is the state transition matrix ruling the probability of going from one state to the other; usually a Markov process is also referred as a *Markov chain*.

Properties of the state transition matrix Since \mathcal{P} has all non-negative entries and is irreducible (because it's associated to a fully connected graph), by the *Person-Frobenius theorem* it holds that

- the largest (in norm) eigenvalue r of \mathcal{P} must satisfy

$$\min_i \sum_j P_{ij} \leq r \leq \max_i \sum_j P_{ij}$$

Since each row sums up exactly to 1, then it follows that the maximum eigenvalue of \mathcal{P} is $r = 1$;

- all other eigenvalue λ of \mathcal{P} are smaller then r , i.e.

$$\lambda_i\{\mathcal{P}\} < 1 \quad i = 2, 3, \dots$$

A Markov chains sets up a probabilistic framework to describe a system's dynamic, considering $x^+ = \Pr \{x^+|x\}$; optimal control instead relies on the deterministic approach $x^+ = f(x, u)$ that we will try to stick with.

Markov reward process

A *Markov reward process (MRP)* is described by a tuple

$$\langle \mathcal{X}, \mathcal{P}, C, \gamma \rangle \quad (3.4)$$

where C is the *cost function* and $\gamma \in (0, 1)$ is the *discount factor*. The first estimates the expected cost at the next time-step given the current state, i.e.

$$c_x = \mathbb{E} \{ \ell_{t+1} | x = x_t \}$$

while the discount factor, as name suggests, is a measure of "how relevant a future cost is with respect to the present". The cost-to-go (*return* in RL jargon) is in fact the *total discounted cost* from t to infinity, i.e.

$$\mathcal{J}_t = \ell_t + \gamma \ell_{t+1} + \gamma^2 \ell_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k \ell_{t+k} \quad (3.5)$$

Here it's clear that the discount factor can be chosen to favor a reduction of the cost "in short time" ($\gamma \rightarrow 0$: myopic evaluation) or in the "long run" ($\gamma \rightarrow 1$: far-sighted evaluation).

Value function In reinforcement learning the value function $V(\cdot)$ represents the cost that we pay starting from a given state x , i.e.

$$V(x) = \mathcal{J}_t(x = x_t)$$

In light of the Bellman's optimality principle and (3.5) it turns out that

$$V(x) = \ell(x) + \gamma V(f(x)) \quad (3.6)$$

Considering in general the elements (3.4) of the Markovian reward process, then it directly follows

$$V = C + \gamma \mathcal{P}V$$

$$\begin{pmatrix} V(1) \\ \vdots \\ V(n) \end{pmatrix} = \begin{pmatrix} \ell(1) \\ \vdots \\ \ell(n) \end{pmatrix} + \gamma \begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & & P_{nn} \end{bmatrix} \begin{pmatrix} V(1) \\ \vdots \\ V(n) \end{pmatrix}$$

It turns out that this is a linear equation in the unknown V , so by simply reversing the terms

$$V = (I - \gamma \mathcal{P})^{-1} C \quad (3.7)$$

This formulation is simple, but is applicable only to Markov chains of small size; when the number of states become larger, the numerical inversion of $I - \gamma \mathcal{P}$ becomes numerically ill conditioned, thus iterative methods to evaluate V are preferred.

Markov decision process

A *Markov decision process (MDP)* is described by the tuple

$$\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, C, \gamma \rangle \quad (3.8)$$

where \mathcal{U} is the (finite) set of control inputs. Denoting

$$P_{xx'}^u = \Pr \{ x' = x_{t+1} | x = x_t, u = u_t \} \quad C_x^u = \ell_t(x = x_t, u = u_t)$$

we define the *policy* π the probability distribution of the control inputs u given the states x , i.e.

$$\pi\{u|x\} = \Pr \{ u = u_t | x = x_t \} \quad (3.9)$$

In case of deterministic policies the policy reduces simply to a function $u = \pi\{x\}$.

Action-value function We denote with $Q^\pi(x, u)$ the *action-value function* that's representing the cost that we pay starting from a state x and applying inputs u based on the control policy π :

$$Q^\pi(x, u) = \mathcal{J}_t(x = x_t, u = u_t, u_{k>t} \simeq \pi) \quad (3.10)$$

and considering the Bellman's theorem

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

where to simplify the notation we write $Q^\pi(x, \pi\{x\})$ simply as $V^\pi(x)$

Optimal value function and policy The *optimal value function* is the one associated with a policy the minimize it's value, i.e.

$$V^*(x) = \min_{\pi} V^\pi(x) \quad (3.11)$$

thus the *optimal action-value function* is

$$Q^*(x, u) = \min_{\pi} Q^\pi(x, u) \quad (3.12)$$

The optimal policy π^* is the one that satisfies $\pi^* \leq \pi$ for any other policy π , where $\pi \leq \pi'$ requires that $V^\pi(x) \leq V^{\pi'}(x)$ for any x . This is achieved by minimizing the optimal action-value function among all possible inputs, i.e.

$$\pi^*(x) = \arg \min_{u \in \mathcal{U}} Q^*(x, u) \quad (3.13)$$

The main idea is that if we have access to the optimal action-value function, then it's possible to generate the optimal control policy.

Bellman optimality equality Given that the optimal value function is

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

and recalling (\dagger) allows us to state the following condition that the optimal solution must obey to:

$$\begin{cases} V^*(x) &= \min_u \ell(x, u) + \gamma V^*(f(x, u)) \\ Q^*(x, u) &= \ell(x, u) + \gamma \min_{u'} Q^*(f(x, u), u') \end{cases} \quad (3.14)$$

Even in this case of deterministic system, there's no close-form solution for the non-linear minimization so iterative algorithms are necessary.

3.2 Dynamic programming

Dynamic programming (DP) is a set of methods that are relying on the full knowledge of the Markovian decision process in order to solve two kinds of problems:

- *prediction*, so given $\langle \mathcal{X}, \mathcal{U}, \mathcal{P}, C, \gamma \rangle$ and a policy π , it computes the value function V^π associated to such policy;
- *control*, so given the MDP it outputs both the optimal policy π^* as well as the related optimal value function V^* .

The latter is of course more interesting in control application as it provides us the best way to act on a system given the cost it's subjected to.

Finite and infinite horizon Dynamic programming was already introduced in sec. 2.1, page 15, in order to find the global optimum of a discretized optimal control problem; in that context we dealt with finite horizon problems for which theory is simpler and algorithms can compute solutions in a finite time.

For what concerns reinforcement learning we usually deal with *infinite horizon* problem for which theory is more complex (but elegant), assuming that both policy and value functions are time independent. Such theory is also a good approximation of problems with "long" (but finite) time horizons.

Bellman's operators Before continuing with dynamic programming, let's first define the *Bellman (expectation backup) operator*

$$T^\pi(V) = C^\pi + \gamma \mathcal{P}^\pi V \quad (3.15)$$

and the *Bellman optimality backup operator*

$$T(V) = \min_{u \in \mathcal{U}} C^\pi + \gamma \mathcal{P}^\pi V \Leftrightarrow (TV)(x) = \min_{u \in \mathcal{U}} \ell(x, u) + \gamma V(f(x, u)) \quad (3.16)$$

3.2.1 Prediction: iterative policy evaluation

Let's consider a prediction problem where it's requested to evaluate the value function V given the policy π . The solution to such problem is performed by iteratively applying the Bellman expectation operator:

$$V_{k+1} = T^\pi(V_k) \quad (3.17)$$

where k is the iteration index; expanded the value function is iteratively updated as

$$V_{k+1}(x) = \ell(x, \pi\{x\}) + \gamma V_k(f(x, \pi\{x\})) \quad \forall x \in \mathcal{X}$$

So in practice we can chose an arbitrary initial value function V_0 since the *iterative policy evaluation* algorithm is guarantee to asymptotically converge, i.e.

$$\lim_{k \rightarrow \infty} V_k = V^\pi$$

Proof of convergence In order to prove the convergence of the iterative policy evaluation we firstly need to show that T^π is contracting. Given two value function V, Z we see that

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

Due to the Person-Frobenius theorem it follows that

$$\gamma \|\mathcal{P}^\pi(V - Z)\|_\infty \leq \gamma \|V - Z\|_\infty \Rightarrow \|T^\pi(V) - T^\pi(Z)\|_\infty \leq \gamma \|V - Z\|_\infty$$

so given any two value function, the relative step obtained by applying the Bellman's expectation operator is upper bounded by the norm difference $V - Z$.

We can now show that given V_0 and V_π , then $V_k \rightarrow V_\pi$ for $k \rightarrow \infty$. Considering

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

so generalizing

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

Since $\gamma \in (0, 1)$, then V_k converges to V_π at a geometric rate.

Lastly we show that V_π is a unique fixed point for the operator T^π ; assuming that V and Z are both fixed point for T^π , i.e. $T^\pi(V) = V$ and $T^\pi(Z) = Z$, then we can show that V must equate Z . Because T^π is contracting, then

$$\|T^\pi(V) - T^\pi(Z)\|_\infty \leq \gamma \|V - Z\|_\infty$$

but since V, Z are fixed points

$$\|T^\pi(V) - T^\pi(Z)\|_\infty = \|V - Z\|_\infty$$

This condition must hold at the same time, so

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

but since $\gamma < 1$, then it must have $\|V - Z\|_\infty = 0$, implying that $V = Z$.

3.2.2 Control: policy and value iteration

Policy iteration In this case we deal with a more complex problem where, given a Markovian decision process, we want to find the optimal policy π^* . The idea is that we start from an arbitrary policy π_0 and we perform until convergence the following two steps:

1. we evaluate the policy π_k , so we compute the value function V^{π_k} ;
2. we improve the policy by acting greedily with respect to the value function V^{π_k} , so

$$\pi_{k+1}(x) = \arg \min_{u \in \mathcal{U}} \ell(x, u) + \gamma V^{\pi_k}(f(x, u)) \quad (3.18)$$

Such algorithm always converges to the optimum policy:

$$\lim_{k \rightarrow \infty} \pi_k = \pi^*$$

Proof of convergence Let's consider the simplified mathematical notation $\pi' \leq \pi$ to represent $V^{\pi'} \leq V^\pi$, we want to show now that $\pi_{k+1} \leq \pi_k$ for any iteration k , implying that we are always going closer to the optimal solution. Calling $\pi = \pi_k$ and $\pi' = \pi_{k+1}$ we see that

$$\begin{aligned} \pi'(x) &= \arg \min_u \left(\ell(x, u) + \gamma V^\pi(f(x, u)) \right) \\ &= \underbrace{\min_u \left(\ell(x, u) + \gamma V^\pi(f(x, u)) \right)}_{i)} \leq \underbrace{\ell(x, \pi(x)) + \gamma V^\pi(f(x, \pi(x)))}_{ii)} \end{aligned}$$

where $i)$ is the cost that we get by following π' for the first step and then following π , while $ii)$ is the one obtained by simply following π . Considering now that

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

since by acting greedily always leads to a cost improvement (or at worst the cost remains equal), then we can see that

$$\begin{aligned} V^\pi(x) &\geq \min_u \ell(x, u) + \gamma V^\pi(f(x, u)) = \ell(x, \pi') + \gamma V^\pi(f(x, \pi')) \\ &\geq \ell(x, \pi') + \gamma Q^\pi(x', \pi') = \ell(x, \pi') + \gamma \ell(x', \pi') + \gamma^2 V^\pi(x'') \\ &\geq \ell(x, \pi') + \gamma \ell(x', \pi') + \gamma^2 Q^\pi(x'', \pi') \\ &\geq \sum_{i=0}^{\infty} \gamma^i \ell(x^{(i)}, \pi') = V^{\pi'}(x) \\ &\geq V^{\pi'}(x) \end{aligned}$$

thus showing that $\pi' \leq \pi$, so at each iteration we can't have a policy worsening.

Modified policy iteration The main issue in policy iteration is that evaluating exactly the policy can take many iterations, so the idea is just to compute an approximation of V^{π_k} using just m_k policy evaluation iterations. Under mild assumptions this eventually converges to V^{π_k} .

With this idea taking $m_k = \infty$ leads to the yet described policy iteration method, while with $m_k = 1$ we get *value iteration*; in the latter case every-time we update V we use a policy that's greedy with respect to the value function itself.

Value iteration The algorithm of *value iteration* simply starts with an arbitrary guess V_0 for the value function and at each iteration

$$V_{k+1}(x) = \min_{u \in \mathcal{U}} \ell(x, u) + \gamma V_k(f(x, u)) \quad \forall x \in \mathcal{X} \quad \Rightarrow V_{k+1} = T(V_k)$$

The value function eventually converges to the optimal value

$$\lim_{k \rightarrow \infty} V_k = V^*$$

and the optimal policy π^* is computed at the end from V^* as

$$\pi^* = \arg \min_{u \in \mathcal{U}} \ell + \gamma V^*$$

TODO: proof

3.3 Model-free prediction

Model-free prediction deals with the problem of finding the value function $V(\cdot)$ of an unknown Markovian decision process by simply acting on the environment and observing the response.

3.3.1 Monte Carlo policy evaluation

Monte Carlo (MC) methods can be used to solve model-free prediction problem and it's based on the simple idea of having the value function that's the mean return. This method comes with the caveat that can be applied only to *episodic* Markov decision processes, so each episode must have an end.

The goal of this method is to learn V^π from episodes of experience under a policy π , so once it has been collected

$$x_1, u_1, \ell_1, x_2, u_2, \ell_2, \dots, x_k \sim \pi$$

Once experience has been collected we can use the equation of the total discounted cost to compute the overall cost-to-go as

$$\mathcal{J}_t = \ell_t + \gamma \ell_{t+1} + \dots + \gamma^{T-1} \ell_{T-1}$$

The value function then is the expected cost-to-go under a policy π ; this can be modeled including also the stochastic part (due to the policy, the MDP and the cost) considering

$$V^\pi(x) = \mathbb{E}_\pi \{ \mathcal{J}_t | x_t = x \} \quad (3.19)$$

Monte Carlo policy evaluation estimates V^π as the average cost-to-go; for what concerns the algorithm the first time a state is visited in an episode we should:

- increment a counter $N(x) = N(x) + 1$;
- increment the total cost of such state $C(x) = C(x) + \mathcal{J}(x)$;
- estimate the value function as $V(x) = V(x)/N(x)$.

By law of large number the estimate $V(x)$ tends to $V^\pi(x)$ as $N(x) \rightarrow \infty$. The idea is that by getting experience from more and more episodes, the stochastic contribution on the value function estimation gets negligible.

It also exists an *every-visit Monte Carlo* variant for which counter and costs are incremented every-time a state is visited (even inside the same episode).

Incremental Monte Carlo update To avoid storing the total cost $C(x)$ as well as the value function estimate $V(x)$ that are highly correlated, we can compute the mean incrementally, in fact

$$\begin{aligned} V_N &= \frac{1}{N} \sum_{i=0}^N \mathcal{J}_i = \frac{\mathcal{J}_N + \sum_{i=0}^{N-1} \mathcal{J}_i}{N} = \frac{\mathcal{J}_N + (N-1)V_{N-1}}{N} \\ &= V_{N-1} + \frac{\mathcal{J}_N - V_{N-1}}{N} \end{aligned}$$

This definition can be easily extended for *non-stationary* problems in order to consider a running mean that allow to embed a *forget factor* α to discount information from older episodes:

$$V(x) = V(x) + \alpha(\mathcal{J} - V(x)) \quad (3.20)$$

3.3.2 Temporal difference learning

Temporal difference learning (TDo) is an alternative to Monte Carlo evaluation for model-free prediction and is based on the cost-to-go estimation based on 1 step look-ahead to update the value function:

$$V(x_t) = V(x_t) + \alpha_t \left(\underbrace{\ell_t + \gamma V(x_{t+1})}_{\text{TD target}} - V(x_t) \right) \quad (3.21)$$

TD error

where the TD target is the estimated return; here we can find an high similarity with the Monte Carlo incremental update (3.20). Moreover choosing $\alpha_t = 1$ and sampling all state uniformly, then TDo is exactly the policy evaluation algorithm presented in sec. 3.2.1, page 39.

Properties The TDo algorithm is a stochastic approximation algorithm; since the Bellman operator 3.15 has a unique fixed point, then if TDo converges and all states are sampled infinitely many times than the estimate must converge to V^π .

Converges is guaranteed is the step-size sequence α_t satisfies the *Robinson-Monroe condition* for which

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

A plausible sequence that satisfies this condition is $\alpha_t = \bar{\alpha}/t$, however in practice people use constant step size since the algorithm works anyway.

Table 3.2: comparison between Monte Carlo and temporal difference methods for solving a model-free prediction problem.

Monte Carlo	temporal difference
must wait until the end of an episode to know the cost	can learn after every step
only works with episodic (terminating) environments	works also in non-terminating environment
the cost-to-go \mathcal{J}_t is an unbiased estimate of $V^\pi(x_t)$	target is a biased estimate of $V^\pi(x_t)$

Comparison with Monte Carlo Table 3.2 mainly compares Monte Carlo and temporal difference methods, showing the pros and the cons of each one.

In general Monte Carlo has an high variance but a zero bias; it's characterized by good convergence properties (even with function approximators) and is not very sensitive to the initialization of the problem. It's also simple and intuitive to understand.

Temporal difference instead has a lower variance but embeds some bias; it's usually more efficient then Monte Carlo and TDo converges to $V^\pi(x)$ (but not always with function approximators); furthermore is more sensitive to the initialization of the estimate.

n-step return The concept of temporal difference can also be extended to a n -step return; TDo in fact can be regarded as a 1-step return, but we can have a general formulation

$$\begin{aligned} n=2 \quad \mathcal{J}_t^{(2)} &= \ell_t + \gamma \ell_{t+1} + \gamma^2 V(x_{t+2}) \\ n \quad \mathcal{J}_t^{(n)} &= \ell_t + \gamma \ell_{t+1} + \dots + \gamma^{n-1} \ell_{t+n-1} + \gamma^n V(x_{t+n}) \end{aligned}$$

Note that by choosing $n = \infty$, the evaluation describes the Monte Carlo method. In general a TD($n-1$) method is updated with the following function:

$$V(x_t) = V(x_t) + \alpha_t (\mathcal{J}_t^{(n)} - V(x_t)) \quad (\circ)$$

Forward view TD λ The λ -cost-to-go \mathcal{J}_t^λ combines all n -step cost-to-go $\mathcal{J}_t^{(n)}$ using as weights the sequence $(1-\lambda)\lambda^{n-1}$, i.e.

$$\mathcal{J}_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \mathcal{J}_t^{(n)} \quad (3.22)$$

As for Monte-Carlo, such temporal difference can be computed only from complete episodes (since the time horizon n theoretically should go to infinity), for this reason the method is referred as *forward view TD λ* ; once \mathcal{J}_t^λ is computed, we can use (\circ) to update the value function.

Backward view TD λ The *backward view TD λ* allows an online update of the value function at every step starting from incomplete sequences. Defining the *eligibility traces* $e_t(\cdot)$ as

$$\begin{aligned} e_0(x) &= 0 \\ e_t(x) &= \gamma \lambda e_{t-1}(x) + \mathbb{1}(x_t = x) \quad \gamma < 0 \end{aligned} \quad (3.23)$$

Table 3.3: classification of control learning methods based on their interaction and optimization performance.

optimization performance	environment interaction	
	yes	no
offline	active learning	non-interactive learning
online	online learning	

Table 3.4: control learning categories based on the quantities that should be learned and the knowledge on the respective Markov decision process.

Markov decision process	learned quantities	
	$V(x)/Q(x,u)$	$V(x)/Q(x,u)$ and $\pi(x)$
known	value iteration	policy iteration
unknown	direct methods	actor-critic methods

then $e(x)$ increase every-time I visit x ; otherwise it decays exponentially. Such function $e(\cdot)$ is a measure of how often and how recently I've visited x ; recalling the TD error $\delta_t = \ell_{t+1} + \gamma V(x_{t+1}) - V(x_t)$, then the update of the value function is given by

$$V(x) = V(x) + \alpha \delta_t e_t(x) \quad (3.24)$$

3.4 Control learning methods

While dealing with *control learning* problems, we have different methods based on the interactivity (if the learned can interact actively with the environment, influencing the future) and the on/off-line performance (if the cost-to-go is compute while learning or after learning), as it can be seen in table 3.3.

As can be seen in table 3.4, control learning problems are further categorized based on the quantities that need to be learned, as well as on the knowledge of the underlying Markov decision process of the environment.

It has to be noted that when the MDP is known, that it's required to learn the $V(x)$, while if the MDP is unknown we usually estimate the action-value function $Q(x,u)$ (instead of V) since it allows to compute the policy

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

As last classification, we have *on-policy learning* when we learn about a given policy π using experience sampled from π itself, while in *off-policy learning* we learn π using experience sampled using another policy $\tilde{\pi}$.

3.4.1 Q learning

Q learning is a direct method³ that iteratively updates an estimate $Q_k(x,u)$ (thus the name of the method) of the optimal action-value function $Q^*(x,u)$. After observing a transition (x,u,ℓ,x') , the algorithm update considering a TD error δ as follows:

$$\begin{aligned} \delta(Q_k) &= \ell + \gamma \min_{u' \in \mathcal{U}} Q(x',u') - Q(x,u) \\ Q_{k+1}(x,u) &= Q_k(x,u) + \alpha_k \delta(Q_k) \end{aligned} \quad (3.25)$$

³so it do not require the a-priori knowledge on the Markov decision process

At stochastic equilibrium, for each pair (x, u) visited infinitely often it must hold

$$\mathbb{E} \{ \delta(Q) \} = 0 \quad \Rightarrow \quad TQ - Q = 0 \quad \Rightarrow \quad Q = Q^*$$

The requirement that (x, u) must be visited infinitely often set a need for exploration, i.e. the learner must be able to learn from any possible condition. Q_k then converges to Q^* when appropriate learning rates α_k are used.

Control strategy While solving an online learning problem, key is choosing the input u that needs to be applied to the system. A strategy that we might use is to be greedy with respect to the action-value function, i.e.

$$u(x) = \arg \min_w Q_k(x, w) \quad (3.26)$$

This choice allows us to always chose the control that leads to an estimated lower cost in the future, however this might not be optimal since exploration is not ensure and could lead to a large loss over time. A good learner in fact must choose also suboptimal controls in order to explore.

For this reason we might use the so called *ε -greedy strategy* for which

- we choose a random control with probability ε ;
- we choose a greedy control (3.26) with probability $1 - \varepsilon$.

This simple technique ensures exploration, and in practice we can have a value of ε that changes over time (typically decreasing over time toward zero).

Another technique that we might use is the *Boltzmann exploration* for which the control is chosen with a probability Π proportional to it's mean, considering so

$$\Pi(u) = \frac{\exp(-\beta Q_t(u))}{\sum_{u' \in \mathcal{U}} \exp(-\beta Q_t(u'))} \quad (3.27)$$

where $Q_t(u)$ is the mean cost obtained applying the control u ; in this way we choose more often controls that are assumed to lead to a lower cost. Furthermore choosing $\beta \rightarrow \infty$, this strategy reduces to the greedy one.

One last technique that we might use is *optimism in the face of uncertainty*, a method that weights the greedy strategy by preferring inputs that have been explored fewer times. We chose in fact the control with the *lowest confidence bound* defined as

$$\text{LCB}(u) = Q_t(u) \quad (3.28)$$

where the uncertainty on $Q_t(u)$ is

$$Q_t(u) - L \sqrt{\frac{2 \log(t)}{n_t(u)}}$$

where $L = \max\{|\ell(u)|\}$ and $n_t(u)$ is the number of times that u was selected.

3.4.2 Actor-critic

The *actor-critic methods* can be regarded as a generalized policy iteration (sec. 3.2.2) extended to unknown Markov decision process. In that case we alternated a policy evaluation with a policy improvement in order to converge to the optimal policy π^* .

The extension to unknown MDPs is carried out using Monte Carlo (sec. 3.3.1) or temporal difference (sec. 3.3.2) to evaluate the policy. Since the exact evaluation of V^π can take infinitely many samples, the idea is to improve the policy π based on an approximation of the value function (so the MC/TD doesn't run until convergence, but for a given amount of time).

The value function is called *critic* because it evaluates the policy (so measures "how good" the learner perform), while in turn the *actor* is the policy itself.

In general the policy used to generate transitions is typically not the same that is evaluated and improved: the *behavior policy* mixes a small amount of exploration into a *target policy*. We note that the generalized policy iteration can generate policies that are worse then the previous ones.

SARSA: critic In the *SARSA method* we use TDo to learn the *critic* (so the action-value function Q) from on-policy samples. After each transition (x, u, ℓ, x', u') ⁴ the update of Q is performed as

$$\begin{aligned}\delta(Q_k) &= \ell(x, u) + \gamma Q(x', u') - Q(x, u) \\ Q_{k+1}(x, u) &= Q_k(x, u) + \alpha_k \delta(Q_k)\end{aligned}\tag{3.29}$$

We note the similarity with Q learning (3.25), but in this case we use $Q(x', u')$ instead of $\min_z Q(x, z)$. Note that if the policy π is fixed, then SARSA and TDo are exactly the same thing.

As for temporal difference, we can extend this concept to multi-step evaluation: using $TD\lambda$ gives the method $SARSA(\lambda)$.

From TD, SARSA inherits the possibility of diverging in off-policy settings.

SARSA: actor The goal of an actor-critic method is also to reach the optimal policy, a goal of the *actor* part. The simplest idea to improve the policy is by acting greedily with respect to the action-value function Q (that can be computed on the fly).

To ensure exploration we might use ε -greedy strategies, but in order to let SARSA converge to the optimal action-value function $Q^*(x, u)$ we need to ensure that

1. we are *Greedy in the Limit with Infinite Exploration (GLIE)* policies, so all state-action (x, u) are visited infinitely many times and the policy converges to the greedy one ($\varepsilon \xrightarrow{x \rightarrow \infty} 0$);
2. the sequence of step sizes α_k satisfies the Robbins-Monroe conditions.

⁴using reinforcement learning naming it becomes (s, a, r, s', a') , thus the name *SARSA*.