
Lecture notes on Control Systems and Reinforcement Learning

Written by
Manuel Hinz

`mh@mssh.dev` or `s6mlhinz@uni-bonn.de`

Lecturer

Prof. Dr. Jochen Garcke
`garcke[at]math.uni-bonn.de`



University of Bonn
Summer semester 2025
Last update: May 15, 2025

Contents

Chapter 0	Manuel's notes	2
0.1	Organization	2
Chapter 1	Introduction to optimal control	3
1.1	State Space Models	4
1.2	Linear State Space Model	5
1.3	State Space Models in continuous Time	5
1.4	Value iteration	13
1.5	Policy iteration	14
1.6	Exploration	15
1.7	Linear Quadratic Regulator, Revisited	16
1.8	Approximate Q -functions	16
1.9	Bandits	17
1.10	Other control formulations	17
1.11	Geometry in continuous time	19
1.12	Optimal control in continuous time	20
1.13	Linear quadratic regulator revisited (once more)	21
Chapter 2	ODE methods for algorithm design	23
2.1	ODE methods for algorithm design	23
2.2	Euler's method once more	25
2.3	Optimization	25
2.4	Qausi stochastic approximation	28
2.5	Approximate Policy Improvement	30
Journal		35
Bibliography		36

Chapter 0:

Manuel's notes

Warning

These are unofficial lecture notes written by a student. They are messy, will almost surely contain errors, typos and misunderstandings and may not be kept up to date! I do however try my best and use these notes to prepare for my exams. Feel free to email me any corrections to mh@mssh.dev or s6mlhinz@uni-bonn.de.
Happy learning!

Many thanks to Vincent for his feedback and some corrections!

General Information

- Basis: [Basis](#)
- Website: <https://ins.uni-bonn.de/teachings/ss-2025-467-v5e1-advanced-topics/>
- Time slot(s): **Tuesday: 14-16** SR 2.035 and **Thursdays: 16-18** SR 2.035
- Exams: ?
- Deadlines: No exercise sheets / tutorials

0.1 Organization

- Focused on ingredients, won't get to the current state of the art
- Some algorithmic / numerical background (Euler method is fine)
- Control Problems (Steering the bike / car)

Start of lecture 01
(10.4.2025)

The main source for this course is [\[2\]](#). We will follow this somewhat closely, especially in the first part of the course!

Chapter 1:

Introduction to optimal control

1. u is the control (input / action)
2. y observations (outputs)
3. $\phi : Y \rightarrow U$ policy
4. ff feed forward control (plan we had)

Interactions with the outside world might be hidden in the observations. Typically ff is in regard to some reference state. There might be some disturbances (holes in the road, ...).

The overall aim is to find a policy ϕ that sticks close to $r(k), k \geq 0$.

t is continuous, k is step
by step / iterative

$$u(k) = u_{\text{ff}}(k) + U_{\text{fb}}(k)$$

where u_{ff} is the planing to reach the overall goal and u_{fb} actual steering, updated "all the time".
Some examples from the book:

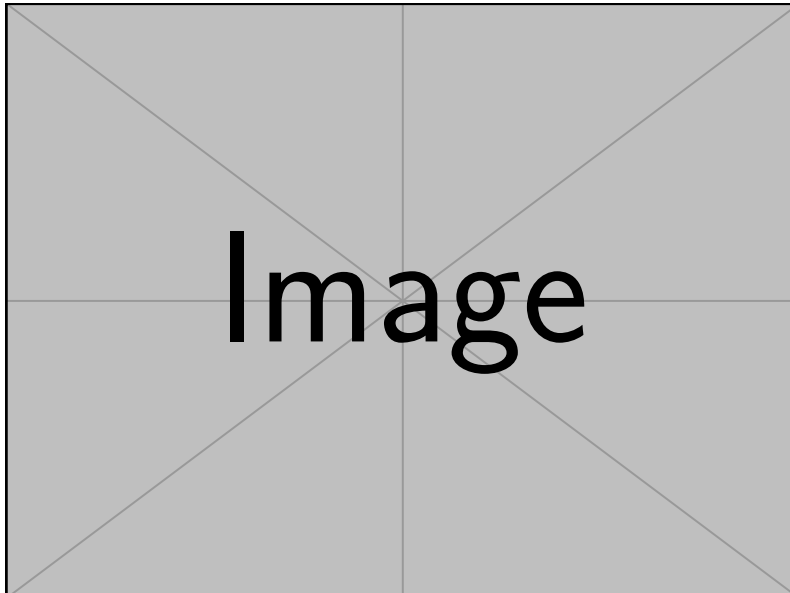


Figure 1.1: Sketch 1.01

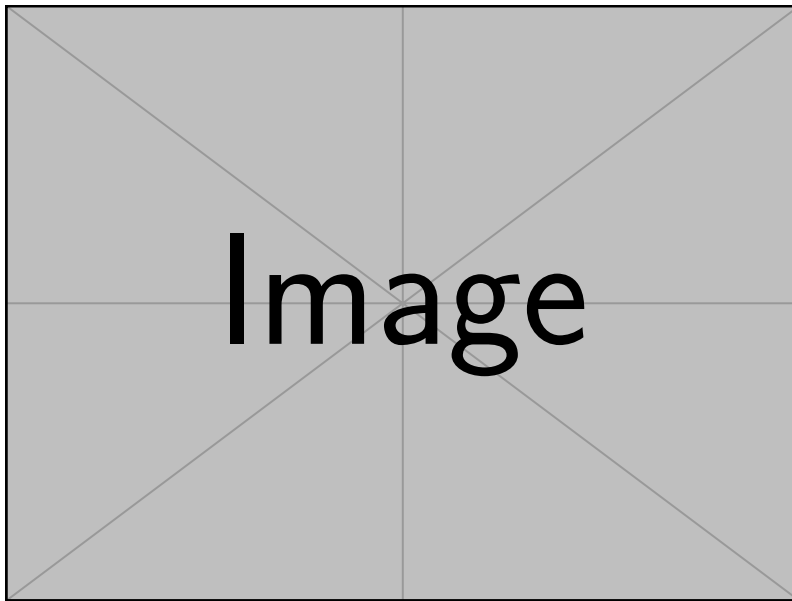


Figure 1.2: Sketch 1.02: Mountain car

Difference: In Reinforcement learning, we don't start with a model / ode.
Some part of reinforcement learning works model-free (i.e. assumes the model only implicitly)

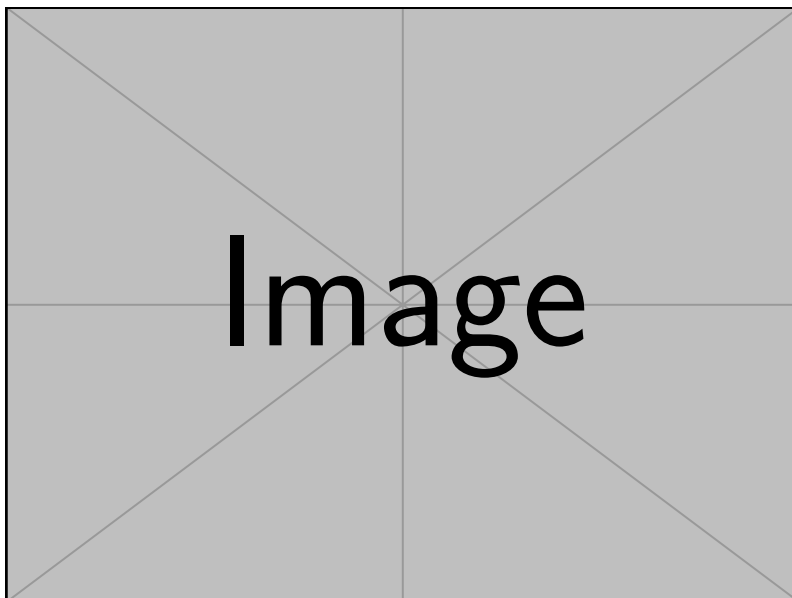


Figure 1.3: Sketch 1.03: cart pole / inverted pendulum

Next example: Acrobot (more than one equilibrium)

1.1 State Space Models

We have some

- state space $X, x \in X$
- action space $U, u \in U$
- action at step $k : u(k) \in U(k)$, i.e. we might have some constraints

- observation space $Y, y \in Y$

Definition 1.1. Given state, action and observation spaces X, U, Y , a state space model is defined by

$$x(k+1) = \mathcal{F}(x(k), u(k)) \quad (1)$$

$$y(k) = \mathcal{C}(x(k), u(k)) \quad (2)$$

$x(k)$ might include the past, might be useful for the stock trading problem

Remark. Overcomplicating problems by loading lots of information into the state space, might make the problem harder!

1.2 Linear State Space Model

$$x(k+1) = Fx(k) + Gu(k) \quad (3)$$

$$y(k) = Cx(k) + Du(k) \quad (4)$$

Remark. The representations (in terms of the matrices) might not be unique!

Common scenario for (3) is to keep $x(k)$ near the origin. You have to think about robustness of the system. Disturbances should be handled by the system.

$$u(k) = -Kx(k).$$

Consider a disturbance under the same control:

$$u(k) = -Kx(k) + v(k)$$

inserting this into (3) yields

$$x(k+1) = (F - GK)x(k) - Gv(k)$$

$$y(k) = (C - DK)x(k) + Dv(k)$$

Closed vs open loop: In closed loops we don't change our course based on observations, while in open loop systems we do.

1.3 State Space Models in continuous Time

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

for $x \in \mathbb{R}^n, u \in \mathbb{R}^m$. We often write u_t, x_t for u, x at time t . If f is linear we get

$$\frac{d}{dt}x = Ax + Bu$$

$$y = Cx + Du$$



Figure 1.4: Sketch 1.04

To discretize we use the forward Euler method. Given time interval Δ

$$x(k+1) = x(k) + \Delta f(x(k), u(k))$$

so in (1) $\mathcal{F}(x, u) = x + \Delta f(x, u)$. Using Taylor

$$x_{t+\Delta} = x_t \Delta f(x, u) + O(\Delta^2)$$

For the linear model we get $F = I + \Delta A$

$$x(k+1) = x(k) + \Delta A x(k) + \underbrace{\Delta B}_{=:G} u(k)$$

For now fix some policy ϕ , so $u(k) = \phi(x(k))$:

$$x(k+1) = \mathcal{F}(x(k))$$

Assumption 1.2. The state space X is equal to \mathbb{R}^n or a closed subset of \mathbb{R}^n .

Definition 1.3. An equilibrium x^e is a state at which is system is frozen:

$$x^e = \mathcal{F}(x^e).$$

Definition 1.4. Given a cost function $C : X \rightarrow \mathbb{R}_+$ and a policy ϕ we define

$$J_\phi(x) = J(x) = \sum_{k=0}^{\infty} C(x(k)), \quad x(0) = x$$

This is called total cost or value function of the policy ϕ .

Given x^e , we usually assume $C(x^e) = 0$. Generally, we consider a discount factor γ^k in front of $C(x(k))$.

Definition 1.5. Denote by $\mathcal{X}(k; x_0)$ the state step k with initial condition x_0 and following fixed policy ϕ . The equilibrium x^e is stable in the sense of Lyapunov if for all $\epsilon > 0 \exists \delta > 0$ s.t. $\|x_0 - x^e\| < \delta$, then

$$\|\mathcal{X}(k; x_0) - \mathcal{X}(k; x^e)\| < \epsilon \forall k \geq 0$$

The same concept with a different sign comes up in RL under the term reward

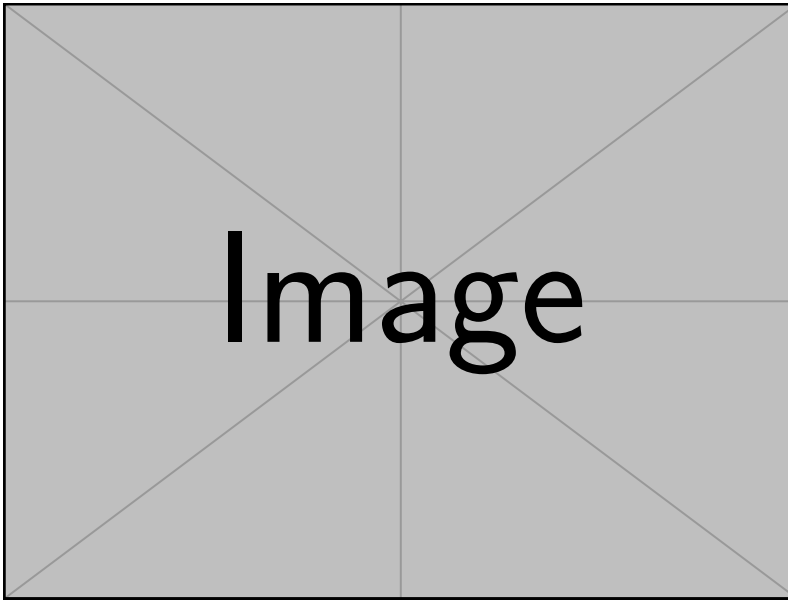


Figure 1.5: Sketch about Lyapunov stability

Definition 1.6. An equilibrium is said to be **asymptotically stable** if x^e is stable in the sense of Lyapunov and for some $\delta_0 > 0$, whenever $\|x_0 - x^e\| < \delta_0$, it follows

$$\lim_{k \rightarrow \infty} \mathcal{X}(k, x_0) = x^e.$$

The set of x_0 for which this holds is the **region of attraction** for x^e . An equilibrium is **globally asymptotically stable** if the region of attraction is X .

Definition 1.7 (Lyapunov function). A function $V : X \rightarrow \mathbb{R}_+$ is called **Lyapunov function**. We frequently assume V is **inf-compact**, i.e.: it holds

$$\forall x^0 \in X : \{x \in X \mid V(x) \leq V(x^0)\} \text{ is a bounded set.}$$

Remark. There is some variability in the definition of Lyapunov functions! We often assume $V(x)$ is large if x is large.

Sublevel sets:

$$S_V(r) = \{x \in X \mid V(x) \leq r\}.$$

One can see with V being inf-compact $S_V(r)$ is either

- empty
- the whole domain X
- a bounded subset of X .

We usually want to avoid this

Usually, $S_V(r) = X$ is impossible, a common assumption is **coersiveness**:

$$\lim_{\|x\| \rightarrow \infty} V(x) = \infty.$$

Example. • $V(x) = x^2$, coercive

- $V(x) = \frac{x^2}{(1+x)^2}$, not coercive, but inf-compact $r > 1 : S_V(r) = \mathbb{R}$, $r < 1 : S_V(r) = [-a, a]$ with $a = \sqrt{\frac{r}{1+r}}$
- $V(x) = e^x$ is neither

Lemma 1.8. Suppose that the cost function C and the value function J from definition 1.5 are non-negative and finite valued.

this is a assumption on the value function

1. $J(x(k))$ is non-increasing in k and $\lim_{k \rightarrow \infty} J(x(k)) = 0$ for each initial condition.
2. In addition let J be continuous, inf-compact and vanishing only at x^e . Then for each initial condition

$$\lim_{k \rightarrow \infty} x(k) = x^e$$

Proof. Consider $J(x) = \sum_{k=0}^{\infty} c(x(k))$, then

$$\begin{aligned} J(x) &= c(x) + \sum_{k=1}^{\infty} c(x(k)) \\ &= c(x) + \sum_{k=0}^{\infty} c(x^+(k)); \quad x^+(0) = \mathcal{F}(x) \\ &= c(x) + J(\mathcal{F}(x)) \end{aligned}$$

This is the dynamic programming principle for a fixed policy. It is also called Bellmann equation. For 1. from this it follows

$$J(x(k+1)) + c(x) - J(x(k)) = 0$$

summing up from $k = 0$ up to $N - 1$

$$\begin{aligned} J(x) &= J(x(N)) + \sum_{k=0}^{N-1} c(x(k)) \\ &\implies \text{non-increasing} \end{aligned}$$

We are separating one step!

This is the same Bellman from the curse of dimensionality!

Taking the limit

$$= \lim_{N \rightarrow \infty} \left[J(x(N)) + \sum_{k=0}^{N-1} c(x(k)) \right] = \left[\lim_{N \rightarrow \infty} J(x(N)) \right] + J(x)$$

using $J(x)$ is finite gives (i).

For 2. with $r = J(x)$, we get $x(k) \in S_J(r) \forall k$. Now suppose $\{x(k_i)\}$ is a convergent subsequence of the trajectory with limit x^∞ . Then $J(x^\infty) = \lim_{i \rightarrow \infty} J(x(k_i)) = 0$ by the continuity of J . We assumed $J(x) = 0 \iff x^e = x \implies x^\infty = x^e$. Finally, the assumption follows, since each convergent subsequence reach the same value x^e .

□

Definition 1.9 (Poisson's inequality). Let $V, c : X \rightarrow \mathbb{R}_+$ and $\eta \geq 0$. Then Poisson's inequality states that

We often assume $\eta = 0$

$$V(\mathcal{F}(x)) \leq V(x) - c(x) + \eta.$$

Proposition 1.10. Suppose the Poisson inequality holds with $\eta = 0$. Additionally V shall be continuous, inf-compact and it shall have a unique minima at x^e . Then x^e is stable in the sense of Lyapunov (sitsoL).

Proof.

$$\bigcap \{S_V(r) \mid r > V(x^e)\} = \{S_V(r)|_{r=V(x^e)}\}^{\text{unique minimizer}} = \{x^e\}.$$

Using compactness we get: For each $\epsilon > 0$, we can find some $r > V(x^e)$ and some $\delta < \epsilon$ s.t.

$$\{x \in X \mid \|x - x^e\| < \delta\} \subset S_V(r) \subset \{x \in X \mid \|x - x^e\| < \epsilon\}$$

If $\|x_0 - x^e\| < \delta$, then $x_0 \in S_V(r)$ and hence $x(k) \in S_V(r)$ since $V(x(k))$ is non-increasing. With the second inclusion we see

$$\|x(k) - x^e\| < \epsilon \forall k$$

This gives sitsoL.

□

Proposition 1.11 (Comparison theorem). *Poisson's inequality implies*

1. For each $N \geq 1$ and $x = x(0)$

$$V(x(N)) + \sum_{k=0}^{N-1} c(x(k)) \leq V(x) + N\eta$$

2. If $\eta = 0$, then $J(x) \leq V(x) \forall x$

3. Assume $\eta = 0$ and V, c are continuous. Suppose that c is inf-compact and vanishes only at the equilibrium x^e . Then x^e is globally asymptotically stable.

We don't write that explicitly, but we don't start in x^e !

Proof. 1.

$$V(x(k+1)) - V(x(k)) + c(x(k)) \leq \eta$$

summing up from 0 to $N-1$:

$$V(x(N)) - V(x(0)) + \sum_{k=0}^{N-1} c(x(k)) \leq N\eta$$

2. for $\eta = 0$ the above is ≤ 0 , so $\sum_{k=0}^{N-1} c(x(k)) \leq V(x(0)) - V(x(N)) \leq V(x(0))$ where the LHS converges to $J(x(0))$ for $N \rightarrow \infty$

3. Show sitsoL, with $\eta = 0$ it follows from definition 1.9 that $V(x) \geq c(x)$, which gives V is also inf-compact. c is vanishing only at x^e , so $V(x(k))$ is strictly decreasing. When $x(k) \neq x^e$, implies $V(x(k)) \downarrow V(x^e)$ for each $x(0)$. Further

This is important!

$$V(x^e) < V(x(0)) \quad \forall x(0) \in X \setminus \{x^e\}.$$

So it is a unique minimum. V has therefore the properties of proposition 1.10, which gives sitsoL. For global: with 1. we get

$$\lim_{k \rightarrow \infty} c(x(k)) = 0$$

and assumptions give us by lemma 1.8 that $x(k) \rightarrow x^e$ as $k \rightarrow \infty$. So, we converge from any initial condition, which gives global asymptotical stability. □

Proposition 1.12. *Suppose that $V(\mathcal{F}(x)) = V(x) - c(x)$. Further, we assume that*

1. J is continuous, inf-compact, vanishing only at x^e

2. V is continuous

Then $J(x) = V(x) - V(x^e)$.

Proof. As before we sum up:

$$V(x(N)) + \underbrace{\sum_{k=0}^{N-1} c(x(k))}_{J(x(N-1)) \xrightarrow{N \rightarrow \infty} J(x)} = V(x).$$

Lemma 1.8 together with the continuity of V implies that

$$V(x(N)) \rightarrow V(x^e) \quad \text{as } N \rightarrow \infty.$$

This gives

$$V(x^e) + J(x) = V(x) \quad \square$$

Start of lecture 03
(17.04.2025)

Example (Linear state space model). Setting $x(k+1) = \mathcal{F}(x(k))$, now with linear dynamics:

$$x(k+1) = Fx(k) = F^{k+1}x(0) = F^{k+1}x.$$

Assume quadratic cost $c(x) = x^\top Sx$, where S is symmetric and positive definite. Observe

$$c(x(k)) = (F^k x)^\top S F^k x$$

Summing up yields

$$J(x) = x^\top \underbrace{\left[\sum_{k=0}^{\infty} (F^k)^\top S F^k \right]}_{=:M} x$$

This satisfies a linear fixed point equation:

$$M = S + F^\top M F \quad (5)$$

This is also called
discrete time
Lyapunov equation

One can show for the linear state space model, that the following are equivalent:

1. the origin is asymptotically stable
2. the origin is globally asymptotically stable
3. Each eigenvalue λ of F satisfies $|\lambda| < 1$
4. (5) admits a solution M positive semi-definite for any S positive semidefinite.

Reference: [1]

Consider 1.1 without y

$$y(k+1) = \mathcal{F}(x(k), u(k))$$

with

$$c : X \times U \rightarrow \mathbb{R}_+.$$

The total cost J_ϕ for a given ϕ given $u(k) = \phi(x(k))$ is

$$J_\phi(x) = \sum_{k=0}^{\infty} c(x(k), u(k)).$$

The optimal value function is the minimum over all controls

$$J^*(x) = \min_{\underline{U}=[u(0), u(1), \dots]} \sum_{k=0}^{\infty} c(x(k), u(k)), \quad x(0) = x \in X \quad (6)$$

This describes the
optimal control policy
(OCP)

Remark. The minimizer might not be unique! In harder settings this might need to be an inf!

Goal: Find a control sequence that achieves the minimum.

Computationally we can't expect to calculate J_ϕ exactly, but we will approximate it.

and the corresponding
policy

Remark. We are in the infinite horizon setting (infinite time steps) to talk about the stability. For this it is important that the equilibrium has cost 0. Without an equilibrium we can also think about discounted value functions

$$J_\phi = \sum_{k=0}^{\infty} \gamma^k c(x(k), u(k))$$

We will see later that it holds for the sequence x^* achieving the minimum

$$J^*(x^*(k)) = c(x^*(k), u^*(k)) + J^*(x^*(k+1))$$

which is definition 1.9 with $\eta = 0$ and equality.

Proposition 1.11 implies, under some conditions, that x^e is globally asymptotically stable.

Under the following assumptions J^* is finite:

1. there is a (target) state x^e that is an equilibria for some control $F(x^e, u^e) = x^e$
2. $c \geq 0, c(x^e, u^e) = 0$
3. for any initial condition $x(0) = x$ there is a control sequence \underline{u} and a time T , such that $x(T) = x^e$ for $x(0) = x$ using control \underline{u} .

This is sometimes called controllability

Example (Linear Quadratic Regulator). Consider linear dynamics 3 from the first lecture with quadratic cost $c(x, u) = x^\top Sx + u^\top Ru$ with S positive semi-definite and R positive definite.

Reminder: $u = -Kx$.

If there is a policy for which J^* is finite, then

$$J^*(x) = x^\top M^* x$$

with M^* positive semi-definite and

$$\phi^*(x) = -K^*(x)$$

with K^* depends on M^*, R, F, G .

and implicitly on c

Bellmann equation

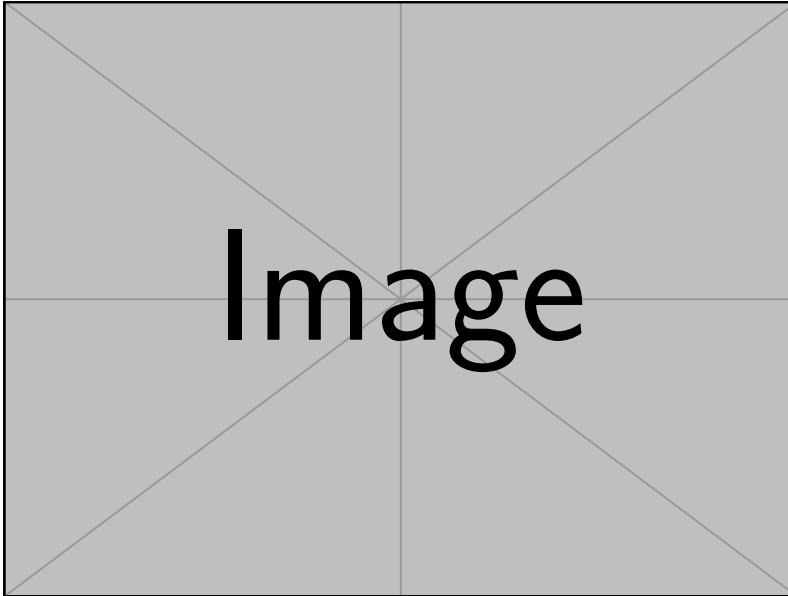


Figure 1.6: Sketch 1.06; Principle of optimality

Observation:

$$\begin{aligned} J^*(x) &= \min_{\underline{u}} \left[\sum_{k=0}^{k_m-1} c(x(k), u(k)) + \sum_{k_m}^{\infty} c(x(k), u(k)) \right] \\ &= \min_{u[0, \dots, k_m-1]} \left[\sum_{k=0}^{k_m-1} c(x(k), u(k)) + \underbrace{\min_{u[k_m, \dots]} \sum_{k_m}^{\infty} c(x(k), u(k))}_{=J^*(x(k_m))} \right] \end{aligned}$$

This gives

$$J^*(x) = \min_{u[0, \dots, k_m-1]} \left[\sum_{k=0}^{k_m-1} c(x(k), u(k)) \right] + J^*(x(k_m)).$$

which can be seen as a kind of fix point equation

With $k_m = 1$ we have shown the following theorem

Theorem 1.13 (Bellmann equation, Dynamic Programming equation). Assume that J^* is finite and optimal control u^* solving (6) exists. Then the value function satisfies

$$J^*(x) = \min_u \{c(x, u) + J^*(\mathcal{F}(x, u))\} \quad (7)$$

Suppose the minimum is unique for each x and let $\phi^*(x)$ denote the minimum feedback law at x . Then the optimal control is expressed as

$$u^*(k) = \phi^*(x^*(k)).$$

Definition 1.14 (Q-function). The function of two variables within the minimum in (7) is called Q-function.

$$Q^*(x, u) = c(x, u) + J^*(\mathcal{F}(x, u))$$

In the optimal case we write Q^* . Thus

$$J^*(x) = \min_{\bar{u}} Q^*(x, \bar{u}).$$

The optimal feedback law is then

$$\phi^*(x) \in \operatorname{argmin}_u Q^*(x, u).$$

The Q-function solves the fixed point equation

$$Q^*(x, u) = c(x, u) + \min_u Q^*(\mathcal{F}(x, u), u).$$

This already gives a hint for an algorithm coming later next lecture.

Remark. In RL the difference is that we don't know the model, we only observe state action pairs. This motivates the Q-function.

Definition, which is not so useful for the analysis, but for the practical application!

Some concepts from Reinforcement Learning

Actors and critic:

Given is a parameterized family of policies $\{\phi^\theta \mid \theta \in \mathbb{R}^d\}$. the actors. For each θ , observe the trajectories by their states x and actions u determined by their policy.

The critic approximates the associated value function \tilde{J}_θ . Aim for the minimum

$$\theta^* = \operatorname{argmin}_\theta \langle v, \tilde{J}_\theta \rangle,$$

where the weight vector $v \geq 0$ reflects the weighting of the states. $v(x)$ is large for *important* states.

scalar product in \mathbb{R}^n (all states?)

Temporal differences:

$$J_\theta(x(k)) = c(x(k), u(k \mid \theta)) + J_\theta(x(k+1))$$

Look for an approximation \hat{J} for which the error is small (w.r.t. the equality above).

Temporal differences are

$$D_{k+1}(\hat{J}) := -\hat{J}(x(k)) + \hat{J}(x(k+1)) + c(x(k), u(k)).$$

After N samples

$$\Gamma(\hat{J}) := \frac{1}{N} \sum_{k=0}^{N-1} D_{k+1}(\hat{J})^2.$$

We can optimize / minimize this.

There is a whole class of TD algorithms and those fit into the actors critic approach!

What changes, or what is the information gain

Start of lecture 04
(22.04.2025)

1.4 Value iteration

We approximate J^* by a sequence of V^k given an initial value function V^0 .

$$V^{k+1}(x) = \min_u \{c(x, u) + V^k(\mathcal{F}(x, u))\}, x \in X, k \geq 0$$

This is called **value iteration** often shortened to VI.

Algorithm 1 Value iteration

Input: Start with an initial value function V^0

Output: Estimates V^{k+1}

$n = 0$

while not good enough **do**

Value function improvement to obtain next value function

$$V^{k+1}(x) = \min_u \{c(x, u) + V^k(\mathcal{F}(x, u))\}, x \in X, k \geq 0$$

end while

For infinite state spaces we will have to fix this algorithm for memory related reasons

Proposition 1.15. Let V^0 be chosen with non-negative entries and $V^0(x^e) = 0$. Further, we assume

1. X, U are finite sets
2. c is non-negative and vanishes only at (x^e, u^e) , and J^* is finite valued.

Then there is $n_0 \geq 1$ such that

$$V^k(x) = J^*(x), x \in X, k \geq n_0.$$

Proof. Let $\phi^*(x)$ be an optimal policy, and let $n_0 \geq 1$ denote the value such that

$$(x^*(k), u^*(k)) = (x^e, u^e)$$

for $k \geq n_0$. This exists since J^* is finite.

Using the principle of optimality (6) we can show

$$V^n(x) = \min_{u[0, \dots, n-1]} \left\{ \sum_{k=0}^{n-1} c(x(k), u(k)) + V^0(x(n)) \right\}, x(0) \in X \quad (8)$$

This gives

$$V^n(x) \leq \sum_{k=0}^{n-1} c(x(k), u(k)) + V^0(x(n)) \text{ for all } u \text{ including } u(k) = \phi^*(k)$$

$$\stackrel{n \geq n_0}{=} J^*(x) + V^0(x^e) = J^*(x)$$

For such n , the inequality must be an equality, due to (8) and the use of the optimal policy. \square

VI provides a sequence of policies ϕ^n

$$\phi^n(x) \in \operatorname{argmin}_u \{c(x, u) + V^n(\mathcal{F}(x, u))\}.$$

If we assume that V^0 is non-negative and satisfies poisson's inequality(1.9) for some $\eta \geq 0$

$$V^0(\mathcal{F}(x, u)) \leq V^0(x) - c(x, \phi^0(x)) + \eta, x \in X$$

then we get the following statement

We really exploit the finiteness!

Proposition 1.16. Suppose that V^0 is non-negative and it holds

$$\begin{aligned} \min_u (c(x, u) + V^0(\mathcal{F}(x, u))) &= \{c(x, u) + V^0(\mathcal{F}(x, u))\} |_{u=\phi^0(x)} \\ &\leq V^0(x) + \eta, \quad x \in X \end{aligned}$$

Then a corresponding bound holds for each n

$$\{c(x, u) + V^n(\mathcal{F}(x, u))\} |_{u=\phi^n(x)} \leq V^n(x) + \eta_n, \quad x \in X,$$

where η_i is non-increasing:

$$\eta \geq \eta_0 \geq \eta_1 \dots$$

Proof. Write $B^n(x) = V^{n+1}(x) - V^n(x)$

$$\eta_n := \sup_x B^n(x).$$

This is (connected to?)
the Bellman error

Value iteration gives

$$\begin{aligned} \{c(x, u) + V^n(\mathcal{F}(x, u))\} |_{u=\phi^n(x)} &= \min_u \{c(x, u) + V^n(\mathcal{F}(x, u))\} \\ &= V^{n+1}(x) = V^n(x) + B^n(x) \\ &\leq V^n(x) + \eta_n \end{aligned}$$

To show that the η are non-increasing, we consider

$$V^1(x) = \{c(x, u) + V^0(\mathcal{F}(x, u))\} |_{u=\phi^0(x)} \stackrel{\text{Assumption}}{\leq} V^0(x) + \eta$$

which gives $B^0(x) \leq \eta \forall x \implies \eta_0 \leq \eta$.

For $n \geq 1$ The trick is using the old control in the second line:

$$\begin{aligned} V^n(x) &= \{c(x, u) + V^{n-1}(\mathcal{F}(x, u))\} |_{u=\phi^{n-1}(x)} \\ V^{n+1}(x) &\leq \{c(x, u) + V^n(\mathcal{F}(x, u))\} |_{u=\phi^{n-1}(x)} \end{aligned}$$

So,

$$V^{n+1}(x) - V^n(x) \leq \{V^n(\mathcal{F}(x, u)) - V^{n-1}(\mathcal{F}(x, u))\} |_{u=\phi^{n-1}(x)} \leq \eta_{n-1}.$$

Hence, $\eta_n = \sup_x B^n(x) \leq \eta_{n-1}$. □

Now consider $\eta = 0$, so for each n

$$\{c(x, u) + V^n(\mathcal{F}(x, u))\} |_{u=\phi^n(x)} \leq V^n(x)$$

with proposition 1.11 it follows

$$J^* \leq V^n(x), \quad x \in X,$$

where J^* is the total cost using policy ϕ^n .

One view of policy iteration is the focus on updating the policy function!

1.5 Policy iteration

Start with an initial policy $\phi^0, n = 0$

- Compute the total cost for the policy ϕ^n , this is called policy evaluation

$$J^n(x) = \sum_{k=0}^{\infty} c(x(k), u(k)), \quad u(k) = \phi^n(x(k)) \forall x \in X$$

- perform **policy improvement** to obtain the next policy

$$\phi^{n+1}(x) \in \underset{u}{\operatorname{argmin}} \{c(x, u) + J^n(\mathcal{F}(x, u))\}, \quad x \in X$$

- while *not good enough*

This is sometimes also called Howard's algorithm.

Remark. The first step is some linearization and the second is the update. Like a generalization of Newton's method

Algorithm 2 Policy iteration

Input: Start with an initial policy ϕ^0

Output: Estimates $J^n(x), \phi^{n+1}(x)$

$n = 0$

while not good enough **do**

 Compute the total cost for the policy ϕ^n , this is called policy evaluation

$$J^n(x) = \sum_{k=0}^{\infty} c(x(k), u(k)), \quad u(k) = \phi^n(x(k)) \quad \forall x \in X$$

 perform **policy improvement** to obtain the next policy

$$\phi^{n+1}(x) \in \underset{u}{\operatorname{argmin}} \{c(x, u) + J^n(\mathcal{F}(x, u))\}, \quad x \in X$$

end while

Proposition 1.17. Suppose that J^0 for ϕ^0 is finite valued. Then for each $n \geq 0$

$$\{c(x, u) + J^n(\mathcal{F}(x, u))\}_{|u=\phi^{n+1}(x)} \leq J^n(x), \quad x \in X$$

and consequently, the value functions are non-increasing

$$J^0(x) \geq J^1(x) \geq \dots$$

Proof. Similar to the proof of proposition 1.16, where the non-increasing sequence again follows from proposition 1.11. \square

Here we always assumed that we can compute everything, especially \mathcal{F} and the infinite sum.

1.6 Exploration

In RL we learn from observations, each state-action pair, new state and observed cost gives us information. We need *good* and *useful* information.

Consider a policy that is not optimal, but has $x(k) \rightarrow x^e$ reasonably rapidly, where we assume $c(x^e, \cdot) = 0$. Typically we have continuity

$$\begin{aligned} \lim_{k \rightarrow \infty} D_{k+1}(\hat{J}) &= \lim_{k \rightarrow \infty} \left[-\hat{J}(x(k)) + \hat{J}(x(k+1)) + c(x(k), u(k)) \right] \\ &= -\hat{J}(x^e) + \hat{J}(x^e) + 0 = 0. \end{aligned}$$

This is not much information, one cannot further improve the policy!

$$\Gamma^\epsilon(\hat{J}, x^i) = \frac{1}{N_\epsilon} \sum_{k=0}^{N_\epsilon-1} [D_{k+1}(\hat{J})]^2, \quad x(0) = x^i$$

To avoid getting *small* information from long trajectories, one can take a couple of shorter ones.

$$\hat{\Gamma}(\hat{J}) = \frac{1}{M} \sum_{i=1}^M \Gamma^\epsilon(\hat{J}; x^i)$$

How to choose x^i is current research. Much of the theoretical research assume that “every state is assumed regularly”, which is nice for results, but not so nice realistic in most applications.

Another way to get more diverse information is to use exploration. Namely one modifies the trajectories, not strictly follows ϕ^n .

$u(k) = \hat{\phi}(x(k), \zeta(k))$, where $\zeta(k)$ is some form of noise. Typically

1. $\hat{\phi}(x(k), \zeta(k)) = \phi^\theta(k)$ for *most* k
2. Choose action to explore the state-action space (e.g. randomly) the other times

this is also sometimes called off-policy and on-policy

Generally, the trajectory to gather information stems from a different policy than the current estimate ϕ^θ . This dilemma is called the exploration-exploitation dilemma.

Start of lecture 05
(24.04.2025)

1.7 Linear Quadratic Regulator, Revisited

We had $J^*(x) = x^\top M^* x$ and quadratic costs, $c(x, u) = x^\top S x + u^\top R u$.

For the Q -function:

$$Q^*(x, u) = c(x, u) + J^*(Fx + Gu).$$

An optimal policy ϕ is a minimum over Q w.r.t. u :

$$0 = \nabla_u Q^*(x, u^*) = 2Ru^* + 2G^\top M^*(Fx + Gu^*)$$

Assuming R is positive definite; then $R + G^\top M^* G$ is positive definite and therefore invertible.

$$K^* = [R + G^\top M^* G]^{-1} G^\top M^* F$$

and

$$\phi^*(x) = -Kx.$$

To obtain M^* we can solve a fixed point equation called the algebraic Riccati equation

$$M^* = F^\top \left(M^* - M^* G [R + G^\top M^* G]^{-1} G^\top M^* F + S \right) \quad (9)$$

This is a hint, we will prob. revisit this later

1.8 Approximate Q -functions

Consider a family of Q -functions $\{Q^\theta \mid \theta \in \mathbb{R}^d\}$ to approximate Q^* . Classically used is a linear parametrization

$$Q^\theta(x, u) = \theta^\top \psi(x, u), \quad \theta \in \mathbb{R}^d$$

where $\psi_i : X \times U \rightarrow \mathbb{R}$, $1 \leq i \leq d$ is some set of basis functions. Given Q^θ we have $\phi^\theta(x) \in \operatorname{argmin}_u Q^\theta(x, u)$, $x \in X$.

Policy iteration for Q -functions:

1. obtain θ^n to get an approximation of Q^{θ^n} where $Q^{\theta^n}(x, u) = c(x, u) + Q^{\theta^n}(x^+, u^+)$, $x^+ = \mathcal{F}(x, u)$, $u^+ = \phi^n(x^+)$
2. define new policy $\phi^{n+1}(x) := \phi^{\theta^n}$

Think kernels, finite element basis,...

Approximation since we do this sample-based in RL

As an alternative, consider dynamic programming equation from definition 1.14:

$$Q^*(x, u) = c(x, u) + \min_{\bar{u}} Q^*(\mathcal{F}(x, u), \bar{u}).$$

We follow a given/ observed state-action trajectory $(x(k), u(k))_{k=0}^N$

$$Q^*(x(k), u(k)) = c(x(k), u(k)) + Q^*(x(k+1), u(k+1))$$

The temporal difference / Bellmann error

$$D_{k+1}(Q^\theta) = -Q^\theta(x(k), u(k)) + c(x(k), u(k)) + Q^\theta(x(k+1), u(k+1))$$

If $Q^\theta = Q^*$ then $D_{k+1}(Q^\theta) = 0 \forall k$. In Q -learning algorithms, one chooses θ^n such that $D_{k+1}(Q^{\theta^n})$ is small in a suitable fashion. So we minimize θ to achieve this, i.e.

$$\Gamma^\epsilon(\theta) = \frac{1}{N} \sum_{i=0}^{N-1} [D_{k+1}(Q^\theta)]^2$$

1.9 Bandits

Theory of multi-armed bandits. One has to accept some loss through exploration in order to achieve(find) the best strategy. One exploits the learned strategy when choosing an action according to it.

In the control of dynamic systems one has for each state x (or $x(k)$) a multi-armed bandit.

1.10 Other control formulations

Discounted cost:

$$J^*(x) = \min_{\underline{u}} \sum_{k=0}^{\infty} \gamma^k c(x(k), u(k)), \quad x(0) \in X$$

where $\gamma \in (0, 1)$ is the discount factor.

Shortest Path Problem: Given $A \subset X$ define $\tau_A := \min\{k \geq 1 \mid x(k) \in A\}$.

$$J^*(x) = \min_u \sum_{k=0}^{\tau_A-1} \gamma^k c(x(k), u(k)), \quad x(0) = x.$$

This is problematic, since we might have longer path with lower cost ...

Proposition 1.18. *If J^* is finite valued, then it is the solution to the dynamic programming equation in the following sense:*

$$J^*(x) = \min_u \{c(x, u) + \gamma 1_{\{\mathcal{F}(x, u) \in A^c\}} J^*(\mathcal{F}(x, u))\}, \quad x \in X$$

where $1_{\{\dots\}}$ denotes an indicator function.

Proof.

$$\begin{aligned} J^*(x) &= \min_{\underline{u}} \left\{ c(x, \underline{u}) + \sum_{k=1}^{\tau_A-1} \gamma^k c(x(k), u(k)) \right\} \\ \tau_A=1 &\implies \sum=0 \quad \min_{u(0)} \left\{ c(x, u(0)) + \gamma 1_{\{x(1) \in A^c\}} + \min_{u[1, \dots, \tau_A]} \left\{ \sum_{k=1}^{\tau_A-1} \gamma^{k-1} c(x(k), u(k)) \right\} \right\} \\ &= \min_{u(0)} \{c(x, u(0)) + \gamma 1_{\{x(1) \in A^c\}} J^*(x(1))\} \end{aligned}$$

$c(x, u(0))$ since we're extracting the first element of the sum

where $x(1) = \mathcal{F}(x, u(0))$. □

To formulate this as a discounted problem

1. modify the cost function $c_A(x, u) = \begin{cases} c(x, u) & x \in A^c \\ 0 & x \in A \end{cases}$
2. modify the state dynamics $\mathcal{F}_A(x, u) = \begin{cases} \mathcal{F}(x, u) & x \in A^c \\ x & x \in A \end{cases}$

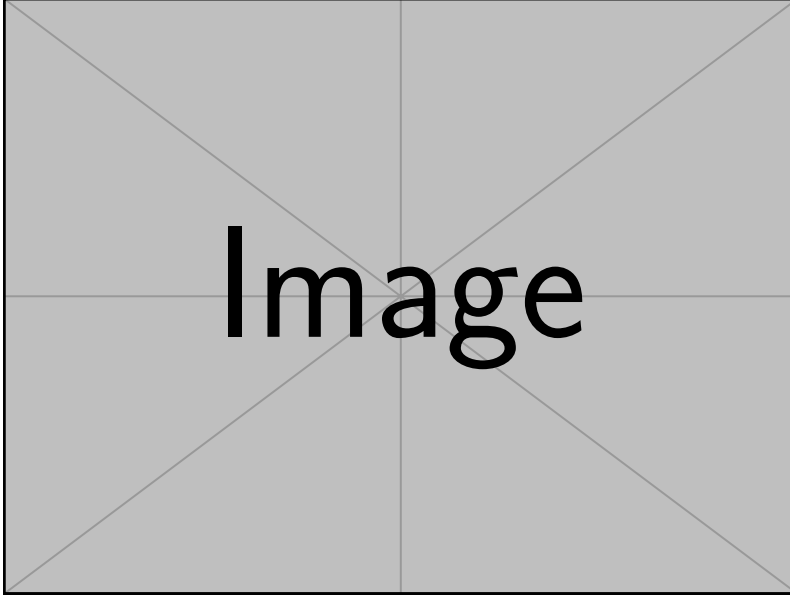


Figure 1.7: Sketch: mountain car value function

Can be numerically very hard, since the value function can be quite discontinuous, but not all value functions are that bad.

Finite Horizon Fix horizon $N \geq 1$ and define

$$J^*(x) = \min_{u[0,N]} \sum_{k=0}^N c(x(k), u(k)), \quad x(0) = x \in X.$$

We can connect to the optimal control problem by

1. enlarging the state space $x^a(k) = (x(k), \tau(k))$, where $\tau(k) = \tau(0) + k$, $k \geq 0$
2. modify the cost function $c^a((x, \tau), u) = \begin{cases} c(x, u) & \tau \leq N \\ 0 & \tau > N \end{cases}$

Then

$$J^*(x^a) = \min_{\underline{u}} \underbrace{\sum_{k=0}^{\infty} c^a(x^a(k), u(k))}_{J^*(x, \tau)}, \quad x^a(0) = (x, 0)$$

The Bellmann equation from theorem 1.13 now becomes

$$J^*(x, \tau) = \min_u \{c(x, u)1_{\{\tau \leq N\}} + J^*(\mathcal{F}(x, u), \tau + 1)\} \quad (10)$$

For $\tau > N$, it follows that $J^*(x, \tau) = 0$. This gives

$$J^*(x, N) = \min_u c(x, u) = \bar{c}(x).$$

kind of a boundary condition

So,

$$J^*(x, N-1) = \min_u \{c(x, u) + \bar{c}(\mathcal{F}(x, u))\}$$

repeating this backwards in time yields

$$J^*(x, 0) = J^*(x^a).$$

For the policy $\phi^*(x, \tau) \in \operatorname{argmin}_u \{c(x, u) + J^*(\mathcal{F}(x, u), \tau + 1)\}, \tau \leq N$ and

$$u^*(k) = \phi^*(x^*(k), k).$$

Start of lecture 06
(29.04.2025)

Model Predictive Control

Here, the policy is computed on-the-fly at each step of the state-action trajectory as a finite horizon problem. The control is

$$u(k) = \phi^{\text{mpc}}(x^*(k)) = \phi^*(x^*(k), 0),$$

where ϕ^* from the finite horizon setting (10) for small N .

Consider

$$J^{\text{mpc}}(x) = \sum_{k=0}^{\infty} c(x(k), u(k)), \quad x(0) = x, u(k) = \phi^{\text{mpc}}(x(k)).$$

Due to the finite horizon we are not optimal ...

Proposition 1.19. Consider $u(k)$ from above with

$$J^*(x; 0) = \min_{u[0, N-1]} \sum_{k=0}^{N-1} c(x(k), u(k)) + V^0(x(N)),$$

where $V^0 : X \rightarrow \mathbb{R}^+$ satisfies the assumption from proposition 1.16 with $\eta = 0$:

$$\min_u \{c(x, u) + V^0(\mathcal{F}(x, u))\} \leq V^0(x).$$

Then the total cost J^{mpc} is finite everywhere.

Proof. Using an equation from proposition 1.15:

$$V^N(x) = \min_{u[0, N-1]} \left\{ \sum_{k=0}^{N-1} c(x(k), u(k)) + V^0(x(N)) \right\}$$

and the definition of J^* from above we get $J^*(x, 0) = V^N(x)$ Proposition 1.16 then gives the bound

$$\{c(x, u) + V(\mathcal{F}(x, u))\}_{|u=\phi^{\text{mpc}}(x)} \leq V(x) = V^N(x)$$

This is also a version of a poisson inequality

From the Comparison theorem 1.11, it follows that J^{mpc} is finite. □

1.11 Geometry in continuous time

Consider $x(k+1) = \mathcal{F}(x(k))$, now in continuous time:

$$\frac{d}{dt}x_t = f(x_t) \text{ or } \frac{d}{dx}x = f(x)$$

$\mathcal{X}(t, x_0)$ is the solution to the differential equation above. Definition 1.5, 1.6 carry over.

$$\lim_{t \rightarrow \infty} \mathcal{X}(t, x_0) = x^e$$

Definition 1.20. A function $V : X \rightarrow \mathbb{R}_0^+$ is called Lyapunov function for global asymptotic stability if the following conditions hold:

- (i) $V \in C^1$
- (ii) V is inf-compact
- (iii) For any solution x , whenever $X_t \neq x^e$

$$\frac{d}{dt}V(x_t) < 0.$$

If $x_t = x^e$, we have $V(x_{t+s}) = V(x^e)$ for all $s \geq 0$, so $\frac{d}{dt}V(x^e) = 0$.

If we look back at the proof of proposition 1.10 and proposition 1.11 (iii), we can see that these also carry over to the continuous case. So we get

Proposition 1.21 (Extension of prop 1.11 (iii)). *If there exists a Lyapunov function after definition V 1.20, then the equilibrium x^e is globally asymptotically stable.*

Since we did not exploit the step-wise nature previously

The continuous version of Poisson's inequality is then

$$\langle \nabla V(x), f(x) \rangle \leq -c(x) + \eta \quad (11)$$

using the chain rule we get

$$\frac{d}{dt}V(x) \leq -c(x) + \eta$$

further observing

$$0 \leq V(x_T) = V(x_0) + \int_0^T \frac{d}{dt}V(X_t)dt \leq V(x_0) + T\eta - \int_0^T c(x_t)dt$$

we have shown

Proposition 1.22 (Continuous Comparison theorem). *If (11) holds for non-negative c, V, η , then we have*

$$V(X_t) + \int_0^T c(x_t)dt \leq V(x) + T\eta, \quad x_0 = x \in X, T > 0 \quad (12)$$

If $\eta = 0$

$$\int_0^\infty c(x_t)dt \leq V(x)$$

the total cost is bounded.

1.12 Optimal control in continuous time

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

with total cost for $\underline{u} = u[0, \infty)$

$$J(\underline{u}) = \int_0^\infty c(x_t, u_t)dt.$$

As before, we minimize over u and want J to be finite. We assume

$$f(x^e, u^e) = 0$$

for some u^e and

$$c(x^e, u^e) = 0$$

which yields that J is finite. As before

$$J^*(x) = \min_u \int_0^\infty c(x_t, u_t)dt, \quad x_0 = x \in X.$$

We extend the Bellmann equation to continuous times

$$\begin{aligned} J^*(x) &= \min_{u[0, \infty]} \left[\int_0^{t_m} c(x_t, u_t)dt + \int_{t_m}^\infty c(x_t, u_t)dt \right] \\ &= \min_{u[0, t_m]} \left[\int_0^{t_m} c(x_t, u_t)dt + \underbrace{\min_{u[t_m, \infty)} \int_{t_m}^\infty c(x_t, u_t)dt}_{J^*(x_{t_m})} \right] \end{aligned}$$

Same principle of optimality: What happens for $t_m \downarrow 0$. We assume $J^* \in C^1$ and write $\Delta x = x_{t_m} - x_0 = x_m - x$. We now use Taylor on the above expression

$$\begin{aligned} J^*(x) &= \min_{u[0, t_m]} \{c(x_t, u_t)t_m + J^*(x) + \nabla J^*(x) \cdot \Delta x + o(t_m)\} \\ \implies 0 &= \min_{u[0, t_m]} \left\{ \underbrace{c(x_t, u_t)}_{\rightarrow 0} \underbrace{\frac{t_m}{t_m}}_{\rightarrow 0} + \nabla J^*(x) \underbrace{\frac{\Delta x}{t_m}}_{\substack{\frac{d}{dt}|_{t=0} \\ = f(x_0, u_0)}} \right\} + \underbrace{o(1)}_{\rightarrow 0} \\ \implies 0 &= \min_u [c(x, u) + \nabla J^*(x) \cdot f(x_0, u_0)] \end{aligned}$$

this is a strong assumption! In principle we would need to talk about viscosity solutions ... Even weak solutions are not enough

Theorem 1.23. *If the value function J^* has continuous derivatives, then it satisfies the Hamilton-Jacobi-Bellmann equation*

$$0 = \min_u [c(x, u) + \nabla J^*(x) \cdot f(x_0, u_0)] \quad (13)$$

The term to minimize has an interpretation as an Hamiltonian

$$H(x, p, u) = c(x, u) + p^\top f(x, u).$$

One can show

Theorem 1.24. *Suppose that an optimal state-action pair exists and that $J^* \in C^1$. Then u_t^* must minimize for each t*

$$\min_u H(x_t^*, p_t^*, u) = H(x_t^*, p_t^*, u_t^*)$$

with $p_t^* = \nabla_x J^*(x_t^*)$.

Remark. *Relaxing away from ∇J^* or ∇J can have theoretical and computational advantages.*

1.13 Linear quadratic regulator revisited (once more)

$$\begin{aligned} \frac{d}{dt}x &= Fx + Gu, \quad x(0) = x_0 \\ c(x, u) &= x^\top Sx + u^\top Ru \end{aligned}$$

everything we observed so far carries over, assuming J^* is finite, we have

$$J^*(x) = x^\top M^* x$$

the HSB (13) gives

$$\begin{aligned} \phi^*(x) &= \operatorname{argmin}_u \{x^\top Sx + u^\top Ru + [2M^*x]^\top [Fx + Gu]\} \\ &= \operatorname{argmin}_u \{u^\top Ru + 2x^\top M^*Gu\} \end{aligned}$$

So,

$$0 = \nabla_u \{u^\top Ru + 2x^\top M^*Gu\}|_{u=\phi^*(x)}$$

and we get

$$\phi^*(x) = -R^{-1}G^\top M^*x$$

and

$$\frac{d}{dt}x^* = [F - GR^{-1}G^\top M^*]x^*.$$

HSB (13) further gives

$$\begin{aligned} 0 &= \{x^\top Sx + u^\top Ru + [2M^*x]^\top [Fx + Gu]\}|_{u=\phi^*(x)} \\ &= x^\top \{S + M^*GR^{-1}G^\top M^*\}x + x^\top \{2M^*F + 2M^*GR^{-1}G^\top M^*\}x \end{aligned}$$

using $2x^\top M^* F x = x^\top [M^* F + F^\top M^*]$ we get

$$\begin{aligned} &= x^\top \{S + M^* F + F^\top M^* - M^* G R^{-1} G^\top M^*\} x \\ &\quad \{S + M^* F + F^\top M^* - M^* G R^{-1} G^\top M^*\} \end{aligned}$$

holds for any x and is symmetric, so it follows M^* is a positive definite solution to the algebraic Riccati equation

$$0 = S + M^* F + F^\top M^* - M^* G R^{-1} G^\top M^*$$

Chapter 2:

ODE methods for algorithm design

2.1 ODE methods for algorithm design

Start of lecture 07
(06.05.2025)

Four steps:

- Formulate the algorithmic goal as the root finding problem

$$\bar{f}(\theta^*) = 0$$

- if necessary, refine the design of \bar{f} to ensure that the associated ODE is **globally asymptotically stable**

$$\frac{d}{dt}\vartheta = \bar{f}(\vartheta)$$

θ for discrete settings, ϑ for continuous settings.
Both do the same job

- Is an **Euler-approximation** appropriate?

$$\theta_{n+1} = \theta_n + \alpha_{n+1} \bar{f}(\theta_n) \quad (1)$$

θ_{n+1} is the next iterate, not the next time step!

- Design an algorithm to approximate (1) based on whatever observation is available.

Remark. The idea is to transfer the global stability from the ODE to the algorithm.

Goal: Construct a vector field f such that ϑ_t converges to the **target** $\theta^* \in \mathbb{R}^d$, where θ^* is an equilibrium

$$f(\theta^*) = 0.$$

In ODE theory one uses so called **Picard-Iteration**

$$\vartheta_t^{n+1} = \theta_0 + \int_0^t f(\vartheta_\tau^n) d\tau, \quad 0 \leq t \leq T \quad (2)$$

based on

$$\vartheta_0 + \int_0^t f(\vartheta_\tau) d\tau, \quad 0 \leq t \leq T. \quad (3)$$

Proposition 2.1. Suppose that the function f is globally Lipschitz continuous:

$$\exists L > 0 : \forall x, y \in \mathbb{R}^d : \|f(x) - f(y)\| \leq L\|x - y\|$$

Then for each θ_0 there exists a unique solution to (3). in the finite time horizon. Moreover, successive approximation is uniformly convergent:

$$\lim_{n \rightarrow \infty} \max_{0 \leq t \leq T} \|\vartheta_t^n - \vartheta_t\| = 0$$

Proposition 2.2 (Grönwall-Bellman-inequality). Let α, β and z be non-negative functions defined

on $[0, T]$, $T > 0$. Assume that β, z are continuous and that

$$z_t \leq \alpha_t + \int_0^t \beta_s z_s ds, \quad 0 \leq t \leq T$$

Then it holds

$$(i) \quad z_t \leq \alpha_t + \int_0^t \alpha_s \beta_s \exp\left(\int_s^t B_r dr\right) ds$$

(ii) if in addition the function α is non-decreasing, then

$$z_t \leq \alpha_t \exp\left(\int_0^t B_s ds\right), \quad 0 \leq t \leq T$$

Proof. Both proofs can be found in any textbook on ODEs. The second one is also found in [2]. \square

Proposition 2.3. Consider $\frac{d}{dt}\vartheta = f(\vartheta)$, $\vartheta_0 = \theta_0$ with f globally Lipschitz. Then

Not that nice, but at least a bound ...

(i) There is a constant B_f depending only on f such that, with $t \geq 0$

$$\|\vartheta_t\| \leq (B_f + \|\vartheta_0\|) e^{Lt} - B_f \quad (4)$$

$$\|\vartheta_t - \vartheta_0\| \leq \|B_f + L\|\vartheta_0\| t e^{Lt} \quad (5)$$

(ii) If there is an equilibrium θ^* , then for each initial condition:

$$\|\vartheta_t - \theta^*\| \leq \|\vartheta_0 - \theta^*\| e^{Lt} \quad (6)$$

Proof. (ii): use 3 to get

$$\vartheta_t - \theta^* = \vartheta_0 - \theta^* + \int_0^t f(\vartheta_\tau) d\tau$$

Since $f(\theta^*) = 0$, we see

$$\begin{aligned} \|f(\vartheta_\tau)\| &= \|f(\vartheta_\tau) - f(\theta^*)\| \\ &\leq L \underbrace{\|\vartheta_\tau - \theta^*\|}_{=: z_\tau} \end{aligned}$$

So

$$z_t \leq z_0 + L \int_0^t z_\tau d\tau.$$

Using proposition 2.2 (ii) with $\beta_t \equiv L$, $\alpha_t \equiv z_0$ we get

$$\|\vartheta_t - \theta^*\| \leq \|\vartheta_0 - \theta_0\| \exp(Lt)$$

(i): take any $\bar{\theta} \in \mathbb{R}^d$ and use the Lipschitz continuity

$$\begin{aligned} \|f(\theta)\| &\leq \|f(\theta) - f(\bar{\theta})\| + \|f(\bar{\theta})\| \\ &\leq L\|\theta - \bar{\theta}\| + \|f(\bar{\theta})\| \\ &\leq L\|\theta\| + L\|\bar{\theta}\| + \|f(\bar{\theta})\|. \end{aligned}$$

For any fixed $\bar{\theta}$, define $B_f = \|\bar{\theta}\| + \|f(\bar{\theta})\|/L$ which gives

$$\|f(\theta)\| \leq L[\|\theta\| + B_f], \quad \theta \in \mathbb{R}^d$$

using (3)

$$\begin{aligned} \|\vartheta_t\| + B_f &\leq \|\vartheta_0\| + B_f + \underbrace{L}_{\beta} \int_0^t \left[\underbrace{\|\vartheta_\tau + B_f\|}_{z_\tau} \right] d\tau \\ &\leq [\|\vartheta_0\| + B_f] \exp(Lt) \end{aligned}$$

where the last step follows by the same trick as in (ii), i.e. by using Grönwall. \square

2.2 Euler's method once more

$$\frac{\hat{\vartheta}_{t_{n+1}} - \hat{\vartheta}_{t_n}}{\alpha_{n+1}} = f(\hat{\vartheta}_{t_n}), \quad \hat{\vartheta}_0 = \vartheta_0 = \theta_0 \quad (7) \quad \begin{array}{l} \text{Explicit Euler, implicit} \\ \text{Euler is nicer to analyze} \end{array}$$

or

$$\hat{\vartheta}_{t_{n+1}} = \hat{\vartheta}_{t_n} + \alpha_{n+1} f(\hat{\vartheta}_{t_n})$$

It can be shown for f globally Lipschitz

$$\max_{0 \leq t \leq T} \|\hat{\vartheta}_t - \vartheta_t\| \leq \underbrace{K(L, T)}_{\text{exponential in } L, T} \max\{\alpha_k \mid t_k < T\} \quad (8)$$

2.3 Optimization

Goal: Find, for some loss function $\Gamma : \mathbb{R}^d \rightarrow \mathbb{R}_+$,

$$\theta^* \in \operatorname{argmin} \Gamma(\theta). \quad (9)$$

Use steepest-descent, formulated as ODE

$$\frac{d}{dt} \vartheta = -\nabla_{\theta} \Gamma(\theta) \quad (10)$$

so called gradient flow.

$$\nabla \Gamma(\theta_0) \perp \{\theta \in \mathbb{R}^d \mid \Gamma(\theta) = \Gamma(\theta_0)\} =: S_{\Gamma}(\theta_0)$$

The gradient flow steers into the interior of $S_{\Gamma}(\theta_0)$.

Definition 2.4. (i) A set $S \subset \mathbb{R}^d$ is convex if it contains all line segments with endpoints in S

(ii) A function $\Gamma : S \rightarrow \mathbb{R}$ with S convex, is called convex if for any $\theta^0, \theta^1 \in S$ and $\rho \in (0, 1)$

$$\Gamma((1 - \rho)\theta^0 + \rho\theta^1) \leq (1 - \rho)\Gamma(\theta^0) + \rho\Gamma(\theta^1)$$

Γ is strictly convex if this inequality is strict whenever $\theta^0 \neq \theta^1$

(iii) If Γ is differentiable, then it is called strongly convex if for $\delta_0 > 0$

$$\langle \nabla \Gamma(\theta) - \nabla \Gamma(\theta^0), \theta - \theta^0 \rangle \geq \delta_0 \|\theta - \theta_0\|^2, \quad \forall \theta, \theta^0 \in S$$

From numerical optimization we know:

Theorem 2.5. Suppose that $\Gamma : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex. Then for given $\theta^0 \in \mathbb{R}^d$

(i) if θ^0 is a local minima, then it is also a global minimum

(ii) if Γ is differentiable at θ^0 , with $\nabla \Gamma(\theta) = 0$, then θ^0 is a global minimum

(iii) if either (i) or (ii) hold, and if Γ is strictly convex, then θ^0 is the unique global minimum

Proposition 2.6. Suppose that Γ is continuously differentiable, convex and coercive, with unique minimizer θ^* . Then the gradient flow

$$\frac{d}{dt} \vartheta = -\nabla \Gamma(\vartheta)$$

is globally asymptotically stable, with unique equilibrium θ^* .

If Γ is strongly convex, then the rate of convergence is exponential

$$\|\vartheta_t - \theta^*\| \leq e^{-\delta_0 t} \|\vartheta_0 - \theta^*\|,$$

where δ_0 comes from theorem 2.5.

Proof. We use as Lyapunov function $V(\theta) = \frac{1}{2}\|\theta - \theta^*\|^2$. From the chain rule

$$\frac{d}{dt}V(\vartheta_t) = -\nabla_{\theta}\Gamma(\vartheta_t)^{\top} [\vartheta_t - \theta^*]$$

By convexity we get the following bound

$$\Gamma(\theta^*) \geq \Gamma(\vartheta_t) + \nabla_{\theta}\Gamma(\vartheta_t)^{\top} [\theta^* - \vartheta_t]$$

using the support condition this becomes

$$\frac{d}{dt}V(\vartheta_t) \leq \Gamma(\theta^*) - \Gamma(\vartheta_t) \leq 0$$

since θ^* is the minimum. The strict inequality (< 0) holds when $\vartheta_t \neq \theta^*$. V fulfills definition 1.20 and proposition 1.21 gives global asymptotic stability.

Under strict convexity

Coercive, therefore
inf-compact

$$\begin{aligned} \frac{d}{dt}V(\vartheta_t) &= - \left[\nabla_{\theta}\Gamma(\vartheta_t) - \underbrace{\nabla_{\theta}\Gamma(\theta^*)}_{=0} \right]^{\top} [\vartheta_t - \theta^*] \\ &\stackrel{\text{strong convexity}}{\leq} -\delta_0 \|\vartheta_t - \theta^*\|^2 = -2\delta_0 V(\vartheta_t) \end{aligned}$$

This implies $V(\vartheta_t) \leq V(\vartheta_0) \exp(-2\delta_0 t) \forall t$ by integrating. \square

Theorem 2.7. *If the Polyak-Lojasiewicz (PL) inequality*

$$\frac{1}{2}\|\nabla\Gamma(\theta)\|^2 \geq \mu|\Gamma(\theta) - \Gamma(\theta^*)| \quad (11)$$

holds then the gradient flow satisfies for each initial ϑ_0

$$\Gamma(\vartheta_t) - \Gamma^* \leq e^{-\mu t}(\Gamma(\vartheta_0) - \Gamma^*).$$

If in addition Γ is coercive, then the solutions are bounded and any limit point θ_{∞} of $\{\vartheta_t\}$ is an optimizer

$$\Gamma(\theta_{\infty}) = \Gamma^*$$

*Used in stochastic
gradient descent*

Proof. We use $V(\theta) = \frac{1}{2}|\Gamma(\theta) - \Gamma^*|$ for the Lyapunov function.

$$\begin{aligned} \implies \frac{d}{dt}V(\vartheta_t) &= \frac{1}{2}\nabla_{\theta}\Gamma(\vartheta_t)^{\top} \frac{d}{dt}\vartheta_t \\ &= -\frac{1}{2}\|\nabla\Gamma(\vartheta_t)\|^2 \leq -\mu V(\vartheta_t) \end{aligned}$$

This implies using the same technique as in the previous proof

$$\begin{aligned} \frac{1}{2}[\Gamma(\vartheta_t) - \Gamma^*] &= V(\vartheta_t) \leq e^{-\mu t}V(\vartheta_0) \\ &= e^{-\mu t} \frac{1}{2}[\Gamma(\vartheta_0) - \Gamma^*] \end{aligned}$$

If Γ is coercive, then trajectories of ϑ evolve in the compact set $S = \{\theta \mid V(\theta) \leq V(\vartheta_0)\}$. If θ_{∞} is a limit point $\theta_{\infty} = \lim_{n \rightarrow \infty} \vartheta_{t_n}$ for $t_n \rightarrow \infty$. Using the continuity of the loss function, this implies optimality:

$$\Gamma(\theta_{\infty}) = \lim_{n \rightarrow \infty} \Gamma(\vartheta_{t_n}) = \Gamma^* \quad \square$$

Consider the Euler method for the gradient flow:

$$\theta_{k+1} = \theta_k - \alpha \nabla\Gamma(\theta_k) \quad (12)$$

Theorem 2.8. Suppose that Γ satisfies

(i) the L -smooth inequality (LSI)

$$\Gamma(\theta') \leq \Gamma(\theta) + [\theta' - \theta]^\top \nabla \Gamma(\theta) + \frac{1}{2}L\|\theta' - \theta\|^2$$

(ii) the PL inequality 11

Then it holds for $\alpha \leq \frac{1}{2}$

$$\Gamma(\theta_k) - \Gamma^* \leq (1 - \alpha\mu)^k [\Gamma(\theta_0) - \Gamma^*].$$

Proof.

$$\begin{aligned} \Gamma(\theta_{k+1}) - \Gamma(\theta_k) &\stackrel{\text{LSI}}{\leq} [\theta_{k+1} - \theta_k]^\top \nabla \Gamma(\theta_k) + \frac{1}{2}L\|\theta_{k+1} - \theta_k\|^2 \\ &\stackrel{12}{=} -\alpha\|\nabla \Gamma(\theta_k)\|^2 + \frac{1}{2}L\alpha^2\|\nabla \Gamma(\theta_k)\|^2 \\ &= (-\alpha + \frac{1}{2}L\alpha^2)\|\nabla \Gamma(\theta_k)\|^2 \end{aligned}$$

If $\alpha \leq \frac{1}{L}$ then $(-\alpha + \frac{1}{2}L\alpha^2) \leq \frac{1}{2}\alpha$

$$\begin{aligned} &\leq -\frac{1}{2}\alpha\|\nabla \Gamma(\theta_k)\|^2 \\ &\stackrel{\text{LSI}}{\leq} -\alpha\mu|\Gamma(\theta_k) - \Gamma^*| \end{aligned}$$

and therefore

$$\Gamma(\theta_{k+1}) - \Gamma^* \leq (1 - \alpha\mu)(\Gamma(\theta_k) - \Gamma^*)$$

after iterating $k - 1$ times we obtain the result. \square

Lemma 2.9. Suppose that $\nabla \Gamma$ is globally Lipschitz

$$\|\nabla \Gamma(\theta') - \nabla \Gamma(\theta)\| \leq L\|\theta' - \theta\|, \quad \forall \theta, \theta' \in S$$

Then

$$(i) \quad |\langle \nabla \Gamma(\theta') - \nabla \Gamma(\theta), \theta' - \theta \rangle| \leq L\|\theta' - \theta\|^2$$

(ii) if S is convex, then Γ is L -smooth

Proof. (i)

$$\begin{aligned} |\langle \nabla \Gamma(\theta') - \nabla \Gamma(\theta), \theta' - \theta \rangle| &\leq \|\nabla \Gamma(\theta') - \nabla \Gamma(\theta)\| \|\theta' - \theta\| \\ &\leq L\|\theta' - \theta\|^2 \end{aligned}$$

(ii) for $\theta', \theta \in S$ denote $S \ni \theta^t := \theta + t(\theta' - \theta)$ and $\xi^t = \Gamma(\theta^t)$.

θ^t in S , since S is convex

$$\begin{aligned} \frac{d}{dt}\xi^t &= \langle \nabla \Gamma(\theta^t), \theta' - \theta \rangle \\ \frac{d}{dt}\xi^t - \frac{d}{dt}\xi^0 &= \langle \nabla \Gamma(\theta^t) - \nabla \Gamma(\theta^0), \theta' - \theta \rangle \\ &\stackrel{(i)}{\leq} tL\|\theta' - \theta\|^2 \end{aligned}$$

Now integrate

$$\begin{aligned} \Gamma(\theta') &= \xi^1 = \xi^0 + \int_0^1 \frac{d}{dt}\xi^t dt \\ &\leq \xi^0 + \frac{d}{dt}\xi^0 + \frac{1}{2}L\|\theta' - \theta\|^2 \\ &= \Gamma(\theta) + \langle \nabla \Gamma(\theta), \theta' - \theta \rangle + \frac{1}{2}L\|\theta' - \theta\|^2 \end{aligned}$$

\square

These are more general version of global Lipschitz and convexity

Remark. *Strong convexity:*

$$\langle \nabla \Gamma(\theta') - \nabla \Gamma(\theta), \theta' - \theta \rangle \geq \delta_0 \|\theta' - \theta\|^2$$

With $D_\Gamma(y \mid x) = \Gamma(y) - \Gamma(x) + \langle \nabla \Gamma(x), y - x \rangle$ is the Bregman divergence.

$$\frac{\mu}{2} \|\theta' - \theta\|^2 \leq D_\Gamma(\theta' \mid \theta) \leq \frac{L}{2} \|\theta' - \theta\|^2$$

This gives a bound on the loss function from both sides ...

2.4 Quasi stochastic approximation

Assume there are observations $\Phi_n \subset \Omega$, which we might consider as realizations of a random variable Φ . We have

$$f : \mathbb{R}^d \times \Omega \rightarrow \mathbb{R}^d$$

$$\bar{f}(\theta) := \mathbb{E}(\underbrace{f(\theta, \Phi)}_{\text{what we observe}}), \theta \in \mathbb{R}^d$$

As before we look for $\bar{f}(\theta^*) = 0$

$$\frac{d}{dt} \vartheta_t = \bar{f}(\vartheta_t)$$

A key assumption is that what happens when following the state dynamics in any step depends only on the current state.

I.e. we have the Markov property

$$\Phi_n = [\cos(\omega n), \sin(\omega n)], \omega > 0$$

Markov chain on the unit circle. We will talk about the probing signal ξ and consider

the book uses Θ instead of $\hat{\theta}$

$$\frac{d}{dt} \hat{\theta}_t = a_t f(\hat{\theta}_t, \xi_t) \quad (13)$$

a quasistochastic approximation(QSA)-ODE, a_t is the step size.

For deterministic probing signals, we mainly consider two examples

Mixture of sin functions

$$\xi_t = \sum_{i=1}^K \overbrace{V^i}^{\in \mathbb{R}^m} \sin(2\pi[\Phi_i + \omega_i t])$$

Mixture of periodic functions, fixed K , phase $\{\Phi_i\}$, frequencies $\{\omega_i\}$.

$$\xi_t = \sum_{i=1}^K V^i [\Phi_i + \omega_i t]_{\text{modulo } 1}$$

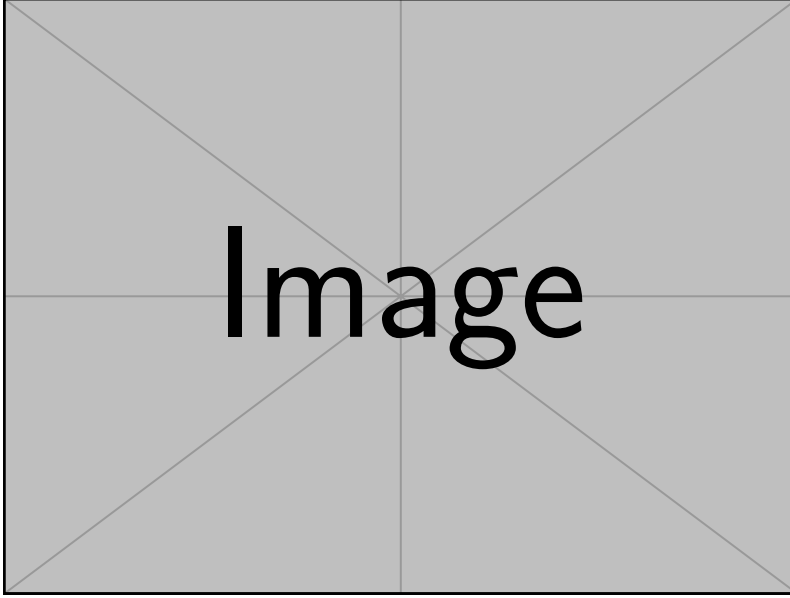


Figure 2.1: Sketch 2.01

These signals have well defined steady-state means and covariance matrices.

Special case: $\xi_t(i) = \sqrt{2} \sin(\omega_i t)$, $1 \leq i \leq m$, $\omega_i \neq \omega_j \forall i \neq j$. Then the steady-state mean

$$\lim_{T \rightarrow \infty} \int_0^T \xi_t dt = 0$$

and covariance

$$\lim_{T \rightarrow \infty} \int_0^T \xi_i \xi_i^\top dt = \text{Id}$$

We now use a slightly different notation $\hat{\theta}$ becomes $\tilde{\theta}$.

Start of lecture 09
(13.05.2025)

$$\frac{d}{dt} \tilde{\theta}_t = a_t f(\tilde{\theta}, \xi_t) \quad (14)$$

a_t non-negative.

Now consider integrating $y : [0, 1] \rightarrow \mathbb{R}$. Basic Monte-Carlo

$$\theta_n = \frac{1}{n} \sum_{i=0}^{n-1} y(\underbrace{\Phi(k)}_{\sim \text{Unif}([0,1])}) \quad (15)$$

A QSA approach is to use the saw tooth function

$$\xi_t = t(\text{modulo } 1).$$

Obtain estimate by

$$\tilde{\theta} = \frac{1}{t} \int_0^t y(\xi_t) dr \quad (16)$$

with a reasonable discretization afterwards.

To use (QSA-ODE (14)) $f(\theta, \xi) = y(\xi) - \theta$ with mean vector field

$$\begin{aligned} \bar{f}(\theta) &= \lim_{T \rightarrow \infty} \int_0^T f(\theta, \xi_t) dt \\ &= \int_0^1 y(\xi_t) dt - \theta \end{aligned}$$

which gives $\theta^* = \int_0^1 y(\xi_t) dt$ as the unique root of \bar{f} . The QSA-ODE 14 is

$$\frac{d}{dt} \tilde{\theta}_t = a_t [y(\xi_t) - \tilde{\theta}_t]$$

(16) can be transformed into

$$\frac{d}{dt} \tilde{\theta}_t = \left[-\frac{1}{t^2} \int_0^t y(\xi_r) dr + \frac{1}{t} y(\xi_t) \right] = \underbrace{\frac{1}{t}}_{\equiv a_t} [y(\xi_t) - \theta_t] \quad (17)$$

Example. $y(\theta) = e^4(\sin(100\theta))$, mean $\theta^* \approx -0.5 \approx -0.48$. Choose $a_t = \frac{g}{1+t}$

2.5 Approximate Policy Improvement

nonlinear state model in continuous time:

$$\frac{d}{dt} x_t = f(x_t, u_t), t \geq 0 \quad (18)$$

$$J^*(x) = \min_{\underline{u}} \int_0^\infty c(x_t, u_t) dt, x = x_0 \quad (19)$$

Given feedback law $u_t = \phi(x_t)$, we have

$$J^\phi(x) = \int_0^\infty c(x_t, \phi(x_t)) dt, x = x_0 \quad (20)$$

Proposition 2.10. *If J is finite, then for each initial condition x_0 and each t*

$$\frac{d}{dt} J(x_t) = -c(x_t)$$

If J is continuously differentiable, then the Lyapunov bound $\frac{d}{dt} V(x_t)$ from definition 1.20 follows with equality

$$\nabla J(x) f(x) = -c(x)$$

Proof. For any $T > 0$, $J(x_0) = \int_0^T c(x_r) dr + J(x_T)$. For $t \geq 0, \delta > 0$ given, use $T = t + \delta$ and $T = t$ and subtract:

$$\begin{aligned} 0 &= J(x_0) - J(x_0) = \int_t^{t+\delta} c(x_r) dr + (J(x_{t+\delta}) - J(x_t)) \\ &= \underbrace{\frac{1}{\delta} \int_t^{t+\delta} c(x_r) dr}_{\xrightarrow{\delta \rightarrow 0} c(x_t)} + \underbrace{\frac{1}{\delta} (J(x_{t+\delta}) - J(x_t))}_{\xrightarrow{\delta \rightarrow 0} \frac{d}{dt} J(x_t)} \\ &\implies \frac{d}{dt} J(x_t) = -c(x_t) \end{aligned}$$

Using the chain rule yields the second equation. □

For J^ϕ we have

$$0 = c(x, \phi(x)) + \nabla J^\phi(x) \cdot f(x, \phi(x))$$

Policy Improvement in continuous time:

$$\phi^+(x) \in \operatorname{argmin}_u \left\{ \underbrace{c(x, u) + \nabla J(x) \cdot f(x, u)}_{\text{need to approximate by } Q^\phi(x, u)} \right\}$$

Now aim for updating of Q -function. Add to the above J^ϕ on both sides

$$J^\phi(x) = J^\phi(x) + c(x, \phi(x)) + \nabla J^\phi(x) \cdot f(x, \phi(x))$$

We solved for the optimal Q -function by using a fixed point equation, with $\underline{Q}^\phi(x) = Q^\phi(x, \phi(x))$ we write

$$Q^\phi(x, u) = \underline{Q}^\phi(x) + c(x, u) + \nabla \underline{Q}^\phi(x) f(x, u).$$

\underline{Q} for the fixed, but optimal choice of u

Consider $\{Q^\theta \mid \theta \in \mathbb{R}^d\}$ family of approximations. Bellman errors (Temporal differences expressions?) gives

$$B^\theta(x_t, u_t) = -Q^\theta(x_t, u_t) + \underline{Q}^\theta(x) + c(x_t, u_t) + \underbrace{\nabla \underline{Q}^\theta(x) f(x_t, u_t)}_{= \frac{d}{dt} Q^\theta(x_t)} \quad (21)$$

Everything on the RHS is can be observed for any state-action pair without knowledge of f . Now, find θ^* that minimizes

$$\|B^\theta\|^2 = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T [B^\theta(x_t, u_t)]^2 dt$$

Choose feedback law with exploration $u_t = \tilde{\phi}(x_t, \xi_t)$. Assuming bounded state trajectories, such that (21) exists, define $\Gamma(\theta) = \frac{1}{2} \|B^\theta\|^2$, we get

$$0 \stackrel{!}{=} \nabla \Gamma(\theta) = \lim_{t \rightarrow \infty} \int_0^T [B^\theta(x_t, u_t)] \nabla_\theta B^\theta(x_t, u_t) dt$$

Gradient flow

$$\frac{d}{dt} \vartheta_t = -\nabla_\theta \Gamma(\vartheta_t)$$

QSA counterpart is (21) with probing signal

$$\frac{d}{dt} \tilde{\theta}_t = -a_t B^{\tilde{\theta}_t}(x_t, u_t) \kappa_t^{\tilde{\theta}_t}$$

with

$$\begin{aligned} \kappa_t^{\tilde{\theta}_t} &= \nabla_\theta B^{\tilde{\theta}_t}(x_t, u_t) \\ &= -\nabla_\theta Q^\theta(x_t, u_t) + \{\nabla_\theta Q^\theta(x_t, \phi(x_t)) + \frac{d}{dt} \nabla_\theta Q^\theta(x_t, \phi(x_t))\} \end{aligned}$$

assuming we can exchange differentiation w.r.t time and w.r.t θ . (QSA-ODE)

$$\frac{d}{dt} \tilde{\theta}_t = a_t f(\tilde{\theta}_t, \xi_t)$$

aim to relate this to

$$\frac{d}{dt} \vartheta_t = \bar{f}(\vartheta_t).$$

Lemma 2.11. Define the change of variables

$$\tau = s_t := \int_0^t a_r dr, \quad t \geq t_0.$$

Let $\{\vartheta_\tau \mid \tau \geq \tau_0\}$ the solution to the ODE above initialized to $\tau_0 = s_{t_0}$ with $\vartheta_{\tau_0} = \tilde{\theta}_{t_0}$. The solution to

$$\frac{d}{dt} \bar{\theta}_t = a_t \bar{f}(\bar{\theta}_t), \quad t \geq t_0, \quad \bar{\theta}_{t_0} = \tilde{\theta}_{t_0}$$

is given by $\bar{\theta}_t = \vartheta_\tau$.

Proof. Change of variables and observing that

$$d\tau = a_t dt.$$

□

Recall $\bar{f}(\theta) := \lim_{T \rightarrow \infty} \int_0^T f(\theta, \xi_t) dt$ for all $\theta \in \mathbb{R}^d$. Remember the temporal transformation

Start of lecture 10
(15.05.2025)

$$\tau = s_t = \int_0^t a_r dr$$

and lemma 2.11. Define $\hat{\theta}_\tau = \tilde{\theta}(s^{-1}(\tau)) = \tilde{\theta}_t|_{t=s^{-1}(\tau)}$. By the chain rule and observing that $d\tau = a_t dt$ yields

$$\frac{d}{d\tau} \hat{\theta}_\tau = \frac{d}{d\tau} \tilde{\theta}(s^{-1}(\tau)) = f(\tilde{\theta}(s^{-1}(\tau)), \xi(s^{-1}(\tau))).$$

$\hat{\theta}, \tilde{\theta}$ differ only by a time scaling, so convergence of the one yields convergence of the other.

Lemma 2.12. Consider the original ODE

$$\frac{d}{dt} \vartheta_t = \bar{f}(\vartheta_t) \quad (22)$$

and assume f is locally Lipschitz with constant L_f . Then there exists a constant B_f depending only on f , such that

Version of proposition 2.3

$$\|\hat{\theta}_t - \hat{\theta}_0\| \leq (B_f + L_f \|\hat{\theta}_0\|) t e^{L_f t}, \quad t \geq 0$$

Proof. Proof of proposition 2.3 in adapted notation. \square

Now, denote by $\vartheta_w^\tau, w \geq \tau$ the unique solution to (22):

$$\frac{\partial}{\partial w} \vartheta_w^\tau = \bar{f}(\vartheta_w^\tau), \quad w \geq \tau, \quad \vartheta_\tau^\tau = \hat{\theta}_\tau$$

with that we get

quasistochastic vs
continuous

1. $\vartheta_{\tau+v}^\tau = \hat{\theta}_\tau + \int_0^{\tau+v} \bar{f}(\vartheta_w^\tau) dw, \quad \tau, v \geq 0$
2. $\hat{\theta}_{\tau+v} = \hat{\theta}_\tau + \int_\tau^{\tau+v} f(\hat{\theta}_w, \xi(s^{-1}(w))) dw, \quad \tau, v \geq 0$

The following assumptions will be used in the following:

QSA1 The process a is non-negative, monotonically decreasing and $\lim_{t \rightarrow \infty} a_t = 0, \int_0^\infty a_r dr = \infty$

it does not go to zero too fast

QSA2 The functions \hat{f}, f are Lipschitz continuous with constant L_f :

$$\begin{aligned} \|\bar{f}(\theta') - \bar{f}(\theta)\| &\leq \|L_f\| \|\theta' - \theta\| \\ \|f(\theta', z) - f(\theta, z)\| &\leq \|L_f\| \|\theta' - \theta\| \end{aligned}$$

for all $\theta, \theta' \in \mathbb{R}^d, z \in \Omega$ and there exists a Lipschitz continuous functions $b_0 : \mathbb{R}^d \rightarrow \mathbb{R}_+$, such that for all $\theta \in \mathbb{R}^d$

Is my probing covering everything: ergodicity, ergodic bound

$$\left\| \int_{t_0}^{t_1} f(\theta, \xi_t) - \bar{f}(\theta) dt \right\| \leq b_0(\theta), \quad 0 \leq t_1 \leq t_1$$

QSA3 The ODE $\frac{d}{dt} \vartheta_t = \bar{f}(\vartheta_t)$ has a globally asymptotically stable equilibrium θ^*

Consider first, arbitrary θ

Lemma 2.13. Assume (QSA1), (QSA2) hold for any fixed $T > 0$ and $\theta \in \mathbb{R}^d$.

There is a connect to the law of large numbers ...

$$\left\| \int_\tau^{\tau+T} [f(\theta, \xi(s^{-1}(w))) - \bar{f}(\theta)] dw \right\| \leq b_0(\theta) \epsilon_\tau^f,$$

where $\epsilon_\tau^f = 3a_t|_{t=s^{-1}(\tau)}$ and b_0 comes from (QSA2).

Proof. Set $\tilde{f}_w(\theta) = f(\theta, \xi_w) - \bar{f}(\theta)$ for each w, θ . Write

large ϵ_t in the book?
Prob. \mathcal{E}

$$E_t = \int_0^t \tilde{f}_w(\theta) dw.$$

By assumptions $\|E_t\| \leq b_0(\theta)$, $t \geq 0$.

$$\begin{aligned} \int_{t_0}^{t_1} a_t \tilde{f}_t(\theta) dt &\stackrel{\text{IbP}}{=} a_t E_t \Big|_{t_0}^{t_1} - \int_{t_0}^{t_1} |a'_t| E_t dt \\ \left\| \int_{t_0}^{t_1} a_t \tilde{f}_t(\theta) dt \right\| &\leq a_{t_0} \|E_{t_0}\| + a_{t_1} \|E_{t_0}\| + \int_{t_0}^{t_1} |a'_t| E_t dt \\ &\stackrel{a \text{ decreasing}}{\leq} 2a_{t_0} b_0(\theta) - b_0(\theta) \int_{t_0}^{t_1} a'_t dt \\ &\leq 3a_{t_0} b_0(\theta) \end{aligned}$$

Set $t_0 = s^{-1}(\tau)$, $t_1 = s^{-1}(\tau + T)$, $t = s^{-1}(w)$, giving $dw = a_t dt$

$$\begin{aligned} \left\| \int_{\tau}^{\tau+T} [f(\theta, \xi(s^{-1}(w))) - \bar{f}(\theta)] dw \right\| &= \left\| \int_{t_0}^{t_1} a_t \tilde{f}_t(\theta) dt \right\| \\ &\leq 3a_{t_0} b_0(\theta) = \epsilon_{\tau}^f b_0(\theta) \end{aligned}$$

□

Proposition 2.14. *Assuming that $\hat{\theta}$ is bounded. Then for any $T > 0$*

$$\lim_{\tau \rightarrow \infty} \sup_{v \in [0, T]} \left\| \overbrace{\int_{\tau}^{\tau+v} [f(\hat{\theta}_w, \xi(s^{-1}(w))) - \bar{f}(\hat{\theta}_w)] dw}^{E_{\tau+v}^{\tau}} \right\| = 0$$

and

$$\lim_{\tau \rightarrow \infty} \sup_{v \in [0, T]} \left\| \hat{\theta}_{\tau+v} - \vartheta_{\tau+v}^{\tau} \right\| = 0$$

Proof. We use piecewise constant approximation, as in Riemannian integration, and set for $\delta > 0$, $\tau_k = \tau + k\delta$, $k \geq 0$

$$E_{\tau+v}^{\tau} = \sum_{k=0}^{n_v-1} \int_{\tau_k}^{\tau_{k+1}} [f(\hat{\theta}_{\tau_k}, \xi(s^{-1}(w))) - \bar{f}(\hat{\theta}_{\tau_k})] dw + \epsilon_v^{\tau},$$

which holds due to (QSA1), Lipschitz condition, $n_v = \lfloor \frac{v}{\delta} \rfloor$. and

$$\|\epsilon_v^{\tau}\| \leq b_L v \delta$$

for some finite constant b_L . Assuming $\hat{\theta}$ is bounded, this bound is uniform in τ . For fixed $\hat{\theta}_{t_k}$ we can use lemma 2.13, so

$$\begin{aligned} \|E_{\tau+v}^{\tau}\| &\leq \sum_{k=0}^{n_v-1} \epsilon_{\tau_k}^f b_0(\hat{\theta}_{t_k}) + b_L v \delta \\ &\leq \epsilon_{\tau}^f \sum_{k=0}^{n_v-1} b_0(\hat{\theta}_{\tau_k}) + b_L v \delta \end{aligned}$$

Let $b < \infty$ denote a constant such that $b_0(\hat{\theta}_{\tau_k}) \leq b \forall \tau$, which we can do since $\hat{\theta}$ is bounded, b_0 Lipschitz.

$$\|E_{\tau+v}^{\tau}\| \leq b \frac{v}{\delta} \underbrace{\epsilon_{\tau}^f}_{\xrightarrow{\tau \rightarrow \infty} 0 \text{ by QSA1}} + b_L v \delta$$

For any $T > 0$

$$\lim_{\tau \rightarrow \infty} \sup_{v \in [0, T]} \|E_{\tau+v}^\tau\| \leq 0 + b_L T \delta$$

Since $\delta > 0$ was arbitrary, we have the first statement.

For the second limit: $E_r^\tau = \vartheta_r^\tau - \hat{\theta}_r$. The pair of identities after lemma 2.12 give using Lipschitz condition from (QSA2) we get

$$E_{\tau+v}^\tau = 0 + \int_{\tau}^{\tau+v} \bar{f}(\hat{\theta}_w) - f(\hat{\theta}_w, \xi(s^{-1}(w))) dw + \underbrace{\int_{\tau}^{\tau+v} [\bar{f}(\vartheta_v^\tau) - \bar{f}(\hat{\theta}_w)] dw}_{\|\dots\| \leq L_f \|E_w^\tau\|}$$

$$\|E_{\tau+v}^\tau\| \leq \delta^\tau + L_f \int_{\tau}^{\tau+v} \|E_w^\tau\| dw,$$

where

$$\delta^\tau := \sup_{\tau' \geq \tau} \max_{0 \leq v \leq T} \left\| \int_{\tau'}^{\tau'+v} [\bar{f}(\hat{\theta}_w) - f(\hat{\theta}_w, \xi(s^{-1}(w)))] dw \right\|$$

Grönwall's lemma gives

$$\|E_{\tau+v}^\tau\| \leq e^{L_f v} \delta^\tau \forall \tau, \quad 0 \leq v \leq 1$$

$\delta^\tau \rightarrow 0$ for $\tau \rightarrow \infty$ due to the first statement. □

Journal

- **Lecture 01:** Covering: Introduction, (linear, continuous) State space models, equilibrium, (Lyapunov, asymptotically) stable, region of attraction, globally asymptotically stable . Starting in ‘[Organization](#)’ on page 2 and ending in ‘[State Space Models in continuous Time](#)’ on page 7. Spanning 5 pages
- **Lecture 02:** Covering: Lyapunov function, inf-compactness and coerciveness, sublevel sets, Poisson’s inequality, comparison theorem, a few propositions connecting the value function, equilibria and Lyapunov functions . Starting in ‘[State Space Models in continuous Time](#)’ on page 7 and ending in ‘[State Space Models in continuous Time](#)’ on page 9. Spanning 2 pages
- **Lecture 03:** Covering: discrete time Lyapunov equation, optimal control policy, controllability, linear quadratic regulator, Bellmann equation, principle of optimality, Q-function and some concepts from Reinforcement Learning . Starting in ‘[State Space Models in continuous Time](#)’ on page 9 and ending in ‘[Some concepts from Reinforcement Learning](#)’ on page 12. Spanning 3 pages
- **Lecture 04:** Covering: Value iteration, policy iteration, exploration-exploitation . Starting in ‘[Some concepts from Reinforcement Learning](#)’ on page 12 and ending in ‘[Exploration](#)’ on page 16. Spanning 4 pages
- **Lecture 05:** Covering: Approximate Q-functions, Bandits, discounted cost, shortest path, finite horizon and translations between them . Starting in ‘[Exploration](#)’ on page 16 and ending in ‘[Other control formulations](#)’ on page 19. Spanning 3 pages
- **Lecture 06:** Covering: Model predictive control, continuous time formulations of previous results . Starting in ‘[Other control formulations](#)’ on page 19 and ending in ‘[Linear quadratic regulator revisited \(once more\)](#)’ on page 22. Spanning 3 pages
- **Lecture 07:** Covering: Picard-Iteration, Grönwall-Bellma inequality, Euler’s method, gradient flows . Starting in ‘[ODE methods for algorithm design](#)’ on page 23 and ending in ‘[Optimization](#)’ on page 25. Spanning 2 pages
- **Lecture 08:** Covering: Polyak-Lojasiewicz inequality, L-smooth inequality, Bregman divergence, quasi stochastic approximation . Starting in ‘[Optimization](#)’ on page 25 and ending in ‘[Qausi stochastic approximation](#)’ on page 29. Spanning 4 pages
- **Lecture 09:** Covering: . Starting in ‘[Qausi stochastic approximation](#)’ on page 29 and ending in ‘[Approximate Policy Improvement](#)’ on page 32. Spanning 3 pages

- Lecture 10: Covering:

Starting in ‘Approximate Policy Improvement’ on page 32 and ending in ‘Approximate Policy Improvement’ on page 34. Spanning 2 pages

Bibliography

- [1] Tamer Basar, Sean Meyn, and William R. Perkins. *Lecture Notes on Control System Theory and Design*. 2024. arXiv: [2007.01367](https://arxiv.org/abs/2007.01367) [math.OC]. URL: <https://arxiv.org/abs/2007.01367>.
- [2] Sean Meyn. *Control Systems and Reinforcement Learning*. Cambridge University Press, 2022.