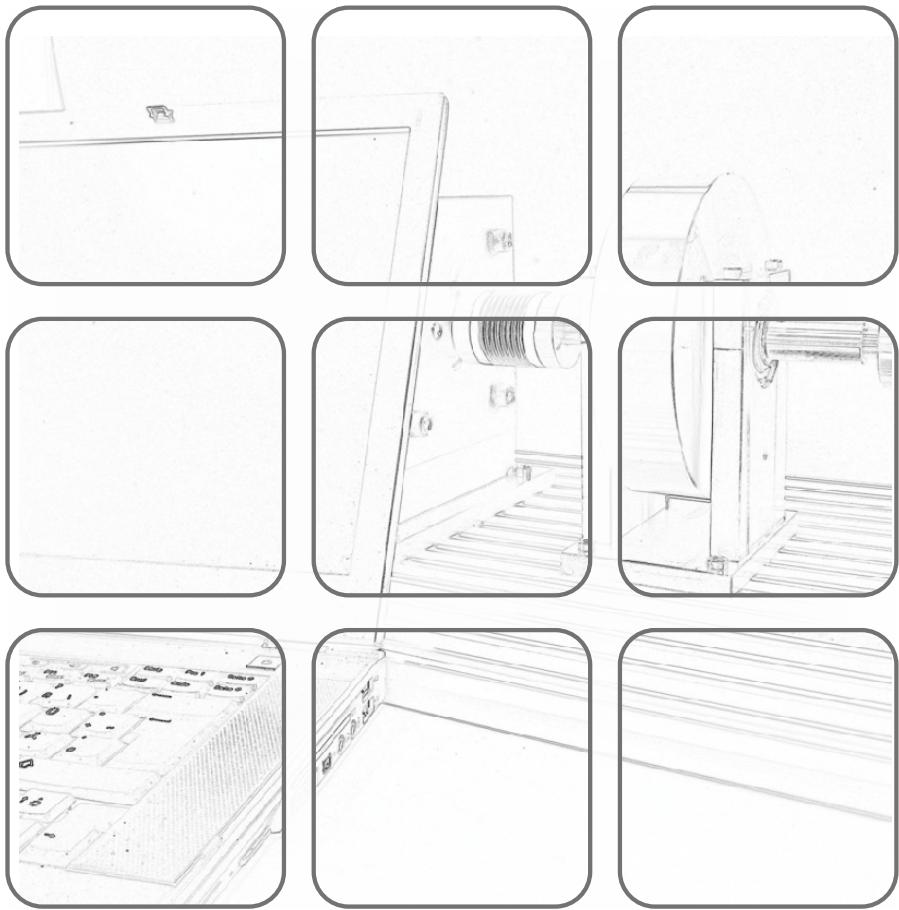


FORTGESCHRITTENE METHODEN DER NICHTLINEAREN REGELUNG

Vorlesung und Übung
Wintersemester 2023/2024

Andreas Deutschmann-Olek, Tobias Glück, Andreas Kugi



Fortgeschrittene Methoden der nichtlinearen Regelung

Vorlesung und Übung
Wintersemester 2023/2024

Andreas Deutschmann-Olek, Tobias Glück, Andreas Kugi

TU Wien
Institut für Automatisierungs- und Regelungstechnik
Gruppe für komplexe dynamische Systeme

Gußhausstraße 27–29
1040 Wien
Telefon: +43 1 58801 – 37615
Internet: <https://www.acin.tuwien.ac.at>

Contents

1 Pfadfolgeregelung	1
1.1 Definition von Pfaden	2
1.1.1 Parametrierte Darstellung von Pfaden	3
1.1.2 Implizite Definition von Pfaden	4
1.1.3 Vergleich und Umrechnung der Darstellungen	4
1.2 Begriffe und Ziele der Pfadfolgeregelung	7
1.3 Pfadfolgeregelung basierend auf exakter Linearisierung	8
1.3.1 Erweiterung der exakten Linearisierung für AI-Systeme	8
1.3.2 Implizit definierte Pfade	13
1.3.3 Parametrierte Pfade	22
1.4 Modellprädiktive Pfadfolgeregelung	31
1.5 Literatur	40
2 Dissipativität und Passivität	41
2.1 Glühsimulator	41
2.2 Einfaches Elektromagnetventil	43
2.3 Systemtheoretisches Konzept	44
2.3.1 Dissipativität	44
2.3.2 Passivität	45
2.3.3 Eigenschaften Passiver Systeme	47
2.3.4 Passivität und Lyapunov-Stabilität	49
2.4 Lineare passive Systeme	50
2.5 Positive Reellheit	53
2.6 Kanonische Form Passiver Systeme	56
2.6.1 Hamiltonsche Systeme	56
2.6.2 Port-Hamiltonsche Systeme	58
2.7 Passivitätsbasierter Reglerentwurf	60
3 Iterative learning control	67
3.1 Fixed-point iterations	69
3.2 Signals, systems, and the frequency domain	72
3.3 Frequency-domain ILC methods on infinite time horizons	79
3.3.1 Analysis of ILC laws	80
3.3.2 Deterministic design methods	82
3.3.3 Stochastically optimal learning laws	85
3.3.4 Q-filtering and robustness	88
3.3.5 Implementation aspects	91

3.4	Discrete-time systems on finite time horizons	96
3.4.1	Lifted system representation	97
3.4.2	ILC as an online optimization strategy	101
3.4.3	Norm-optimal ILC strategies	102
3.5	Literatur	104
4	Stochastic optimal control and reinforcement learning	105
4.1	Theoretical foundation	106
4.1.1	Stochastic dynamical systems	106
4.1.2	Rollouts, rewards and return	108
4.1.3	Different terminologies and interaction models	109
4.1.4	Markov Decision Process	109
4.1.5	Control policy	115
4.1.6	Value functions	116
4.1.7	Bellman Expectation Equation	118
4.1.8	Bellman Optimality Equation	120
4.2	Model-based reinforcement learning	122
4.2.1	Infinite horizon Dynamic Programming	122
4.2.2	Finite horizon Dynamic Programming	127
4.3	Model-free reinforcement learning	133
4.3.1	Generalized Policy Iteration	134
4.3.2	Approximate solution methods	139
4.3.3	Policy Gradients	140
4.3.4	Stochastic Policy Gradient	141
4.3.5	Actor-Critic Methods	144
4.3.6	Off-Policy Policy Gradient	144
4.3.7	Proximal Policy Optimization	146
4.4	Literatur	150
A	Basics of probability theory	152
A.1	Variables	152
A.2	Random variables	152
A.3	Probability distribution	153
A.4	Expectation	153
A.5	Probability of an event	153
A.6	Conditional expectation	154
A.7	Law of total expectation	154
A.8	Rules of probability	154
A.9	The chain rule of conditional probabilities	155
A.10	Bernoulli distribution	155
A.11	Normal distribution	155

1 Pfadfolgeregelung

Die Aufgabenstellungen in der Regelungstechnik lassen sich unter anderem in *Arbeitspunktstabilisierung*, *Trajektorienfolgeregelung* und *Pfadfolgeregelung* einteilen. Während die ersten beiden Aufgabenstellungen bereits für sehr viele Systemklassen erforscht und systematisch gelöst sind, stellt die Pfadfolgeregelung ein vergleichsweise junges Forschungsfeld dar.

Das Ziel der Arbeitspunktstabilisierung ist die Stabilisierung eines konstanten oder eines (im Vergleich zur Systemdynamik) nur sehr langsam veränderlichen Arbeitspunktes. Im Gegensatz dazu beschäftigt sich die Trajektorienfolgeregelung mit der Stabilisierung einer explizit zeitabhängigen Solltrajektorie. Typischerweise werden als Regelgrößen die Ausgangsgrößen des Systems verwendet. Die Solltrajektorie gibt nicht nur einen geometrischen Verlauf im Ausgangsraum vor, sondern auch zu welchem Zeitpunkt sich die Ausgangsgrößen des Systems an welchem Punkt des geometrischen Verlaufs befinden sollen.

Im Gegensatz dazu geht es bei der Pfadfolgeregelung darum, dass *a priori* nur der geometrische Verlauf ohne zeitlicher Information (im Weiteren als Pfad bezeichnet) definiert ist. In diesem Kapitel soll stets der Fall betrachtet werden, dass der Pfad für die Ausgangsgrößen des Systems vorgegeben ist. Das primäre Ziel des Pfadfolgereglers besteht darin, den Pfad (asymptotisch) zu stabilisieren, d. h. die Ausgangsgrößen des Systems zum Pfad zu bewegen. Die genaue Stelle entlang des Pfades ist von vornherein nicht definiert. Darauf hinaus soll auch (sofern es die Systemstruktur erlaubt) die Bewegung entlang des Pfades gezielt beeinflusst werden können.

Durch diese Strategie kann ein Manko der Trajektorienfolgeregelung behoben werden. Um dies zu erläutern, wird ein System, das sich zu einem gewissen Zeitpunkt \tilde{t} genau auf der Solltrajektorie befindet, betrachtet. Selbst im idealen Fall ist bei der Trajektorienfolgeregelung nicht gesichert, dass das System auch auf der Solltrajektorie bleibt. Dazu muss nämlich auch die zeitliche Zuordnung stimmen, d. h. gemäß der Solltrajektorie muss sich das System zum Zeitpunkt \tilde{t} genau am richtigen Punkt aufhalten. Stimmt diese Zuordnung nicht, d. h. sollte sich das System zum Zeitpunkt \tilde{t} eigentlich an einem anderen Ort der Solltrajektorie befinden, dann wird diese im Allgemeinen verlassen und dann wieder asymptotisch angefahren. Dadurch, dass es bei der Pfadfolgeregelung *a priori* keine zeitliche Parametrierung des Pfades gibt, wird, falls das System genau auf dem Pfad startet (bzw. auf einer entsprechenden Teilmenge des Zustandsraumes, siehe Abschnitt 1.3), dieser auch nicht mehr verlassen. Diese Eigenschaft bezeichnet man als *Invarianzeigenschaft* eines Pfadfolgereglers, d. h. sobald sich das System einmal auf dem Pfad befindet, gilt dies auch in Zukunft (im nominellen, ungestörten Fall).

Eine typische Anwendung einer Pfadfolgeregelung stellen CNC¹-Fräsmaschinen dar.

¹Computerized Numerical Control.

Mit diesen können üblicherweise komplexe geometrische Formen gefräst werden, die z. B. durch Splines dargestellt sind. Diese definieren eine geometrische Kurve, der das Fräswerkzeug folgen soll. Um eine geforderte Genauigkeit zu erreichen, ist es das primäre Ziel, der geometrischen Kurve (dem Pfad) möglichst gut zu folgen. Erst an zweiter Stelle steht die Geschwindigkeit dieser Verfolgung. Diese wird in diesem Zusammenhang als Vorschubgeschwindigkeit bezeichnet und soll natürlich auch möglichst genau eingehalten werden. Selbst eine ideal eingehaltene Vorschubgeschwindigkeit nützt aber nichts, wenn die Kontur des fertigen Werkstückes die geforderten Genauigkeitsanforderungen nicht erfüllt. Aus dieser Argumentation leitet sich die Priorisierung der Ziele ab. Darüber hinaus möchte man (im nominellen, ungestörten Fall) nicht, dass die ideale Kontur wieder verlassen wird, sobald sich das Werkzeug einmal auf ihr befindet. Dies soll auch gelten, wenn bspw. durch Inhomogenitäten im Material die aktuelle Vorschubgeschwindigkeit nicht exakt der Vorgabe entspricht. Bei einer Trajektorienfolgeregelung würde dies nach sich ziehen, dass die zeitliche Zuordnung entlang des Pfades nicht erfüllt ist, womit die Kontur wieder verlassen werden würde. Die Pfadfolgeregelung behebt diesen Nachteil und eignet sich damit besser für diese Anwendung.

Ein Pfadfolgeregelung kann auf Grundlage von unterschiedlichsten Regelungsverfahren entwickelt werden. In diesem Kapitel werden Konzepte zur Pfadfolgeregelung basierend auf der exakten Eingangs-Ausgangs- und Eingangs-Zustandslinearisierung (vgl. [1.0]) sowie der modellprädiktiven Regelung vorgestellt. Es werden Systeme mit dem Zustand $\mathbf{x} \in \mathbb{R}^n$ und Ausgangsgrößen der Form $\mathbf{y} = \mathbf{h}(\mathbf{x}) \in \mathbb{R}^p$ mit $p \geq 2$ betrachtet. Wie bereits erwähnt, soll die Pfadfolgeregelung stets nur im Ausgangsraum \mathbb{R}^p des Systems, d. h. für die Ausgangsgrößen \mathbf{y} des Systems, durchgeführt werden.

1.1 Definition von Pfaden

Als *Pfad \mathcal{P}* wird eine *Kurve im Ausgangsraum des Systems* bezeichnet, die entweder offen oder geschlossen sein kann, siehe Abbildung 1.1.

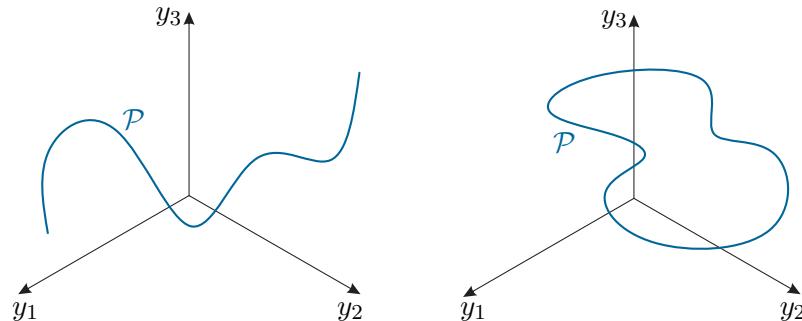


Figure 1.1: Offene und geschlossene Kurve im Ausgangsraum.

Dieser Pfad soll in weiterer Folge vom Pfadfolgeregelung stabilisiert werden, d. h. die Ausgangsgrößen \mathbf{y} des Systems sollen sich zu der Kurve bewegen bzw. sich auf dieser aufhalten. Es gibt die Möglichkeiten der *parametrisierten* und *impliziten* Definition von Pfaden, die im Folgenden näher beleuchtet werden sollen.

1.1.1 Parametrisierte Darstellung von Pfaden

Bei der Parameterdarstellung eines Pfades wird die Kurve über einen reellen Kurven- oder *Pfadparameter* θ dargestellt. Der Kurvenparameter ist Element einer Menge \mathcal{T} , die abgeschlossen (bspw. $\mathcal{T} = [\theta_0, \theta_1]$) oder offen (bspw. $\mathcal{T} = \mathbb{R}$) sein kann. Die Abbildung

$$\sigma(\theta) : \mathcal{T} \rightarrow \mathbb{R}^p \quad (1.1)$$

weist jedem Element aus \mathcal{T} explizit ein Element des Ausgangsraumes zu. Mit Hilfe von (1.1) kann nun der Pfad \mathcal{P} in der Form

$$\mathcal{P} = \{\bar{y} \in \mathbb{R}^p \mid \bar{y} = \sigma(\theta), \theta \in \mathcal{T}\} \quad (1.2)$$

definiert werden. In weiterer Folge wird angenommen, dass die Komponenten $\sigma_i(\theta)$ stets genügend oft stetig differenzierbar sind und \mathcal{P} gemäß (1.2) eine reguläre Kurve [1.0] darstellt, d. h.

$$\sigma'(\bar{\theta}) = \frac{\partial \sigma}{\partial \theta}(\bar{\theta}) \neq \mathbf{0} \quad \forall \bar{\theta} \in \mathcal{T}. \quad (1.3)$$

Der Pfad \mathcal{P} gemäß (1.2) kann je nach Wahl von $\sigma(\theta)$ und \mathcal{T} sowohl offen als auch geschlossen sein.

Example 1.1. Als ein einfaches Beispiel betrachte man für $p = 2$ einen elliptischen Pfad definiert durch

$$\sigma_e(\theta) = \begin{bmatrix} y_{10} + a \cos(\theta) \\ y_{20} + b \sin(\theta) \end{bmatrix}. \quad (1.4)$$

Wie in Abbildung 1.2 dargestellt, ist durch y_{10} und y_{20} der Versatz der Ellipse vom Koordinatenursprung gegeben. Die Parameter $a > 0$ und $b > 0$ stellen die Längen der großen und kleinen Halbachsen dar.

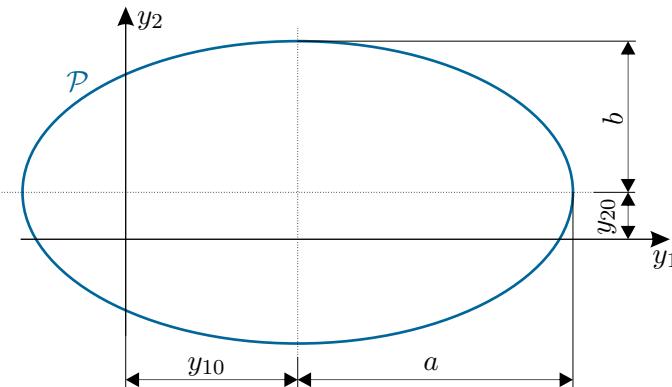


Figure 1.2: Elliptischer Pfad.

Der Definitionsbereich \mathcal{T} des Pfadparameters kann entweder durch $\mathcal{T} = \mathbb{R}$ oder bspw. durch $\mathcal{T} = [0, 2\pi]$ gegeben sein, wobei in beiden Fällen ein geschlossener Pfad resultiert. Für z. B. $\mathcal{T} = [0, \pi]$ resultiert ein offener Pfad. Die partielle Ableitung

folgt zu

$$\frac{\partial \sigma_e}{\partial \theta} = \begin{bmatrix} -a \sin(\theta) \\ b \cos(\theta) \end{bmatrix}. \quad (1.5)$$

Da die beiden Komponenten von (1.5) für beliebige Werte von θ nicht gleichzeitig verschwinden, definiert (1.4) eine reguläre Kurve.

1.1.2 Implizite Definition von Pfaden

Die implizite Definition eines Pfades basiert auf einem Gleichungssystem in Abhängigkeit der Ausgangsgrößen \mathbf{y} des Systems

$$\mathbf{s}(\mathbf{y}) = \mathbf{0} \quad (1.6)$$

mit einer stetigen Abbildung $\mathbf{s} : \mathbb{R}^p \rightarrow \mathbb{R}^{p-1}$. Der Pfad folgt als die Menge

$$\mathcal{P} = \{\bar{\mathbf{y}} \in \mathbb{R}^p \mid \mathbf{s}(\bar{\mathbf{y}}) = \mathbf{0}\}, \quad (1.7)$$

wobei speziell im Hinblick auf die Pfadfolgeregelung gefordert wird, dass die Funktion \mathbf{s}

$$\text{rang}\left(\frac{\partial \mathbf{s}}{\partial \mathbf{y}}(\bar{\mathbf{y}})\right) = p-1 \quad \forall \bar{\mathbf{y}} \in \mathcal{P} \quad (1.8)$$

erfüllt. Der Einfachheit halber sei weiters vorausgesetzt, dass die Komponenten $s_i, i = 1, \dots, p-1$ bis zur jeweils erforderlichen Ordnung stetig differenzierbar sind. Auch hier gilt, dass der Pfad \mathcal{P} gemäß (1.7) je nach Wahl von $\mathbf{s}(\mathbf{y})$ sowohl offen als auch geschlossen sein kann.

Example 1.2. Als Illustration dient wieder ein elliptischer Pfad für $p = 2$ definiert durch

$$s_e(\mathbf{y}) = \frac{(y_1 - y_{10})^2}{a^2} + \frac{(y_2 - y_{20})^2}{b^2} - 1, \quad (1.9)$$

der in Abbildung 1.2 dargestellt ist. Für die partielle Ableitung folgt

$$\frac{\partial s_e}{\partial \mathbf{y}} = \begin{bmatrix} \frac{2(y_1 - y_{10})}{a^2} & \frac{2(y_2 - y_{20})}{b^2} \end{bmatrix}. \quad (1.10)$$

Da auf dem Pfad für $y_1 = y_{10}$ gilt, dass $y_2 \neq y_{20}$ ist (und vice versa), folgt, dass auf \mathcal{P} nie beide Komponenten von (1.10) simultan verschwinden, womit (1.8) erfüllt ist.

1.1.3 Vergleich und Umrechnung der Darstellungen

Die Parameterdarstellung und implizite Definition von Pfaden sind natürlich bis zu einem gewissen Grad äquivalent. Dennoch kann eine der beiden Darstellungen in manchen Situationen geeigneter sein als die andere. Als Beispiel dazu betrachte man einen 8-förmigen Pfad gegeben durch die so genannte Lemniskate von Gerono, siehe Abbildung 1.3. Diese ist definiert durch die Parameterdarstellung

$$\boldsymbol{\sigma}_l(\theta) = \begin{bmatrix} \cos(\theta) \\ \cos(\theta) \sin(\theta) \end{bmatrix} \quad (1.11)$$

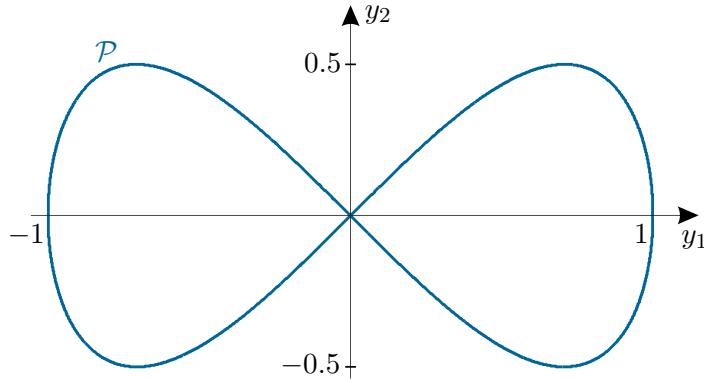


Figure 1.3: Lemniskate von Gerono.

mit $\mathcal{T} = [0, 2\pi)$ bzw. die implizite Darstellung

$$s_l(\mathbf{y}) = y_1^4 - y_1^2 + y_2^2 . \quad (1.12)$$

Die partiellen Ableitungen folgen als

$$\frac{\partial \sigma_l}{\partial \theta}(\theta) = \begin{bmatrix} -\sin(\theta) \\ -\underbrace{\sin^2(\theta) + \cos^2(\theta)}_{\cos(2\theta)} \end{bmatrix} \quad (1.13a)$$

$$\frac{\partial s_l}{\partial \mathbf{y}} = \begin{bmatrix} 4y_1^3 - 2y_1 & 2y_2 \end{bmatrix} . \quad (1.13b)$$

Offensichtlich ist für $\mathbf{y} = \mathbf{0} \in \mathcal{P}$ die Bedingung (1.8) verletzt, womit sich (1.12) nicht für die Pfadfolgeregelung entlang der gesamten Lemniskate eignet. Diese Problematik tritt bei der Parameterdarstellung gemäß (1.11) nicht auf, weil die beiden Komponenten von $\frac{\partial \sigma_l}{\partial \theta}(\theta)$ gemäß (1.13a) für alle $\theta \in \mathcal{T}$ nie gleichzeitig verschwinden. Damit eignet sich die Parameterdarstellung von \mathcal{P} besser für die Pfadfolgeregelung entlang einer Lemniskate.

In vielen Fällen können die Parameterdarstellung und implizite Definition eines Pfades \mathcal{P} auch ineinander umgerechnet werden. Nachfolgend sind die beiden Richtungen erläutert.

1. Parameterdarstellung → implizite Definition: Diese Umrechnung beruht auf der Elimination des Pfadparameters. Gemäß dem *Satz über implizite Funktionen* kann auf Grund von Bedingung (1.3) lokal stets θ aus einer der Gleichungen $\sigma_i(\theta) = y_i$, $i \in \{1, \dots, p\}$ errechnet werden, für die $\frac{\partial \sigma_i}{\partial \theta} \neq 0$ gilt. Daraus folgt $\theta = \sigma_i^{-1}(y_i)$, was eingesetzt in alle anderen Gleichungen $\sigma_j(\theta) = y_j$, $j \in \{1, 2, \dots, p\} \setminus i$ die implizite Definition

$$\sigma_j(\sigma_i^{-1}(y_i)) - y_j = 0, \quad j \in \{1, 2, \dots, p\} \setminus i \quad (1.14)$$

des Pfades in der Form (1.6) liefert.

2. implizite Definition → Parameterdarstellung: Eine Möglichkeit besteht darin, $p - 1$ Ausgangsgrößen aus den Gleichungen $\mathbf{s}(\mathbf{y}) = \mathbf{0}$ in Abhängigkeit einer verbleibenden Ausgangsgröße y_i , $i \in \{1, \dots, p\}$ zu berechnen, was wiederum gemäß dem Satz

über implizite Funktionen und Bedingung (1.8) möglich ist. Damit erhält man $\tilde{\mathbf{y}} = \mathbf{s}^{-1}(y_i)$, wobei $\tilde{\mathbf{y}} \in \mathbb{R}^{p-1}$ alle Ausgangsgrößen bis auf y_i enthält. Anschließend setzt man $\theta = y_i$, womit unmittelbar die Parameterdarstellung des Pfades in der Form (1.1)

$$\begin{bmatrix} \tilde{\mathbf{y}} \\ y_i \end{bmatrix} = \begin{bmatrix} \mathbf{s}^{-1}(\theta) \\ \theta \end{bmatrix} \quad (1.15)$$

folgt.

Example 1.3. Zur Illustration betrachte man wieder den elliptischen Pfad mit den Definitionen gemäß (1.4) bzw. (1.9). Der Einfachheit halber gelte $y_{10} = y_{20} = 0$.

1. Parameterdarstellung → implizite Definition: Die Berechnung von θ aus der ersten Gleichung von (1.4) liefert $\theta = \arccos\left(\frac{y_1}{a}\right)$, was eingesetzt in die zweite Gleichung von (1.4) die gesuchte implizite Darstellung

$$y_2 - b \sin(\theta) = y_2 - b \sin\left(\arccos\left(\frac{y_1}{a}\right)\right) = y_2 - b \sqrt{1 - \frac{y_1^2}{a^2}} = 0 \quad (1.16)$$

ergibt.

2. implizite Definition → Parameterdarstellung: Ein möglicher Ansatz besteht darin, aus $s_e(\mathbf{y}) = 0$ die Ausgangsgröße y_1 in der Form

$$y_1 = \pm a \sqrt{1 - \frac{y_2^2}{b^2}} \quad (1.17)$$

zu berechnen. Über Setzen von $\theta = y_2$ ergibt sich eine Parameterdarstellung zu

$$\boldsymbol{\sigma}_e(\theta) = \begin{bmatrix} \pm a \sqrt{1 - \frac{\theta^2}{b^2}} \\ \theta \end{bmatrix} \quad (1.18)$$

mit $\mathcal{T} = [-b, b]$. Durch Ausnutzen von $\sqrt{1 - x^2} = \cos(\arcsin(x))$ erhält man

$$\boldsymbol{\sigma}_e(\theta) = \begin{bmatrix} \pm a \cos\left(\arcsin\left(\frac{\theta}{b}\right)\right) \\ \theta \end{bmatrix}. \quad (1.19)$$

Substituiert man $\arcsin\left(\frac{\theta}{b}\right) = \tilde{\theta}$ folgt die Beziehung $\theta = b \sin(\tilde{\theta})$ und damit die alternative Parameterdarstellung

$$\tilde{\boldsymbol{\sigma}}_e(\tilde{\theta}) = \begin{bmatrix} \pm a \cos(\tilde{\theta}) \\ b \sin(\tilde{\theta}) \end{bmatrix}, \quad (1.20)$$

die bis auf die Wahl des Vorzeichens identisch zu (1.4) ist.

1.2 Begriffe und Ziele der Pfadfolgeregelung

Für alle folgenden Ausführungen werden AI-Systeme (affine input-Systeme)

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \quad (1.21a)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) \quad (1.21b)$$

mit $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^m$ und $\mathbf{y} \in \mathbb{R}^p$ betrachtet. Es gelte $m \geq p$ womit im Allgemeinen mehr Eingangs- als Ausgangsgrößen erlaubt sind. Für das System (1.21) sei wie in [1.0] vorausgesetzt, dass $\mathbf{f}(\mathbf{x})$ und $\mathbf{g}_j(\mathbf{x})$, die die Spalten von $\mathbf{G}(\mathbf{x}) \in \mathbb{R}^{n \times m}$ darstellen, glatte Vektorfelder sind. Weiters sind die Komponenten von $\mathbf{h}(\mathbf{x})$ glatte Funktionen.

Für die Erläuterung der grundlegenden Konzepte und Ziele der Pfadfolgeregelung wird von einem Pfad \mathcal{P} im Ausgangsraum des Systems ausgegangen. Aus geometrischer Sicht kann jede Bewegung des Ausgangs \mathbf{y} des Systems in *einen transversalen und einen tangentialen* Teil im Hinblick auf den Pfad \mathcal{P} zerlegt werden. Der transversale Teil beschreibt den Anteil der Bewegung, der quer zum Pfad verläuft. Ist dieser Teil identisch null, befindet sich der Ausgang \mathbf{y} des Systems exakt am Pfad. Der tangentiale Teil der Bewegung ist jener, der in Pfadrichtung verläuft. Grob gesprochen kann er als Position des Ausgangs entlang des Pfades gedeutet werden.

Für die Formulierung der Ziele der Pfadfolgeregelung wird die Abbildung $\|\mathbf{y}\|_{\mathcal{P}} : \mathbb{R}^p \rightarrow \mathbb{R}_{\geq 0}$ benötigt, die jedem Punkt \mathbf{y} im Ausgangsraum eine nichtnegative reelle Zahl zuordnet, die den kürzesten Abstand zum Pfad angibt, d. h. $\|\mathbf{y}\|_{\mathcal{P}} = \min_{\bar{\mathbf{y}} \in \mathcal{P}} \|\bar{\mathbf{y}} - \mathbf{y}\|$. Weiters wird noch folgende Definition benötigt.

Definition 1.1 (Gesteuert positiv invariante Menge [1.0]). Eine Teilmenge \mathcal{M} des Zustandsraumes von (1.21) wird *gesteuert positiv invariant* genannt, falls es eine Abbildung $\mathbf{k} : \mathcal{M} \rightarrow \mathbb{R}^m$ so gibt, dass \mathcal{M} eine positiv invariante Menge des autonomen Systems

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{k}(\mathbf{x}) \quad (1.22)$$

ist.

In der englischsprachigen Literatur wird eine gesteuert positiv invariante Menge als *controlled invariant set* bezeichnet. Sie stellt die Verallgemeinerung der positiv invarianten Menge gemäß [1.0] auf Systeme mit Eingangsgrößen dar.

Damit können die Ziele der Pfadfolgeregelung für einen Pfad \mathcal{P} im Ausgangsraum formuliert werden.

- (Z1) *Asymptotische Konvergenz zu \mathcal{P} :* Der Ausgang \mathbf{y} des Systems soll asymptotisch zum Pfad \mathcal{P} konvergieren. Das heißt, $\|\mathbf{y}(t)\|_{\mathcal{P}} \rightarrow 0$ für $t \rightarrow \infty$.
- (Z2) *Invarianzeigenschaft:* Befindet sich der Zustand des Systems $\mathbf{x}(t_0)$ zu einem Zeitpunkt t_0 in einer gesteuert positiv invarianten Teilmenge Γ^* von

$$\Gamma = \{\bar{\mathbf{x}} \in \mathbb{R}^n \mid \mathbf{h}(\bar{\mathbf{x}}) \in \mathcal{P}\}, \quad (1.23)$$

dann fordert die Invarianzeigenschaft, dass für alle $t \geq t_0$ die Eigenschaft $\|\mathbf{y}(t)\|_{\mathcal{P}} = 0$ gelten soll.

(Z3) *Tangentielle Bewegung:* Es soll eine applicationsspezifische Bewegung des Systems auf dem Pfad \mathcal{P} erzielt werden.

Zum Ziel (Z2) sei erwähnt, dass die Menge Γ im Allgemeinen keine gesteuert positiv invariante Menge darstellt. Das Ziel (Z3) kann erst dann genauer spezifiziert werden, wenn der Pfad festgelegt ist und hängt damit sehr stark von der betrachteten Anwendung ab. Darüber hinaus spielen unter anderem auch die Systemstruktur und die Anzahl der Eingangsgrößen eine Rolle. Diese bestimmen auch, welche Bewegung auf dem Pfad überhaupt realisiert werden kann.

1.3 Pfadfolgeregelung basierend auf exakter Linearisierung

Ein Ansatz für den systematischen Entwurf von Pfadfolgereglern beruht auf der Verwendung von Konzepten der exakten Eingangs-Ausgangs- und Eingangs-Zustandslinearisierung. Das System wird dabei in eine spezifische Normalform transformiert. In den transformierten Koordinaten ist der weitere Entwurf eines (Pfadfolge-)Reglers für die Bewegung zum/auf dem Pfad wesentlich vereinfacht. Bezuglich der Herangehensweise sind gewisse Unterschiede für parametrierte und implizit definierte Pfade gegeben. Die gesuchte bzw. resultierende Normalform ist aber qualitativ die selbe.

1.3.1 Erweiterung der exakten Linearisierung für AI-Systeme

In [1.0] ist die exakte Eingangs-Ausgangs- und Eingangs-Zustandslinearisierung für AI-Systeme

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \quad (1.24a)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) \quad (1.24b)$$

mit $p = m$, d. h. gleich vielen Eingangs- wie Ausgangsgrößen, erläutert. Im Folgenden soll das Konzept der exakten Linearisierung auf Systeme der Form (1.24) mit $m \geq p$ erweitert werden. Um auf die Erweiterung hinzuweisen, werden in weiterer Folge Systeme der Form (1.24) mit mehr Eingangs- als Ausgangsgrößen (d. h. $m > p$) als *nichtquadratische Systeme* bezeichnet. Alle kommenden Resultate gelten aber auch für den Spezialfall $p = m$ (quadratisches System) mit den gleichen Resultaten wie in [1.0].

Definition 1.2 (Relativer Grad eines nichtquadratischen AI-Systems). Das nichtquadratische System (1.24) hat den vektoriellen relativen Grad $\{r_1, r_2, \dots, r_p\}$ mit $r = \sum_{j=1}^p r_j \leq n$ an der Stelle $\bar{\mathbf{x}}$, wenn

(A) $L_{\mathbf{g}_j} L_{\mathbf{f}}^k h_i(\mathbf{x}) = 0, j = 1, \dots, m, i = 1, \dots, p, k = 0, \dots, r_i - 2$ für alle \mathbf{x} in einer Umgebung \mathcal{U} von $\bar{\mathbf{x}}$ und

(B) die $(p \times m)$ -Entkopplungsmatrix

$$\mathbf{D}(\mathbf{x}) = \begin{bmatrix} L_{\mathbf{g}_1} L_{\mathbf{f}}^{r_1-1} h_1(\mathbf{x}) & L_{\mathbf{g}_2} L_{\mathbf{f}}^{r_1-1} h_1(\mathbf{x}) & \cdots & L_{\mathbf{g}_m} L_{\mathbf{f}}^{r_1-1} h_1(\mathbf{x}) \\ L_{\mathbf{g}_1} L_{\mathbf{f}}^{r_2-1} h_2(\mathbf{x}) & L_{\mathbf{g}_2} L_{\mathbf{f}}^{r_2-1} h_2(\mathbf{x}) & \cdots & L_{\mathbf{g}_m} L_{\mathbf{f}}^{r_2-1} h_2(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ L_{\mathbf{g}_1} L_{\mathbf{f}}^{r_p-1} h_p(\mathbf{x}) & L_{\mathbf{g}_2} L_{\mathbf{f}}^{r_p-1} h_p(\mathbf{x}) & \cdots & L_{\mathbf{g}_m} L_{\mathbf{f}}^{r_p-1} h_p(\mathbf{x}) \end{bmatrix} \quad (1.25)$$

vollen Rang an der Stelle $\mathbf{x} = \bar{\mathbf{x}}$ hat, d. h. $\text{rang}(\mathbf{D}(\bar{\mathbf{x}})) = p$.

Besitzt ein nichtquadratisches System (1.24) den vektoriellen relativen Grad $\{r_1, r_2, \dots, r_p\}$ dann kommen in den zeitlichen Ableitungen der Ausgangsgrößen $y_j^{(k)}$, $j = 1, \dots, p$, $k = 0, \dots, r_j - 1$, keine Eingangsgrößen u_i , $i = 1, \dots, m$ vor. Darüber hinaus gilt

$$\begin{bmatrix} y_1^{(r_1)} \\ \vdots \\ y_{p-1}^{(r_{p-1})} \\ y_p^{(r_p)} \end{bmatrix} = \underbrace{\begin{bmatrix} L_{\mathbf{f}}^{r_1} h_1(\mathbf{x}) \\ \vdots \\ L_{\mathbf{f}}^{r_{p-1}} h_{p-1}(\mathbf{x}) \\ L_{\mathbf{f}}^{r_p} h_p(\mathbf{x}) \end{bmatrix}}_{\mathbf{b}(\mathbf{x})} + \mathbf{D}(\mathbf{x}) \underbrace{\begin{bmatrix} u_1 \\ \vdots \\ u_{m-1} \\ u_m \end{bmatrix}}_{\mathbf{u}}. \quad (1.26)$$

In einer Umgebung von $\bar{\mathbf{x}}$ kann damit mit Hilfe des Zustandsregelgesetzes

$$\mathbf{u} = \mathbf{D}^T(\mathbf{x}) (\mathbf{D}(\mathbf{x}) \mathbf{D}^T(\mathbf{x}))^{-1} (\mathbf{v} - \mathbf{b}(\mathbf{x})) \quad (1.27)$$

ein *exakt lineares Eingangs-Ausgangsverhalten* vom neuen Eingang $\mathbf{v}^T = [v_1 \ \cdots \ v_p]$ zum Ausgang $\mathbf{y}^T = [y_1 \ \cdots \ y_p]$ in Form von p Integratorketten der Länge r_j , $j = 1, \dots, p$,

$$\begin{bmatrix} y_1^{(r_1)} \\ \vdots \\ y_{p-1}^{(r_{p-1})} \\ y_p^{(r_p)} \end{bmatrix} = \begin{bmatrix} v_1 \\ \vdots \\ v_{p-1} \\ v_p \end{bmatrix} \quad (1.28)$$

erzeugt werden. Die Wahl eines Diffeomorphismus in der Form

$$\mathbf{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_n \end{bmatrix} = \begin{bmatrix} \xi \\ \eta \end{bmatrix} = \begin{bmatrix} \xi_{1,1} \\ \xi_{1,2} \\ \vdots \\ \xi_{1,r_1} \\ \xi_{2,1} \\ \vdots \\ \xi_{p,1} \\ \vdots \\ \xi_{p,r_p} \\ \eta_1 \\ \vdots \\ \eta_{n-r} \end{bmatrix} = \Phi(\mathbf{x}) = \begin{bmatrix} h_1(\mathbf{x}) \\ L_f h_1(\mathbf{x}) \\ \vdots \\ L_f^{r_1-1} h_1(\mathbf{x}) \\ h_2(\mathbf{x}) \\ \vdots \\ L_f^{r_2-1} h_2(\mathbf{x}) \\ \vdots \\ h_p(\mathbf{x}) \\ \vdots \\ L_f^{r_p-1} h_p(\mathbf{x}) \\ \phi_{r+1}(\mathbf{x}) \\ \vdots \\ \phi_n(\mathbf{x}) \end{bmatrix} \quad (1.29)$$

und die Zustandsrückführung (1.27) transformieren das System (1.24) in die Byrnes-Isidori Normalform, die im Allgemeinen die Form

$$\dot{\xi} = \mathbf{A}\xi + \mathbf{B}\mathbf{v} \quad (1.30a)$$

$$\dot{\eta} = \tilde{\mathbf{q}}(\eta, \xi) + \tilde{\mathbf{P}}(\eta, \xi)\mathbf{v} \quad (1.30b)$$

aufweist. Das Teilsystem (1.30a) liegt in Brunovsky Normalform vor.

In [1.0] sind die Bedingungen angeführt unter denen die Methode der exakten Eingangs-Ausgangslinearisierung zu einem stabilen geschlossenen Kreis führt. Unter analogen Anforderungen kann auch in der Normalform (1.30) ein Regler entworfen werden, der einen stabilen geschlossenen Kreis liefert. Allerdings weist die Verwendung der Zustandsrückführung (1.27) einen Nachteil auf. Das System (1.24) in den ursprünglichen Koordinaten weist m Freiheitsgrade in Form der Eingangsgrößen $\mathbf{u} \in \mathbb{R}^m$ auf. Wählt man die Zustandsrückführung (1.27), dann gilt dies für das System (1.30) in transformierten Koordinaten nicht mehr. Es stehen nur mehr p neue Eingangsgrößen $\mathbf{v} \in \mathbb{R}^p$ gemäß der Anzahl der zu regelnden Ausgangsgrößen $\mathbf{y} \in \mathbb{R}^p$ zur Verfügung. Damit verliert man Freiheitsgrade, die man nutzen könnte, um bspw. das System (1.30b) gezielt zu beeinflussen. Nachfolgend wird gezeigt, dass dieser Nachteil behoben werden kann und eine Normalform existiert, die für die Pfadfolgeregelung geeigneter ist. Für die Formulierung des Ergebnisses wird noch nachstehende Definition benötigt.

Definition 1.3 (Äquivalenz unter statischer Rückführung). Zwei dynamische Systeme der Form

$$\Sigma : \dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \quad (1.31a)$$

$$\hat{\Sigma} : \dot{\hat{\mathbf{x}}} = \hat{\mathbf{f}}(\hat{\mathbf{x}}) + \hat{\mathbf{G}}(\hat{\mathbf{x}})\hat{\mathbf{u}} \quad (1.31b)$$

mit $\mathbf{x}, \hat{\mathbf{x}} \in \mathbb{R}^n$ und $\mathbf{u}, \hat{\mathbf{u}} \in \mathbb{R}^m$ sind lokal äquivalent unter statischer Rückführung an einer Stelle $\bar{\mathbf{x}}$ falls in einer Umgebung \mathcal{U} von $\bar{\mathbf{x}}$

(A) eine reguläre Zustandsrückführung

$$\mathbf{u} = \boldsymbol{\alpha}(\mathbf{x}) + \boldsymbol{\Lambda}(\mathbf{x})\hat{\mathbf{u}} \quad (1.32)$$

mit glatten Abbildungen $\boldsymbol{\alpha}$ und $\boldsymbol{\Lambda}$ sowie $\boldsymbol{\Lambda}(\mathbf{x})$ invertierbar für alle \mathbf{x} in der Umgebung \mathcal{U} von $\bar{\mathbf{x}}$ und

(B) ein Diffeomorphismus $\hat{\mathbf{x}} = \boldsymbol{\Phi}(\mathbf{x})$

so existieren, dass

$$\hat{\mathbf{f}}(\hat{\mathbf{x}}) = \frac{\partial \boldsymbol{\Phi}}{\partial \mathbf{x}}(\mathbf{x})(\mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\boldsymbol{\alpha}(\mathbf{x})) \Big|_{\mathbf{x}=\boldsymbol{\Phi}^{-1}(\hat{\mathbf{x}})} \quad (1.33a)$$

$$\hat{\mathbf{G}}(\hat{\mathbf{x}}) = \frac{\partial \boldsymbol{\Phi}}{\partial \mathbf{x}}(\mathbf{x})\mathbf{G}(\mathbf{x})\boldsymbol{\Lambda}(\mathbf{x}) \Big|_{\mathbf{x}=\boldsymbol{\Phi}^{-1}(\hat{\mathbf{x}})}. \quad (1.33b)$$

Man sagt, die beiden Systeme Σ und $\hat{\Sigma}$ sind äquivalent unter statischer Rückführung auf \mathcal{U} .

Exercise 1.1. Rechnen Sie die Beziehungen (1.33) nach.

Der folgende Satz aus [1.0] liefert nun die gewünschte Normalform.

Theorem 1.1 (Exakte Eingangs-Ausgangslinearisierung für nichtquadratische AI-Systeme). Angenommen, das System (1.24) hat vektoriellen relativen Grad $\{r_1, r_2, \dots, r_p\}$ mit $r = \sum_{j=1}^p r_j \leq n$ an der Stelle $\bar{\mathbf{x}}$. Dann gilt, dass (1.24) in einer Umgebung \mathcal{U} von $\bar{\mathbf{x}}$ äquivalent unter statischer Rückführung zu einem System der Form

$$\dot{\boldsymbol{\xi}} = \mathbf{A}\boldsymbol{\xi} + \mathbf{B}\mathbf{v}_1 \quad (1.34a)$$

$$\dot{\boldsymbol{\eta}} = \mathbf{f}^0(\boldsymbol{\eta}, \boldsymbol{\xi}) + \mathbf{G}_1(\boldsymbol{\eta}, \boldsymbol{\xi})\mathbf{v}_1 + \mathbf{G}_2(\boldsymbol{\eta}, \boldsymbol{\xi})\mathbf{v}_2 \quad (1.34b)$$

ist. Dabei gilt $\boldsymbol{\xi} \in \mathbb{R}^r$, $\boldsymbol{\eta} \in \mathbb{R}^{n-r}$, $\mathbf{v}_1 \in \mathbb{R}^p$ und $\mathbf{v}_2 \in \mathbb{R}^{m-p}$. Das Teilsystem (1.34a) liegt in Brunovsky Normalform vor. In den neuen Koordinaten gilt mit einer Partitionierung von $\boldsymbol{\xi}$ gemäß (1.29), dass $y_i = \xi_{i,1}$, $i = 1, \dots, p$.

Proof. Analog zu [1.0] gilt auch hier, dass ein Diffeomorphismus $\mathbf{z} = \boldsymbol{\Phi}(\mathbf{x})$ gemäß (1.29) in einer Umgebung \mathcal{U} von $\bar{\mathbf{x}}$ existiert. Damit folgt die Systemdynamik in

transformierten Koordinaten zu

$$\begin{bmatrix} \dot{\xi}_{1,r_1} \\ \dot{\xi}_{2,r_2} \\ \vdots \\ \dot{\xi}_{p,r_p} \end{bmatrix} = \begin{bmatrix} L_f^{r_1} h_1(\mathbf{x}) \\ L_f^{r_2} h_2(\mathbf{x}) \\ \vdots \\ L_f^{r_p} h_p(\mathbf{x}) \end{bmatrix}_{\mathbf{x}=\Phi^{-1}(\mathbf{z})} + \mathbf{D}(\mathbf{x})|_{\mathbf{x}=\Phi^{-1}(\mathbf{z})} \mathbf{u} \quad (1.35a)$$

$$\dot{\xi}_{i,j} = \xi_{i,j+1}, \quad i = 1, \dots, p, \quad j = 1, \dots, r_i - 1 \quad (1.35b)$$

$$\dot{\eta} = \mathbf{q}(\boldsymbol{\eta}, \boldsymbol{\xi}) + \mathbf{P}(\boldsymbol{\eta}, \boldsymbol{\xi}) \mathbf{u} \quad (1.35c)$$

mit

$$q_i(\boldsymbol{\eta}, \boldsymbol{\xi}) = L_f \phi_{r+i}(\mathbf{x})|_{\mathbf{x}=\Phi^{-1}(\mathbf{z})}, \quad i = 1, \dots, n - r \quad (1.35d)$$

$$P_{i,j}(\boldsymbol{\eta}, \boldsymbol{\xi}) = L_g_j \phi_{r+i}(\mathbf{x})|_{\mathbf{x}=\Phi^{-1}(\mathbf{z})}, \quad i = 1, \dots, n - r, \quad j = 1, \dots, m. \quad (1.35e)$$

Gemäß der Annahme folgt, dass die Entkopplungsmatrix \mathbf{D} in einer Umgebung \mathcal{U} von $\bar{\mathbf{x}} = \Phi^{-1}(\bar{\mathbf{z}})$ vollen Rang besitzt. Sei $\mathbf{N}(\mathbf{x}) \in \mathbb{R}^{m \times (m-p)}$ eine Matrix, deren Spalten den Nullraum von $\mathbf{D}(\mathbf{x})$ aufspannen. Damit lässt sich die in der Umgebung \mathcal{U} von $\bar{\mathbf{x}}$ reguläre Matrix $\boldsymbol{\Lambda}(\mathbf{x}) = [\mathbf{M}(\mathbf{x}) \ \mathbf{N}(\mathbf{x})]$ mit $\mathbf{M}(\mathbf{x}) = \mathbf{D}^T(\mathbf{x}) (\mathbf{D}(\mathbf{x}) \mathbf{D}^T(\mathbf{x}))^{-1} \in \mathbb{R}^{m \times p}$ definieren. Mit

$$\boldsymbol{\alpha}(\mathbf{x}) = -\boldsymbol{\Lambda}(\mathbf{x}) \begin{bmatrix} L_f^{r_1} h_1(\mathbf{x}) \\ L_f^{r_2} h_2(\mathbf{x}) \\ \vdots \\ L_f^{r_p} h_p(\mathbf{x}) \\ \mathbf{0} \end{bmatrix}, \quad (1.36)$$

wobei der Nullvektor in (1.36) $m - p$ Zeilen aufweist, folgt die gesuchte Zustandsrückführung

$$\mathbf{u} = \boldsymbol{\alpha}(\mathbf{x})|_{\mathbf{x}=\Phi^{-1}(\mathbf{z})} + \boldsymbol{\Lambda}(\mathbf{x})|_{\mathbf{x}=\Phi^{-1}(\mathbf{z})} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix}, \quad (1.37)$$

die die gewünschte Normalform (1.34) in $(\boldsymbol{\xi}, \boldsymbol{\eta})$ -Koordinaten liefert. \square

Exercise 1.2. Rechnen Sie nach, dass die im Beweis von Satz 1.1 angegebene Rückführung tatsächlich auf die gewünschte Normalform führt.

In der Normalform (1.34) kommen im Gegensatz zu (1.30) nicht nur die neuen Eingangsgrößen $\mathbf{v}_1 \in \mathbb{R}^p$ sondern auch $\mathbf{v}_2 \in \mathbb{R}^{m-p}$ vor. Das bedeutet, dass die zusätzlichen Freiheitsgrade im Fall $m > p$ nicht verloren gehen. Sie können in Form von \mathbf{v}_2 ggf. dafür genutzt werden, das Teilsystem (1.34b) zu beeinflussen. Basierend auf der Normalform (1.34) werden nachfolgend Pfadfolgeregler für implizit definierte und parametrierte Pfade entworfen.

1.3.2 Implizit definierte Pfade

Ein möglicher Ansatz zum Entwurf von Pfadfolgeregulern für implizit definierte Pfade besteht darin, neue fiktive Ausgangsgrößen

$$\hat{\mathbf{y}} = \hat{\mathbf{h}}(\mathbf{x}) = \mathbf{s}(\mathbf{h}(\mathbf{x})) \quad (1.38)$$

mit $\hat{\mathbf{y}} \in \mathbb{R}^{p-1}$ zu definieren. Für die nachfolgenden Ausführungen sind zwei Annahmen essenziell.

1. Es wird vorausgesetzt, dass das System $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$ mindestens $p - 1$ Eingangsgrößen aufweist, d. h. $m \geq p - 1$, $p > 1$.
2. Das System $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$ mit der Ausgangsgröße $\hat{\mathbf{y}}$ gemäß (1.38) hat in Γ^* einen wohldefinierten vektoriellen relativen Grad gemäß Definition 1.2.

Unter diesen Annahmen kann Satz 1.1 angewandt und das System $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$ mit dem Ausgang (1.38) auf die Normalform (1.34) transformiert werden. Diese Normalform stellt sich im gegebenen Zusammenhang in der Form

$$\dot{\boldsymbol{\xi}} = \mathbf{A}\boldsymbol{\xi} + \mathbf{B}\mathbf{v}^\dagger \quad (1.39a)$$

$$\dot{\boldsymbol{\eta}} = \mathbf{f}^0(\boldsymbol{\eta}, \boldsymbol{\xi}) + \mathbf{G}^\dagger(\boldsymbol{\eta}, \boldsymbol{\xi})\mathbf{v}^\dagger + \mathbf{G}^\parallel(\boldsymbol{\eta}, \boldsymbol{\xi})\mathbf{v}^\parallel \quad (1.39b)$$

mit $\mathbf{v}^\dagger \in \mathbb{R}^{p-1}$ und $\mathbf{v}^\parallel \in \mathbb{R}^{m-(p-1)}$ dar. Die Darstellung (1.39) wird in der englischsprachigen Literatur als *Transverse Normal Form* bezeichnet und soll hier analog dazu *transversale Normalform* genannt werden. Die zugehörige Zustandsrückführung stellt sich gemäß (1.32) in der Form

$$\mathbf{u} = \boldsymbol{\alpha}(\mathbf{x}) + \boldsymbol{\Lambda}(\mathbf{x}) \begin{bmatrix} \mathbf{v}^\dagger \\ \mathbf{v}^\parallel \end{bmatrix} \quad (1.40)$$

dar.

Das Differenzialgleichungssystem (1.39a) stellt die *transversale Dynamik*, d. h. die Dynamik quer zum Pfad \mathcal{P} , dar. Dies begründet sich in der Tatsache, dass durch die Koordinaten $\xi_{i,1}$, $i = 1, \dots, p - 1$ gemäß (1.35a) und (1.35b) direkt die fiktiven Ausgangsgrößen $\hat{\mathbf{y}}$ in der Form

$$\begin{bmatrix} \xi_{1,1} \\ \xi_{2,1} \\ \vdots \\ \xi_{p-1,1} \end{bmatrix} = \hat{\mathbf{h}}(\mathbf{x}) \quad (1.41)$$

gegeben sind. Falls

$$\xi_{i,1} = 0, \quad i = 1, \dots, p - 1 \quad (1.42)$$

gilt, befindet sich der Ausgang \mathbf{y} des Systems exakt am Pfad \mathcal{P} . Ist (1.42) nicht für alle $\xi_{i,1}$, $i = 1, \dots, p - 1$ erfüllt, bewegt sich der Ausgang \mathbf{y} abseits vom Pfad. Die verbleibenden Koordinaten in $\boldsymbol{\xi}$ stellen entsprechend dem vektoriellen relativen Grad die Ableitungen von $\hat{\mathbf{h}}(\mathbf{x})$ dar. Damit ist unmittelbar einsichtig, dass die Koordinaten $\boldsymbol{\xi}$ der

Dynamik quer (transversal) zum Pfad zugeordnet sind. Die transversale Dynamik (1.39a) liegt entsprechend Satz 1.1 in Brunovsky Normalform vor, d. h. das System (1.39a) ist vollständig erreichbar. Damit ist das Ziel (Z1) der Pfadfolgeregelung leicht durch den Entwurf eines asymptotisch stabilisierenden linearen Zustandsreglers in der Form

$$v_i^\dagger = -k_{i,0} \int_{t_0}^t \xi_{i,1} dt - \sum_{j=1}^{r_i} k_{i,j} \xi_{i,j} \quad i = 1, \dots, p-1 \quad (1.43)$$

mit geeigneten Koeffizienten $k_{i,j} > 0$ erfüllbar. Das Regelgesetz (1.43) beinhaltet einen Integralanteil, um z. B. mit Modellabweichungen besser umgehen zu können und damit besser für einen praktischen Einsatz geeignet zu sein.

Um die Dynamik (1.39b) zu erhalten, müssen die Funktionen $\phi_{r+i}(\mathbf{x})$, $i = 1, \dots, n-r$ zur Komplettierung des Diffeomorphismus

$$\begin{bmatrix} \boldsymbol{\xi} \\ \boldsymbol{\eta} \end{bmatrix} = \boldsymbol{\Phi}(\mathbf{x}) = \begin{bmatrix} \boldsymbol{\Phi}_\xi(\mathbf{x}) \\ \boldsymbol{\Phi}_\eta(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Phi}_\xi(\mathbf{x}) \\ \phi_{r+1}(\mathbf{x}) \\ \vdots \\ \phi_n(\mathbf{x}) \end{bmatrix} \quad (1.44)$$

gewählt werden. Da für $\boldsymbol{\xi} = \mathbf{0}$ unabhängig vom Wert von $\boldsymbol{\eta}$ die Ausgangsgrößen \mathbf{y} exakt am Pfad liegen, stellt (1.39b) die *tangentielle Dynamik* dar. Die Koordinaten $\eta_i = \phi_{r+i}(\mathbf{x})$, $i = 1, \dots, n-r$ werden entsprechend dem vorliegenden Pfad und der applikationsspezifischen Ziele (Bewegung auf \mathcal{P}) gewählt. Dabei ist noch die zusätzliche Nebenbedingung zu beachten, dass durch (1.44) tatsächlich ein Diffeomorphismus gegeben ist. Im Fall $m > p-1$ sind die Eingangsgrößen \mathbf{v}^\parallel in (1.39b) vorhanden und damit können geeignete Regler entworfen werden, um (Z3) zu erfüllen.

Das noch offene Ziel der Pfadfolgeregelung ist (Z2), die Invarianzeigenschaft. Befindet sich der Zustand des Systems zum Zeitpunkt t_0 in der Menge

$$\Gamma = \{\bar{\mathbf{x}} \in \mathbb{R}^n \mid \mathbf{s}(\mathbf{h}(\bar{\mathbf{x}})) = \mathbf{0}\}, \quad (1.45)$$

d. h. die Ausgangsgrößen befinden sich exakt am Pfad \mathcal{P} , ist es im Allgemeinen nicht garantiert, dass dies für alle $t > t_0$ der Fall ist. Dafür müssen auch alle Ableitungen von $\mathbf{h}(\mathbf{x})$ zum Zeitpunkt t_0 null sein. Aus dieser Überlegung folgt die Definition der Menge

$$\Gamma^* = \{\bar{\mathbf{x}} \in \mathbb{R}^n \mid \boldsymbol{\Phi}_\xi(\bar{\mathbf{x}}) = \mathbf{0}\} \subset \Gamma. \quad (1.46)$$

Gilt nun, dass $\mathbf{x}(t_0) \in \Gamma^*$, dann kann durch die Wahl $\mathbf{v}^\dagger \equiv \mathbf{0}$ erreicht werden, dass $\boldsymbol{\xi} \equiv \mathbf{0} \forall t > t_0$ wie leicht an (1.39a) erkennbar ist. Damit gilt, dass die Menge Γ^* gesteuert positiv invariant ist und der Pfad \mathcal{P} für alle Zeiten $t > t_0$ nicht mehr verlassen wird, was die Erfüllung der Invarianzeigenschaft (Z2) impliziert.

Damit ist gezeigt, dass basierend auf der Normalform (1.39), die unter den getroffenen Annahmen stets systematisch erreichbar ist, einfach ein Pfadfolgeregler für die Erfüllung der Ziele (Z1)–(Z3) entworfen werden kann.

Exercise 1.3. Zeigen Sie, dass mit der gezeigten Vorgehensweise die Erfüllung der Ziele (Z1)–(Z3) nicht widersprüchlich ist, d. h. das Erreichen eines Ziels beeinflusst nicht die anderen beiden.

Die Menge Γ^* gemäß (1.46) stellt eine Untermannigfaltigkeit des Zustandsraumes dar. In der englischsprachigen Literatur wird sie gelegentlich als *zero path error manifold* oder auch *path following manifold* bezeichnet. Die Dynamik des Systems auf dieser Mannigfaltigkeit ist durch

$$\dot{\boldsymbol{\eta}} = \mathbf{f}^0(\boldsymbol{\eta}, \mathbf{0}) + \mathbf{G}^{\parallel}(\boldsymbol{\eta}, \mathbf{0})\mathbf{v}^{\parallel} \quad (1.47)$$

gegeben. Sie entspricht der internen Dynamik des Systems auf dem Pfad \mathcal{P} .

Example 1.4 (Fahrzeug mit variabler Vorwärtsgeschwindigkeit). Die gezeigte Theorie soll für die Pfadfolgeregelung eines Fahrzeuges, beschrieben durch das mathematische Modell

$$\dot{\mathbf{x}} = \frac{d}{dt} \begin{bmatrix} x \\ y \\ \varphi \\ v_l \end{bmatrix} = \underbrace{\begin{bmatrix} x \\ y \\ 0 \\ 0 \end{bmatrix}}_{\mathbf{f}_v(\mathbf{x})} + \underbrace{\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ v_l & 0 \\ 0 & 1 \end{bmatrix}}_{\mathbf{G}_v(\mathbf{x})} \mathbf{u} \quad (1.48a)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) = \begin{bmatrix} x \\ y \end{bmatrix}, \quad (1.48b)$$

angewandt werden. Die Stellgröße u_1 kann mit dem Lenkwinkel in Verbindung gebracht werden während u_2 direkt die zeitliche Ableitung der Vorwärtsgeschwindigkeit $v_l \neq \text{konst. angibt}$. Der Winkel φ beschreibt die Orientierung des Fahrzeuges relativ zur x -Achse. Der Ausgang \mathbf{y} ist durch die x - und y -Position des Fahrzeuges gegeben.

Im konkreten Fall soll der Pfadfolgeregelung das Fahrzeug zu und entlang einer Ellipse gegeben durch die implizite Gleichung

$$s_e(\mathbf{y}) = \frac{(y_1 - y_{10})^2}{a^2} + \frac{(y_2 - y_{20})^2}{b^2} - 1 \quad (1.49)$$

Führen. Als Ziel (Z3) wird für diese Anwendung gewählt, dass die Vorwärtsgeschwindigkeit auf dem Pfad einer hinreichend oft stetig differenzierbaren Solltrajektorie $v_{l,d}(t)$ folgen soll.

Der erste Schritt besteht darin, die am Beginn von Abschnitt 1.3.2 getroffenen Annahmen zu verifizieren. Die Voraussetzung $m \geq p - 1 = 1$ ist offensichtlich erfüllt. Für die Überprüfung der zweiten Annahme wird die fiktive Ausgangsgröße

$$\hat{y} = s_e(\mathbf{h}(\mathbf{x})) = \frac{(x - x_{10})^2}{a^2} + \frac{(y - y_{20})^2}{b^2} - 1 = \hat{h}(\mathbf{x}) \quad (1.50)$$

benötigt. Nun muss überprüft werden, ob das System (1.48a) mit der Ausgangsgröße (1.50) einen wohldefinierten relativen Grad in Γ^* aufweist. Dies soll gemeinsam mit

der Berechnung von Γ^* gemäß (1.46) im nächsten Schritt durchgeführt werden. Dazu berechnet man zuerst die Ableitungen von (1.50) in der Form

$$\dot{y} = L_{\mathbf{f}_v} \hat{h}(\mathbf{x}) + \underbrace{L_{\mathbf{g}_{v,1}} \hat{h}(\mathbf{x}) u_1}_{0} + \underbrace{L_{\mathbf{g}_{v,2}} \hat{h}(\mathbf{x}) u_2}_{0} = 2v_l \left(\frac{(x - y_{10}) \cos \varphi}{a^2} + \frac{(y - y_{20}) \sin \varphi}{b^2} \right) \quad (1.51a)$$

$$\ddot{y} = L_{\mathbf{f}_v}^2 \hat{h}(\mathbf{x}) + L_{\mathbf{g}_{v,1}} L_{\mathbf{f}_v} \hat{h}(\mathbf{x}) u_1 + L_{\mathbf{g}_{v,2}} L_{\mathbf{f}_v} \hat{h}(\mathbf{x}) u_2 = \quad (1.51b)$$

$$= 2v_l^2 \left(\frac{\cos^2 \varphi}{a^2} + \frac{\sin^2 \varphi}{b^2} \right) + \underbrace{\left[2v_l^2 \left(\frac{(y-y_{20}) \cos \varphi}{b^2} - \frac{(x-y_{10}) \sin \varphi}{a^2} \right) \right]}_{\mathbf{D}(\mathbf{x})}^T \mathbf{u}, \quad (1.51c)$$

wobei in \ddot{y} der Eingang \mathbf{u} erscheint. Damit ist Γ^* durch

$$\Gamma^* = \left\{ \bar{\mathbf{x}} \in \mathbb{R}^4 \mid \hat{h}(\bar{\mathbf{x}}) = L_{\mathbf{f}_v} \hat{h}(\bar{\mathbf{x}}) = 0 \right\} \quad (1.52)$$

gegeben. Sei $\bar{\mathbf{x}} = [\bar{x} \ \bar{y} \ \bar{\varphi} \ \bar{v}_l]^T$ ein Punkt, der in Γ^* liegt. Es muss nun noch nachgewiesen werden, dass $\mathbf{D}(\bar{\mathbf{x}})$ vollen Rang für alle $\bar{\mathbf{x}} \in \Gamma^*$ aufweist. Die Bedingung $L_{\mathbf{f}_v} \hat{h}(\bar{\mathbf{x}}) = 0$ gemäß (1.51a) impliziert, dass sich das Fahrzeug tangential zum Pfad bewegt. Dies kann verifiziert werden, indem der Winkel der Tangente an die Ellipse, der durch

$$\bar{\varphi} = \arctan \left(\frac{(\bar{x} - y_{10}) b^2}{-(\bar{y} - y_{20}) a^2} \right) \quad (1.53)$$

gegeben ist, in (1.51a) eingesetzt wird. Drückt man nun beispielsweise aus $\hat{h}(\bar{\mathbf{x}}) = 0$ \bar{x} in Abhängigkeit von \bar{y} aus und setzt diese Beziehung zusammen mit (1.53) in \mathbf{D} ein, dann erhält man

$$\mathbf{D}(\bar{\mathbf{x}}) = \begin{bmatrix} -\frac{2\bar{v}_l^2 \sqrt{a^2(\bar{y}-y_{20})^2+b^4-b^2(\bar{y}-y_{20})^2}}{ab^2} & 0 \end{bmatrix} \quad (1.54)$$

für alle $\bar{\mathbf{x}} \in \Gamma^*$. Damit ist gezeigt, dass die Entkopplungsmatrix für alle $\bar{v}_l \neq 0$ und $y_{20} - b \leq \bar{y} \leq y_{20} + b$ Rang 1 aufweist und damit das System (1.48a) mit der Ausgangsgröße (1.50) einen wohldefinierten relativen Grad $r = 2$ in Γ^* aufweist.

Damit kann direkt Satz 1.1 angewandt werden. Der erste Teil des Diffeomorphismus für die Transformation in die transversale Normalform lautet

$$\Phi_\xi(\mathbf{x}) = \begin{bmatrix} \hat{h}(\mathbf{x}) \\ L_{\mathbf{f}_v} \hat{h}(\mathbf{x}) \end{bmatrix}. \quad (1.55a)$$

Im Hinblick auf das Ziel (Z3) wird die Zustandstransformation mit

$$\Phi_{\eta}(\mathbf{x}) = \begin{bmatrix} v_l \\ \varphi \end{bmatrix} \quad (1.55b)$$

zu

$$\begin{bmatrix} \xi \\ \eta \end{bmatrix} = \Phi(\mathbf{x}) = \begin{bmatrix} \Phi_{\xi}(\mathbf{x}) \\ \Phi_{\eta}(\mathbf{x}) \end{bmatrix} \quad (1.55c)$$

ergänzt.

Exercise 1.4. Verifizieren Sie, dass es sich bei (1.55) tatsächlich um einen Diffeomorphismus in Γ^* handelt.

Für die Berechnung der exakt linearisierenden Zustandsrückführung (1.37) wird zuerst die Matrix $\mathbf{N}(\mathbf{x})$ in der Form

$$\mathbf{N}(\mathbf{x}) = \begin{bmatrix} -\frac{2\left(\frac{(x-y_{10})\cos\varphi}{a^2} + \frac{(y-y_{20})\sin\varphi}{b^2}\right)}{2v_l^2\left(\frac{(y-y_{20})\cos\varphi}{b^2} - \frac{(x-y_{10})\sin\varphi}{a^2}\right)} \\ 1 \end{bmatrix} \quad (1.56)$$

gewählt. Aus $\mathbf{D}(\mathbf{x})$ gemäß (1.51c) lassen sich direkt $\mathbf{M}(\mathbf{x}) = \mathbf{D}^T(\mathbf{x})\left(\mathbf{D}(\mathbf{x})\mathbf{D}^T(\mathbf{x})\right)^{-1} \in \mathbb{R}^{2 \times 1}$ und in weiterer Folge $\Lambda(\mathbf{x}) = [\mathbf{M}(\mathbf{x}) \quad \mathbf{N}(\mathbf{x})]$ errechnen. Damit folgt die Zustandsrückführung zu

$$\mathbf{u} = \boldsymbol{\alpha}(\mathbf{x}) + \Lambda(\mathbf{x}) \begin{bmatrix} v^{\dagger} \\ v^{\parallel} \end{bmatrix} \quad (1.57)$$

mit

$$\boldsymbol{\alpha}(\mathbf{x}) = -\Lambda(\mathbf{x}) \begin{bmatrix} L_{f_v}^2 \hat{h}(\mathbf{x}) \\ 0 \end{bmatrix}. \quad (1.58)$$

Die transversale Normalform ergibt sich im vorliegenden Fall zu

$$\dot{\xi} = \begin{bmatrix} \xi_2 \\ v^{\dagger} \end{bmatrix} \quad (1.59a)$$

$$\dot{\eta} = \mathbf{f}^0(\boldsymbol{\eta}, \xi) + \mathbf{g}^{\dagger}(\boldsymbol{\eta}, \xi)v^{\dagger} + \mathbf{g}^{\parallel}(\boldsymbol{\eta}, \xi)v^{\parallel} \quad (1.59b)$$

wobei die Terme in \mathbf{f}^0 , \mathbf{g}^{\dagger} und \mathbf{g}^{\parallel} bereits sehr umfangreich sind. Im Hinblick auf (Z3) ist vor allem die Einschränkung von (1.59b) auf Γ^* entsprechend (1.47) interessant. Diese interne Dynamik auf der Ellipse ergibt sich zu

$$\dot{\boldsymbol{\eta}} = \begin{bmatrix} v^{\parallel} \\ \pm \frac{\eta_1^2(a^4(\cos\eta_2)^4 - 2a^2b^2(\cos\eta_2)^4 + b^4(\cos\eta_2)^4 - 2a^4(\cos\eta_2)^2 + 2a^2b^2(\cos\eta_2)^2 + a^4)\sin\eta_2}{b^2\sqrt{a^4\eta_1^2(\sin\eta_2)^2(a^2(\sin\eta_2)^2 - b^2(\sin\eta_2)^2 + b^2)}} \end{bmatrix} \quad (1.60)$$

wobei das Vorzeichen von $\dot{\eta}_2$ davon abhängig ist, in welcher Richtung die Ellipse abgefahren wird. Die erste Zeile von (1.60) lautet $\dot{v}_l = v^{\parallel}$ womit sich (Z3) mit dem Regelgesetz

$$v^{\parallel} = \dot{v}_{l,d}(t) - k^{\parallel}(v_l - v_{l,d}(t)) \quad k^{\parallel} > 0 \quad (1.61)$$

erzielen lässt. Die transversale Dynamik wird mit dem Regelgesetz

$$v^{\perp} = -k_1^{\perp}\xi_1 - k_2^{\perp}\xi_2 \quad k_1^{\perp}, k_2^{\perp} > 0 \quad (1.62)$$

stabilisiert. Es sei noch angemerkt, dass man durch das Erreichen des Ziels (Z3) zumindest auf der Ellipse sicherstellen kann, dass $v_l \neq 0$ gilt. Damit ist gewährleistet, dass die Entkopplungsmatrix $\mathbf{D}(\mathbf{x})$ stets vollen Rang in Γ^* aufweist.

Für die im Folgenden dargestellten Ergebnisse werden die Parameter $a = 4$, $b = 2.5$, $y_{10} = 1$, $y_{20} = -1.5$, $k_1^{\perp} = 0.25$, $k_2^{\perp} = 1$ und $k^{\parallel} = 0.5$ verwendet. Weiters gelte $v_{l,d}(t) = 0.5 + (\sin(\frac{t}{2}))^2$. In Abbildung 1.4 ist die Trajektorie des Fahrzeuges für den Anfangswert $\mathbf{x}_0 = [5.5 \quad 1 \quad \frac{3\pi}{2} \quad 0.45]^T$ dargestellt. An der Startposition des Fahrzeuges ist ein Pfeil eingezeichnet, der die anfängliche Orientierung zusätzlich verdeutlichen soll.

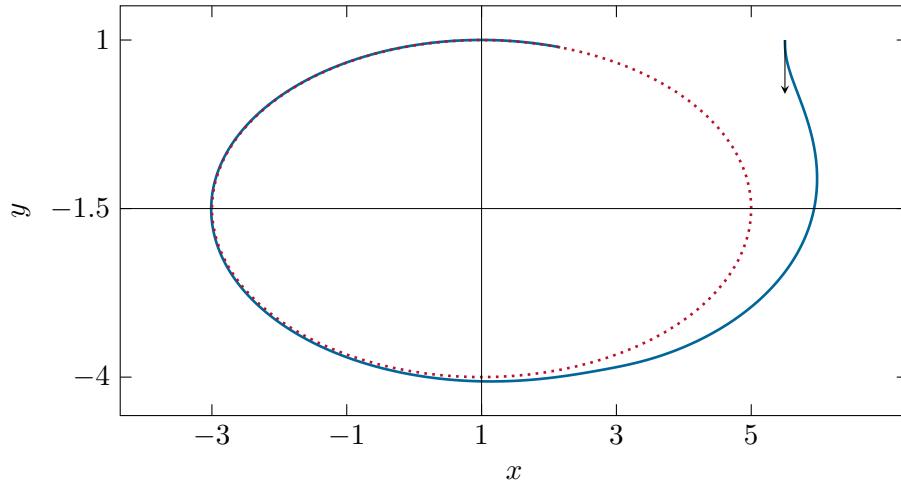


Figure 1.4: Trajektorie des Fahrzeuges mit $v_l \neq \text{konst.}$ für die Pfadfolgeregelung einer Ellipse.

In Abbildung 1.5 sind die zugehörigen Zeitverläufe im geschlossenen Kreis ersichtlich. Die Zustände ξ der transversalen Dynamik konvergieren zu null womit der Zustand \mathbf{x} in die Menge Γ^* gezwungen wird. Gleichzeitig ist ersichtlich, dass $\eta_1 = v_l$ der Solltrajektorie $v_{l,d}(t)$ nachgeführt wird.

Abbildung 1.6 zeigt Trajektorien des Fahrzeuges für verschiedene Anfangswerte. Wenn \mathbf{x}_0 in Γ^* liegt (das Fahrzeug befindet sich auf der Ellipse und ist tangential dazu ausgerichtet) bleibt das Fahrzeug für alle zukünftigen Zeiten auf der Ellipse. Dies zeigt, dass die Invarianzeigenschaft (Z2) erfüllt ist. Je nach Startposition und -orientierung ist die Richtung, in welcher die Ellipse abgefahren wird, eine andere.

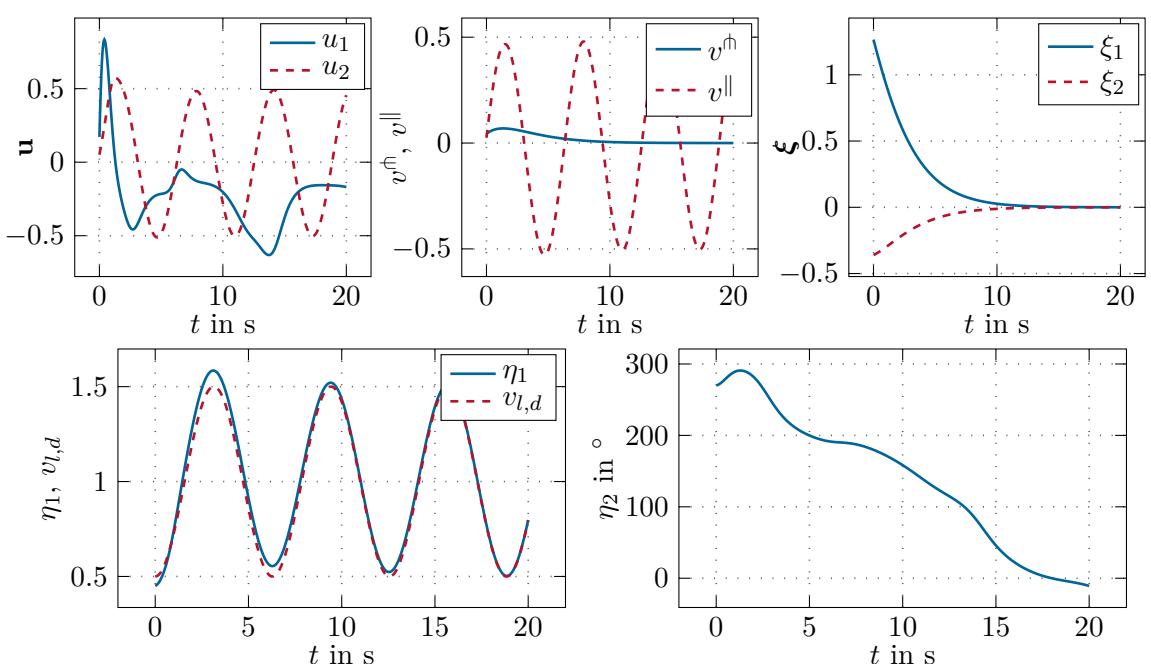


Figure 1.5: Zeitverläufe im geschlossenen Kreis bei der Pfadfolgeregelung einer Ellipse für das Fahrzeug mit $v_l \neq \text{konst.}$

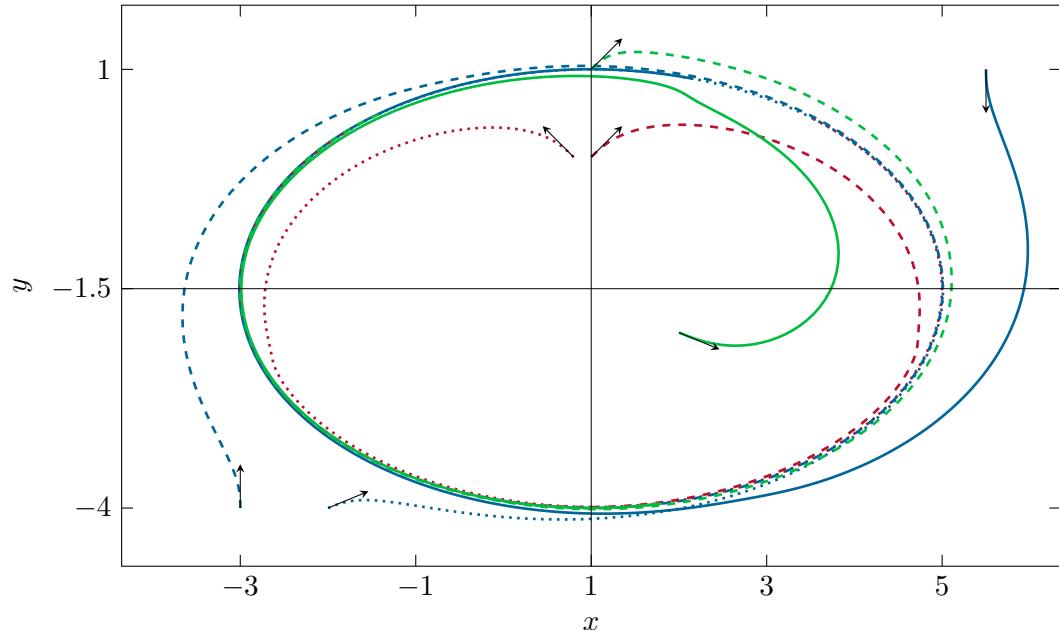


Figure 1.6: Trajektorien des Fahrzeuges mit $v_l \neq \text{konst.}$ für die Pfadfolgeregelung einer Ellipse und verschiedenen Anfangswerten.

Example 1.5 (Fahrzeug mit konstanter Vorwärtsgeschwindigkeit). Im Gegensatz zu Beispiel 1.4 soll hier das Fahrzeug mit konstanter Vorwärtsgeschwindigkeit $v_l = \text{konst.} > 0$ beschrieben durch die Differenzialgleichungen

$$\dot{\mathbf{x}} = \frac{d}{dt} \begin{bmatrix} x \\ y \\ \varphi \end{bmatrix} = \underbrace{\begin{bmatrix} v_l \cos \varphi \\ v_l \sin \varphi \\ 0 \end{bmatrix}}_{\mathbf{f}(\mathbf{x})} + \underbrace{\begin{bmatrix} 0 \\ 0 \\ v_l \end{bmatrix}}_g u \quad (1.63a)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) = \begin{bmatrix} x \\ y \end{bmatrix} \quad (1.63b)$$

betrachtet werden. Berechnet man wieder die Ableitungen von (1.50) in der Form

$$\dot{\hat{y}} = L_f \hat{h}(\mathbf{x}) + \underbrace{L_g \hat{h}(\mathbf{x})}_0 u = 2v_l \left(\frac{(x - y_{10}) \cos \varphi}{a^2} + \frac{(y - y_{20}) \sin \varphi}{b^2} \right) \quad (1.64a)$$

$$\ddot{\hat{y}} = L_f^2 \hat{h}(\mathbf{x}) + L_g L_f \hat{h}(\mathbf{x}) u = \quad (1.64b)$$

$$= 2v_l^2 \left(\frac{\cos^2 \varphi}{a^2} + \frac{\sin^2 \varphi}{b^2} \right) + \underbrace{\left(2v_l^2 \left(\frac{(y - y_{20}) \cos \varphi}{b^2} - \frac{(x - y_{10}) \sin \varphi}{a^2} \right) \right)}_{d(\mathbf{x})} u \quad (1.64c)$$

sieht man, dass

1. die erste Ableitung von \hat{y} ident zu (1.51a) ist und
2. der Vorfaktor von u durch den ersten Eintrag der Entkopplungsmatrix $\mathbf{D}(\mathbf{x})$ in (1.51c) gegeben ist.

Damit bleibt Γ^* qualitativ unverändert und der Ausgang \hat{y} besitzt wieder einen wohldefinierten relativen Grad in Γ^* .

Damit kann wieder direkt Satz 1.1 angewandt werden, dessen Aussagen auf Grund von $m = p - 1 = 1$ zu den Ergebnissen in [1.0] ident sind. Die Transformation in die transversale Normalform kann wieder zu

$$\Phi(\mathbf{x}) = \begin{bmatrix} \hat{h}(\mathbf{x}) \\ L_f \hat{h}(\mathbf{x}) \\ \varphi \end{bmatrix} \quad (1.65)$$

gewählt werden. Auf Grund von $m = p - 1$ ist kein v^{\parallel} vorhanden. Daher wird auch die Matrix \mathbf{N} zur leeren Matrix und \mathbf{M} reduziert sich zu $m(\mathbf{x}) = \frac{d(\mathbf{x})}{d(\mathbf{x})d(\mathbf{x})} = \frac{1}{d(\mathbf{x})} = \lambda(\mathbf{x})$.

Damit folgt die exakt linearisierende Zustandsrückführung zu

$$u = \alpha(\mathbf{x}) + \lambda(\mathbf{x})v^\dagger = \frac{2v_l^2(a^2(\cos\varphi)^2 - (\cos\varphi)^2b^2 - a^2) + a^2b^2v^\dagger}{2v_l^2(-\sin(\varphi)b^2x + \sin(\varphi)b^2y_{10} + \cos(\varphi)a^2y - \cos(\varphi)a^2y_{20})}. \quad (1.66)$$

Die transversale Normalform lautet

$$\dot{\xi} = \begin{bmatrix} \xi_2 \\ v^\dagger \end{bmatrix} \quad (1.67a)$$

$$\dot{\eta} = f^0(\eta, \xi) + g^\dagger(\eta, \xi)v^\dagger \quad (1.67b)$$

wobei die tangentiale Dynamik (1.67b) nicht beeinflussbar ist. Damit können in diesem Fall keine Ziele (Z3) berücksichtigt werden. Die Trajektorien des Fahrzeugs im geschlossenen Kreis sehen qualitativ sehr ähnlich zu denen von Beispiel 1.4 aus.

Example 1.6. Wieder soll das Fahrzeug mit konstanter Vorwärtsgeschwindigkeit v_l in der Form (1.63) betrachtet werden. Diesmal soll ein Pfadfolgeregler für die Verfolgung einer Lemniskate gegeben durch die implizite Gleichung

$$s_l(\mathbf{y}) = y_1^4 - y_1^2 + y_2^2 \quad (1.68)$$

entworfen werden. Dazu wird wieder die fiktive Ausgangsgröße

$$\hat{y} = s_l(\mathbf{h}(\mathbf{x})) = x^4 - x^2 + y^2 = \hat{h}(\mathbf{x}) \quad (1.69)$$

angeschrieben und festgestellt, ob diese einen relativen Grad in Γ^* aufweist. Dazu werden die Ableitungen von (1.69)

$$\dot{\hat{y}} = \mathbf{L}_f \hat{h}(\mathbf{x}) + \underbrace{\mathbf{L}_g \hat{h}(\mathbf{x}) u}_0 = (4x^3 - 2x)v_l \cos(\varphi) + 2yv_l \sin(\varphi) \quad (1.70a)$$

$$\ddot{\hat{y}} = \mathbf{L}_f^2 \hat{h}(\mathbf{x}) + \mathbf{L}_g \mathbf{L}_f \hat{h}(\mathbf{x}) u = \quad (1.70b)$$

$$= (12x^2 - 2)v_l^2(\cos\varphi)^2 + 2v_l^2(\sin\varphi)^2 + \underbrace{2v_l^2((x - 2x^3)\sin(\varphi) + y\cos(\varphi))u}_{d(\mathbf{x})} \quad (1.70c)$$

berechnet. Wie man unmittelbar erkennt, verschwindet der Vorfaktor $d(\mathbf{x})$ von u für beliebige $\begin{bmatrix} 0 & 0 & \bar{\varphi} \end{bmatrix}^T \in \Gamma^*$ womit die Voraussetzung verletzt ist, dass das System mit der Ausgangsgröße \hat{y} einen relativen Grad in Γ^* aufweist. Mit dem resultierenden Regelgesetz könnte nicht der gesamte Pfad abgefahren werden. Wie im folgenden Abschnitt 1.3.3 gezeigt wird, kann basierend auf der parametrierten Darstellung der Lemniskate aber trotzdem ein Pfadfolgeregler entworfen werden, der es ermöglicht, die gesamte Lemniskate abzufahren.

1.3.3 Parametrierte Pfade

Das Ziel dieses Abschnitts ist der Entwurf von Pfadfolgereglern für parametrierte Pfade definiert durch eine Abbildung $\sigma(\theta) : \mathcal{T} \rightarrow \mathbb{R}^p$. Dazu sollen auch hier ähnliche Annahmen wie in Abschnitt 1.3.2 getroffen werden, d. h.

1. für das System $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$ gilt $m \geq p - 1$, $p > 1$, und
2. das System $\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u}$ mit der Ausgangsgröße $\mathbf{y} = \mathbf{h}(\mathbf{x})$ hat in Γ^* einen wohldefinierten vektoriellen relativen Grad $\{r_1, r_2, \dots, r_p\}$. Da $m \geq p - 1$ erlaubt ist, muss im Folgenden eine Unterscheidung getroffen werden.
 - a) Falls $m \geq p$ gilt, kann direkt die Definition 1.2 des vektoriellen relativen Grades verwendet werden.
 - b) Für den Fall $m = p - 1$ bleibt diese Definition unverändert bis auf die Tatsache, dass $\text{rang}(\mathbf{D}(\bar{\mathbf{x}})) = m$ gelten muss.

Um die zeitliche Evolution des Pfadparameters θ beschreiben zu können, wird ein *Hilfssystem* der Form

$$\theta^{(\hat{r})} = v_\theta \quad (1.71)$$

mit einer neuen *fiktiven Eingangsgröße* v_θ und $\hat{r} = \max\{r_1, r_2, \dots, r_p\}$ definiert. Im Folgenden werden der Pfadparameter und seine zeitlichen Ableitungen zu einem Vektor $\zeta = [\theta \ \dot{\theta} \ \ddot{\theta} \ \dots \ \theta^{(\hat{r}-1)}]^T$ zusammengefasst und die Darstellung von (1.71) als ein lineares, zeitinvariantes System

$$\dot{\zeta} = \mathbf{F}\zeta + \mathbf{h}v_\theta \quad (1.72a)$$

mit

$$\mathbf{F} = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \ddots & \\ 0 & 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & 0 & \cdots & 0 \end{bmatrix} \in \mathbb{R}^{\hat{r} \times \hat{r}} \quad \text{und} \quad \mathbf{h} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \in \mathbb{R}^{\hat{r} \times 1} \quad (1.72b)$$

verwendet.

Für den Fall $m = p - 1$ werden die oben getroffenen Annahmen noch erweitert. Dazu definiert man die Indexmenge $\mathcal{I} = \{i \mid r_i = \hat{r}\}$. Basierend auf dieser Indexmenge kann der Vektor $\mathbf{d}_\theta \in \mathbb{R}^p$ aufgebaut werden, für den gilt, dass

$$d_{\theta,i} = \begin{cases} -\sigma'_i(\theta) & i \in \mathcal{I} \\ 0 & i \notin \mathcal{I} \end{cases} \quad (1.73)$$

Es soll nun nicht nur gelten, dass die Entkopplungsmatrix \mathbf{D} zum Ausgang \mathbf{y} in Γ^* den Rang $m = p - 1$ aufweist, sondern auch, dass die Matrix

$$[\mathbf{D} \ \mathbf{d}_\theta] \in \mathbb{R}^{(p \times p)} \quad (1.74)$$

den Rang p besitzt.

Für den Entwurf eines Pfadfolgereglers wird das erweiterte System

$$\dot{\chi} = \begin{bmatrix} \dot{\mathbf{x}} \\ \dot{\zeta} \end{bmatrix} = \begin{bmatrix} \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \\ \mathbf{F}\zeta + \mathbf{h}v_\theta \end{bmatrix} = \mathbf{f}_e(\chi) + \mathbf{G}_e(\chi)\boldsymbol{\nu} \quad (1.75a)$$

mit $\boldsymbol{\nu}^T = [\mathbf{u}^T \ v_\theta]$ und dem Fehlerausgang

$$\mathbf{e} = \mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}(\zeta_1) = \mathbf{h}_e(\chi) \quad (1.75b)$$

betrachtet.

Lemma 1.1. Das System (1.75) weist unter den getroffenen Annahmen einen wohldefinierten vektoriellen relativen Grad $\{r_1, r_2, \dots, r_p\}$ auf.

Exercise 1.5. Zeigen Sie Lemma 1.1.

Für den Fall $m \geq p$ kann unter den getroffenen Annahmen stets eine spezielle Normalform gefunden werden, die den Entwurf von Pfadfolgeregbern für parametrierte Pfade wesentlich vereinfacht. Dieses Ergebnis ist im nachfolgenden Satz formuliert und beruht auf Lemma 1.1 und Satz 1.1.

Theorem 1.2. Es sei $m \geq p$ und alle getroffenen Annahmen seien erfüllt. Dann gilt, dass das System (1.75) äquivalent unter statischer Rückführung zu einem System der Form

$$\dot{\xi} = \mathbf{A}\xi + \mathbf{B}\mathbf{v}_P \quad (1.76a)$$

$$\dot{\eta} = \mathbf{f}^0(\eta, \xi, \eta_\theta) + \mathbf{G}^P(\eta, \xi, \eta_\theta)\mathbf{v}_P + \mathbf{G}^R(\eta, \xi, \eta_\theta)\mathbf{v}_R \quad (1.76b)$$

$$\dot{\eta}_\theta = \mathbf{F}\eta_\theta + [\mathbf{0} \ \mathbf{h}]\mathbf{v}_R \quad (1.76c)$$

ist. Dabei gilt $\xi \in \mathbb{R}^r$, $\eta \in \mathbb{R}^{n-r}$, $\eta_\theta \in \mathbb{R}^{\hat{r}}$, $\mathbf{v}_P \in \mathbb{R}^p$ und $\mathbf{v}_R \in \mathbb{R}^{m-p+1}$ mit $r = \sum_{j=1}^p r_j$. Das Teilsystem (1.76a) liegt in Brunovsky Normalform vor.

Proof. Die Herleitung beruht auf Satz 1.1 und liefert gleichzeitig die Vorschrift, wie die Transformation in die Normalform (1.76) zu errechnen ist. Es wird die spezielle Struktur des Systems (1.75) genutzt, nämlich dass die Dynamiken des Pfadparameters und des ursprünglichen Systems nur über den Fehlerausgang \mathbf{e} verknüpft sind. Um die Normalform (1.76) zu erhalten, wird dem System (1.75) eine weitere skalare Ausgangsgröße

$$y_\theta = \zeta_1 = h_\theta(\chi) \quad (1.77)$$

hinzugefügt. Es gilt, dass basierend auf Lemma 1.1 auch das System

$$\dot{\chi} = \mathbf{f}_e(\chi) + \mathbf{G}_e(\chi)\boldsymbol{\nu} \quad (1.78a)$$

$$\mathbf{e} = \mathbf{h}_e(\chi) \quad (1.78b)$$

$$y_\theta = h_\theta(\chi) \quad (1.78c)$$

einen wohldefinierten vektoriellen relativen Grad $\{r_1, \dots, r_p, \hat{r}\}$ besitzt.

Exercise 1.6. Zeigen Sie diese Aussage.

Damit kann direkt Satz 1.1 auf das System (1.78) angewandt werden. Der erste Teil der Zustandstransformation in die Normalform (1.76) lautet

$$\xi = \Phi_\xi = \begin{bmatrix} h_{e,1}(\chi) \\ L_{\mathbf{f}_e} h_{e,1}(\chi) \\ \vdots \\ L_{\mathbf{f}_e}^{r_1-1} h_{e,1}(\chi) \\ h_{e,2}(\chi) \\ \vdots \\ L_{\mathbf{f}_e}^{r_2-1} h_{e,2}(\chi) \\ \vdots \\ h_{e,p}(\chi) \\ \vdots \\ L_{\mathbf{f}_e}^{r_p-1} h_{e,p}(\chi) \end{bmatrix} \quad (1.79a)$$

$$\eta_\theta = \Phi_{\eta_\theta} = \begin{bmatrix} h_\theta(\chi) \\ L_{\mathbf{f}_e} h_\theta(\chi) \\ \vdots \\ L_{\mathbf{f}_e}^{\hat{r}-1} h_\theta(\chi) \end{bmatrix} = \zeta, \quad (1.79b)$$

die mit den Funktionen

$$\eta = \Phi_\eta = \begin{bmatrix} \phi_{\eta,1}(\chi) \\ \vdots \\ \phi_{\eta,n-r}(\chi) \end{bmatrix} \quad (1.79c)$$

zu einem Diffeomorphismus komplettiert wird. Die zeitliche Ableitung von (1.79) kann gemäß Satz 1.1 in der Form

$$\frac{d}{dt} \begin{bmatrix} \xi \\ \eta_\theta \end{bmatrix} = \bar{\mathbf{A}} \begin{bmatrix} \xi \\ \eta_\theta \end{bmatrix} + \bar{\mathbf{B}} \begin{bmatrix} \bar{\mathbf{v}}_1 \\ v_\theta \end{bmatrix} \quad (1.80a)$$

$$\dot{\eta} = \mathbf{f}^0(\eta, \xi, \eta_\theta) + \mathbf{G}_1(\eta, \xi, \eta_\theta) \begin{bmatrix} \bar{\mathbf{v}}_1 \\ v_\theta \end{bmatrix} + \mathbf{G}_2(\eta, \xi, \eta_\theta) \mathbf{v}_2 \quad (1.80b)$$

mit $\bar{\mathbf{A}}$ und $\bar{\mathbf{B}}$ in Brunovsky Normalform dargestellt werden. Durch Setzen von $\mathbf{v}_P = \bar{\mathbf{v}}_1$ und $\mathbf{v}_R = [\mathbf{v}_2^T \ v_\theta]^T$ folgt die gewünschte Normalform (1.76). Die Zustand-

srückführung stellt sich in der Form

$$\boldsymbol{\nu} = \boldsymbol{\alpha}(\boldsymbol{\chi}) + \boldsymbol{\Lambda}(\boldsymbol{\chi}) \begin{bmatrix} \bar{\mathbf{v}}_1 \\ v_\theta \\ \mathbf{v}_2 \end{bmatrix} = \boldsymbol{\alpha}(\boldsymbol{\chi}) + \bar{\boldsymbol{\Lambda}}(\boldsymbol{\chi}) \begin{bmatrix} \mathbf{v}_{\mathcal{P}} \\ \mathbf{v}_{\mathcal{R}} \end{bmatrix} \quad (1.81)$$

dar, wobei die letzte Zeile in (1.81) $v_\theta = v_\theta$ lautet. \square

Im Fall $m = p - 1$ kann ebenfalls eine Normalform analog zu (1.76) erhalten werden. Dazu wendet man basierend auf Lemma 1.1 Satz 1.1 direkt auf das System (1.75) an und erhält

$$\dot{\boldsymbol{\xi}} = \mathbf{A}\boldsymbol{\xi} + \mathbf{B}\mathbf{v}_{\mathcal{P}} \quad (1.82a)$$

$$\dot{\boldsymbol{\eta}} = \mathbf{f}^0(\boldsymbol{\eta}, \boldsymbol{\xi}) + \mathbf{G}^{\mathcal{P}}(\boldsymbol{\eta}, \boldsymbol{\xi})\mathbf{v}_{\mathcal{P}} \quad (1.82b)$$

mit $\boldsymbol{\xi} \in \mathbb{R}^r$, $\boldsymbol{\eta} \in \mathbb{R}^{n+\hat{r}-r}$ und $\mathbf{v}_{\mathcal{P}} \in \mathbb{R}^p$. Die Normalform (1.82) weist im Gegensatz zu (1.76) weniger Struktur auf. Die interne Dynamik mit den Koordinaten $\boldsymbol{\eta}$ kann im Allgemeinen nicht mehr in zwei Teile aufgespalten werden, von denen einer direkt dem Hilfssystem entspricht. Beide Normalformen (1.76) und (1.82) können analog zu (1.39) als transversale Normalformen bezeichnet werden. Die Koordinaten $\boldsymbol{\xi}$ sind wieder der transversalen Dynamik zugeordnet. Dementsprechend repräsentieren die Koordinaten $\boldsymbol{\eta}$ die tangentiale Dynamik.

Im Fall $m \geq p$ erkennt man anhand von (1.76), dass sich die Trajektorie des Hilfssystems (und damit des Pfadparameters θ) stets gezielt vorgeben lässt. Da für $\boldsymbol{\xi} = \mathbf{0}$ unabhängig vom Wert von $\boldsymbol{\zeta}$ die Ausgangsgrößen \mathbf{y} exakt am Pfad liegen, ist es sinnvoll die Dynamik des Pfadparameters der tangentialen Dynamik zuzuordnen.

Analog zu Abschnitt 1.3.2 kann auch hier die Menge Γ^* in der Form

$$\Gamma^* = \left\{ \bar{\boldsymbol{\chi}} \in \mathbb{R}^{n+\hat{r}} \mid \Phi_{\boldsymbol{\xi}}(\bar{\boldsymbol{\chi}}) = \mathbf{0} \right\} \subset \Gamma \quad (1.83)$$

definiert werden. Die Ziele (Z1)–(Z3) können ebenfalls so wie in Abschnitt 1.3.2 durch den Entwurf von Reglern für die transversale und (falls möglich) tangentiale Dynamik erreicht werden.

Bei der Implementierung des resultierenden Pfadfolgereglers muss die Dynamik (1.71) des Pfadparameters im Digitalrechner mit integriert werden. Dafür sind auch geeignete Anfangswerte $\boldsymbol{\zeta}(t_0)$ erforderlich. In der Literatur [1.0] wird häufig empfohlen, $\dot{\theta}(t_0) = \ddot{\theta}(t_0) = \dots = \theta^{(\hat{r}-1)}(t_0) = 0$ und

$$\theta(t_0) = \arg \min_{\theta \in \mathcal{T}} \|\mathbf{h}(\mathbf{x}(t_0)) - \boldsymbol{\sigma}(\theta)\| \quad (1.84)$$

zu wählen. Man beachte aber, dass im Allgemeinen die Wahl von $\boldsymbol{\zeta}(t_0)$ maßgeblich beeinflusst, ob die Invarianzeigenschaft (Z2) erfüllt ist, siehe dazu auch das Beispiel 1.7.

Example 1.7 (Fahrzeug mit variabler Vorwärtsgeschwindigkeit). Als Beispiel dient zunächst wieder das Fahrzeug mit variabler Vorwärtsgeschwindigkeit mit dem mathematischen Modell (1.48). Für dieses gilt $m = p = 2$. Es soll ein Pfadfolgeregler für

die Verfolgung einer Lemniskate gegeben durch die parametrierte Darstellung

$$\boldsymbol{\sigma}_l(\theta) = \begin{bmatrix} \cos(\theta) \\ \cos(\theta) \sin(\theta) \end{bmatrix} \quad (1.85)$$

entworfen werden. Das Ziel (Z3) wird im konkreten Fall so spezifiziert, dass die Lemniskate alle 20 Sekunden komplett abgefahren werden soll. Zunächst wird der vektorielle relative Grad des Systems (1.48) ermittelt. Durch Bilden der Ableitungen

$$\dot{\mathbf{y}} = \begin{bmatrix} v_l \cos(\varphi) \\ v_l \sin(\varphi) \end{bmatrix} \quad (1.86a)$$

$$\ddot{\mathbf{y}} = \underbrace{\begin{bmatrix} -v_l^2 \sin(\varphi) & \cos(\varphi) \\ v_l^2 \cos(\varphi) & \sin(\varphi) \end{bmatrix}}_{\mathbf{D}(\mathbf{x})} \mathbf{u} \quad (1.86b)$$

erhält man einen vektoriellen relativen Grad von $\{2, 2\}$. Die Entkopplungsmatrix ist für alle Punkte im Zustandsraum mit $v_l \neq 0$ regulär. Damit folgen $\hat{r} = 2$ und das Hilfssystem für den Pfadparameter in der Form

$$\dot{\boldsymbol{\zeta}} = \underbrace{\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}}_{\mathbf{F}} \boldsymbol{\zeta} + \underbrace{\begin{bmatrix} 0 \\ 1 \end{bmatrix}}_{\mathbf{h}} v_\theta \quad (1.87)$$

mit $\boldsymbol{\zeta} = [\theta \quad \dot{\theta}]^T$.

Grundlage des Entwurfes ist das erweiterte System

$$\begin{aligned} \dot{\boldsymbol{\chi}} &= \mathbf{f}_e(\boldsymbol{\chi}) + \mathbf{G}_e(\boldsymbol{\chi})\boldsymbol{\nu} = \mathbf{f}_e(\boldsymbol{\chi}) + [\mathbf{g}_{e,1}(\boldsymbol{\chi}) \quad \mathbf{g}_{e,2}(\boldsymbol{\chi}) \quad \mathbf{g}_{e,3}(\boldsymbol{\chi})]\boldsymbol{\nu} \\ &= \begin{bmatrix} v_l \cos \varphi \\ v_l \sin \varphi \\ 0 \\ 0 \\ \zeta_2 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ v_l & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ v_\theta \end{bmatrix} \end{aligned} \quad (1.88a)$$

mit den Ausgangsgrößen

$$\mathbf{e} = \mathbf{h}_e(\boldsymbol{\chi}) = \mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}_l(\zeta_1) = \begin{bmatrix} x - \cos(\zeta_1) \\ y - \cos(\zeta_1) \sin(\zeta_1) \end{bmatrix} \quad (1.88b)$$

$$y_\theta = h_\theta(\boldsymbol{\chi}) = \zeta_1 . \quad (1.88c)$$

Da gleich viele Eingangs- wie Ausgangsgrößen vorliegen und für den vektoriellen relativen Grad $r_1 + r_2 + \hat{r} = 6$ gilt, reduziert sich die weitere Vorgehensweise auf eine Eingangs-Zustandslinearisierung gemäß [1.0]. Die Koordinaten der transversalen Normalform ergeben sich gemäß (1.79) zu

$$\xi = \begin{bmatrix} x - \cos(\zeta_1) \\ v_l \cos(\varphi) + \sin(\zeta_1)\zeta_2 \\ y - \cos(\zeta_1)\sin(\zeta_1) \\ v_l \sin(\varphi) + ((\sin \zeta_1)^2 - (\cos \zeta_1)^2)\zeta_2 \end{bmatrix} \quad \boldsymbol{\eta}_\theta = \begin{bmatrix} \zeta_1 \\ \zeta_2 \end{bmatrix}. \quad (1.89)$$

Die exakt linearisierende Zustandsrückführung lautet

$$\begin{aligned} \boldsymbol{\nu} &= \begin{bmatrix} L_{\mathbf{g}_{e,1}} L_{\mathbf{f}_e} h_{e,1}(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,2}} L_{\mathbf{f}_e} h_{e,1}(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,3}} L_{\mathbf{f}_e} h_{e,1}(\boldsymbol{\chi}) \\ L_{\mathbf{g}_{e,1}} L_{\mathbf{f}_e} h_{e,2}(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,2}} L_{\mathbf{f}_e} h_{e,2}(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,3}} L_{\mathbf{f}_e} h_{e,2}(\boldsymbol{\chi}) \\ L_{\mathbf{g}_{e,1}} L_{\mathbf{f}_e} h_\theta(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,2}} L_{\mathbf{f}_e} h_\theta(\boldsymbol{\chi}) & L_{\mathbf{g}_{e,3}} L_{\mathbf{f}_e} h_\theta(\boldsymbol{\chi}) \end{bmatrix}^{-1} \left(\begin{bmatrix} \mathbf{v}_P \\ v_\theta \end{bmatrix} - \begin{bmatrix} L_{\mathbf{f}_e}^2 h_{e,1}(\boldsymbol{\chi}) \\ L_{\mathbf{f}_e}^2 h_{e,2}(\boldsymbol{\chi}) \\ L_{\mathbf{f}_e}^2 h_\theta(\boldsymbol{\chi}) \end{bmatrix} \right) \\ &= \begin{bmatrix} -\frac{\sin \varphi}{v_l^2} & \frac{\cos \varphi}{v_l^2} & \frac{\cos(\varphi)(2(\cos \zeta_1)^2 - 1) + \sin(\varphi) \sin(\zeta_1)}{v_l^2} \\ \cos \varphi & \sin \varphi & \frac{\sin(\varphi)(2(\cos \zeta_1)^2 - 1) - \cos(\varphi) \sin(\zeta_1)}{1} \\ 0 & 0 & 1 \end{bmatrix} \left(\begin{bmatrix} \mathbf{v}_P \\ v_\theta \end{bmatrix} - \begin{bmatrix} \cos(\zeta_1)\zeta_2^2 \\ 4\cos(\zeta_1)\sin(\zeta_1)\zeta_2^2 \\ 0 \end{bmatrix} \right). \end{aligned} \quad (1.90)$$

Diese transformiert gemeinsam mit (1.89) das System (1.88) in die Normalform (1.76), die sich im vorliegenden Fall zu

$$\dot{\xi} = \begin{bmatrix} \xi_2 \\ v_{P,1} \\ \xi_4 \\ v_{P,2} \end{bmatrix} \quad \dot{\boldsymbol{\eta}}_\theta = \begin{bmatrix} \eta_{\theta,2} \\ v_\theta \end{bmatrix} \quad (1.91)$$

ergibt. Da $r_1 + r_2 = r = 4 = n$ gilt, gibt es keinen Zustand $\boldsymbol{\eta}$ bzw. keine Dynamik (1.76b). Durch die Verwendung der Regelgesetze

$$v_{P,1} = -k_{P,11}\xi_1 - k_{P,12}\xi_2 \quad (1.92a)$$

$$v_{P,2} = -k_{P,21}\xi_3 - k_{P,22}\xi_4 \quad (1.92b)$$

$$v_\theta = -k_\theta(\eta_{\theta,2} - \dot{\theta}_d) \quad (1.92c)$$

mit $k_{P,ij}, k_\theta > 0$ kann der Pfad stabilisiert und gleichzeitig mit $\dot{\theta}_d = \pm \frac{2\pi}{20}$ das Ziel (Z3) erreicht werden.

Für die nachfolgenden Ergebnisse werden die Parameterwerte $k_{P,11} = k_{P,21} = 0.25$, $k_{P,12} = k_{P,22} = 1$ und $k_\theta = 0.5$ verwendet. Am Anfangszeitpunkt $t_0 = 0$ wird $\zeta_1(0)$ aus (1.84) bestimmt. Damit die Invarianzeigenschaft erfüllt ist, muss zusätzlich nicht

nur das Fahrzeug auf dem Pfad und tangential dazu starten, sondern es muss auch $\zeta_2(0)$ entsprechend (1.89) so gewählt werden, dass die Gleichungen

$$v_l(0) \cos(\varphi(0)) + \sin(\zeta_1(0))\zeta_2(0) = 0 \quad (1.93a)$$

$$v_l(0) \sin(\varphi(0)) + ((\sin \zeta_1(0))^2 - (\cos \zeta_1(0))^2)\zeta_2(0) = 0 \quad (1.93b)$$

erfüllt sind. Da bei den nachfolgend gezeigten Simulationsergebnissen das Fahrzeug ohnehin nicht am Pfad startet, wird, sofern nicht anders angegeben, $\zeta_2(0) = 0$ gewählt.

In Abbildung 1.7 ist die Trajektorie des Fahrzeuges für den Anfangswert $\mathbf{x}_0 = [5.5 \quad 1 \quad \frac{3\pi}{2} \quad 0.45]^T$ dargestellt.

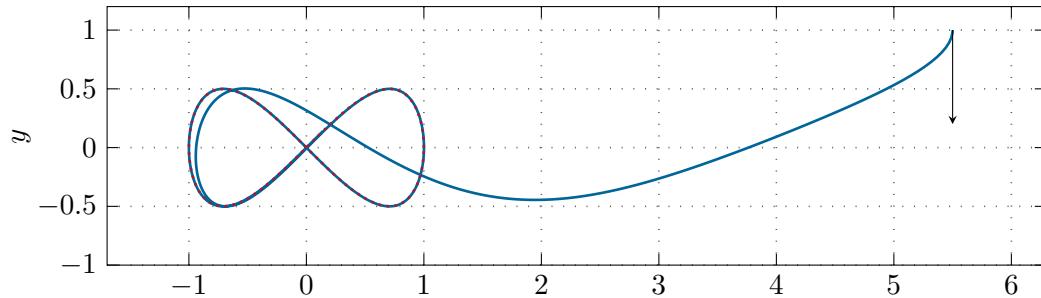


Figure 1.7: Trajektorie des Fahrzeuges mit $v_l \neq \text{konst.}$ für die Pfadfolgeregelung einer Lemniskate.

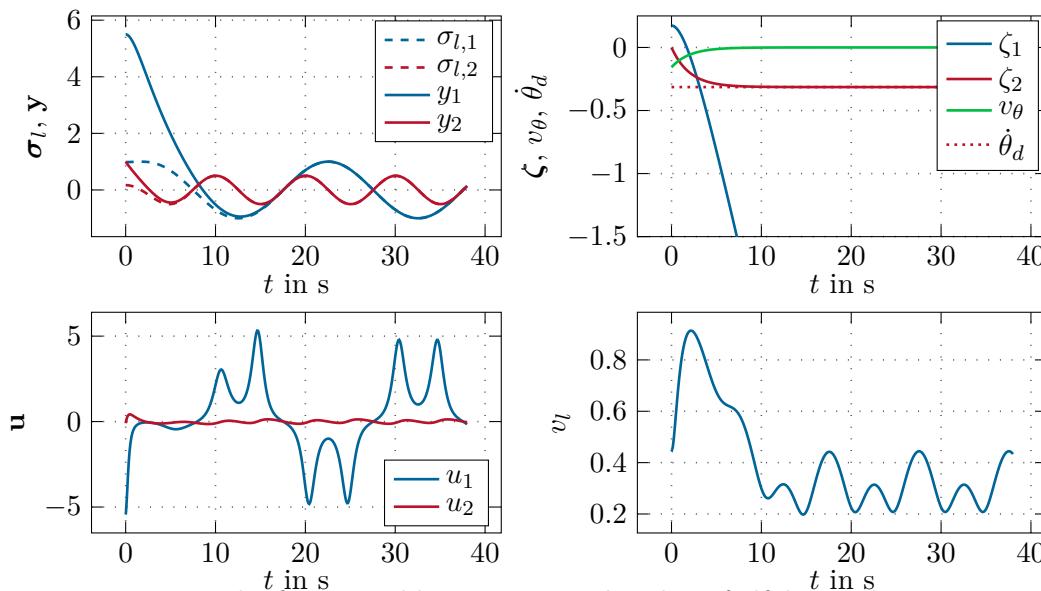


Figure 1.8: Zeitverläufe im geschlossenen Kreis bei der Pfadfolgeregelung einer Lemniskate für das Fahrzeug mit $v_l \neq \text{konst.}$

Die zugehörigen Zeitverläufe im geschlossenen Kreis sind in Abbildung 1.8 ersichtlich. Der Fehlerausgang wird zu null geregelt was durch die Konvergenz von \mathbf{y} zu σ_l ersichtlich ist. Die zeitliche Ableitung des Pfadparameters $\dot{\theta} = \zeta_2$ konvergiert zum

gewünschten Wert $\dot{\theta}_d$. Dadurch wächst im Laufe der Pfadverfolgung der Wert des Pfadparameters $\theta = \zeta_1$ über alle Grenzen. Bedingt durch die Tatsache, dass $\dot{\theta}$ konstant gehalten wird, schwankt die Längsgeschwindigkeit v_l des Fahrzeuges auch bei exakter Verfolgung der Lemniskate, siehe dazu auch Beispiel 1.8.

Abbildung 1.9 zeigt Trajektorien des Fahrzeuges für verschiedene Anfangswerte.

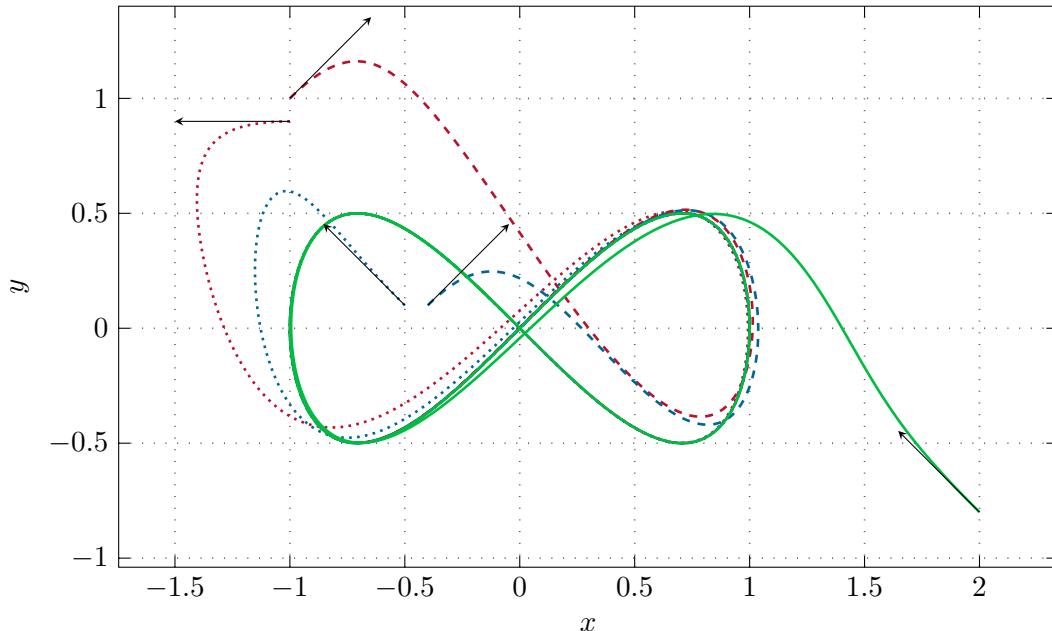


Figure 1.9: Trajektorien des Fahrzeuges mit $v_l \neq \text{konst.}$ für die Pfadfolgeregelung einer Lemniskate und verschiedenen Anfangswerten.

Example 1.8 (Fahrzeug mit konstanter Vorwärtsgeschwindigkeit). Zur Illustration des Falles $m = p - 1$ soll das Fahrzeug mit konstanter Vorwärtsgeschwindigkeit (1.63) herangezogen werden. Das Ziel ist der Entwurf eines Pfadfolgereglers für die Verfolgung einer parametrisierten Lemniskate. Für die Ableitungen der Ausgangsgröße \mathbf{y} gilt

$$\dot{\mathbf{y}} = \begin{bmatrix} v_l \cos \varphi \\ v_l \sin \varphi \end{bmatrix} \quad (1.94a)$$

$$\ddot{\mathbf{y}} = \underbrace{\begin{bmatrix} -v_l^2 \sin \varphi \\ v_l^2 \cos \varphi \end{bmatrix}}_{\mathbf{d}(\mathbf{x})} u . \quad (1.94b)$$

Für $v_l \neq 0$ und beliebige Werte φ gilt, dass nicht beide Komponenten von \mathbf{d} gleichzeitig verschwinden. Damit gilt $\text{rang}(\mathbf{d}) = m = 1$ und $r_1 = r_2 = 2$. Daraus folgt wie in Beispiel 1.7 $\hat{r} = 2$. Da der Fall $m = p - 1$ vorliegt, müssen noch die erweiterten Annahmen aus Abschnitt 1.3.3 überprüft werden. Mit der Indexmenge $\mathcal{I} = \{1, 2\}$

gilt $\mathbf{d}_\theta = -\boldsymbol{\sigma}'_l(\theta)$. Es ist gefordert, dass die erweiterte Matrix

$$\begin{bmatrix} -v_l^2 \sin \varphi & \sin \theta \\ v_l^2 \cos \varphi & (\sin \theta)^2 - (\cos \theta)^2 \end{bmatrix} \quad (1.95)$$

in der Menge Γ^* Rang $p = 2$ besitzt. In Γ^* gilt, dass sich das Fahrzeug tangential zur Lemniskate bewegen muss, d. h.

$$\varphi = \arctan\left(\frac{\sigma'_{l,2}(\theta)}{\sigma'_{l,1}(\theta)}\right) = \arctan\left(\frac{-(\sin \theta)^2 + (\cos \theta)^2}{-\sin \theta}\right). \quad (1.96)$$

Setzt man diese Beziehung in (1.95) ein, ergibt sich die Determinante der so entstandenen Matrix zu

$$v_l^2 \sqrt{4(\cos \theta)^4 - 5(\cos \theta)^2 + 2}. \quad (1.97)$$

Diese ist für alle $\theta \in \mathbb{R}$ ungleich null, womit die Regularität der erweiterten Matrix (1.95) entlang des gesamten Pfades nachgewiesen ist.

Das erweiterte System für den Reglerentwurf lautet

$$\dot{\chi} = \mathbf{f}_e(\chi) + \mathbf{G}_e(\chi)\boldsymbol{\nu} = \begin{bmatrix} v_l \cos \varphi \\ v_l \sin \varphi \\ 0 \\ \zeta_2 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ v_l & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u \\ v_\theta \end{bmatrix} \quad (1.98a)$$

mit den Ausgangsgrößen

$$\mathbf{e} = \mathbf{h}_e(\chi) = \mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}_l(\zeta_1) = \begin{bmatrix} x - \cos(\zeta_1) \\ y - \cos(\zeta_1) \sin(\zeta_1) \end{bmatrix}. \quad (1.98b)$$

Um den Fehlerausgang \mathbf{e} auf null zu regeln, wird für (1.98) in weiterer Folge exakte Eingangs-Ausgangslinearisierung vorgenommen. Es gilt

$$\ddot{\mathbf{e}} = \begin{bmatrix} \cos(\zeta_1) \zeta_2^2 \\ 4 \cos(\zeta_1) \sin(\zeta_1) \zeta_2^2 \end{bmatrix} + \begin{bmatrix} -v_l^2 \sin \varphi & \sin \zeta_1 \\ v_l^2 \cos \varphi & (\sin \zeta_1)^2 - (\cos \zeta_1)^2 \end{bmatrix} \boldsymbol{\nu}, \quad (1.99a)$$

wobei die Entkopplungsmatrix natürlich mit der erweiterten Matrix (1.95) übereinstimmt. Deren Regularität wurde bereits nachgewiesen womit die exakt linearisierende Rückführung

$$\boldsymbol{\nu} = \begin{bmatrix} -v_l^2 \sin \varphi & \sin \zeta_1 \\ v_l^2 \cos \varphi & (\sin \zeta_1)^2 - (\cos \zeta_1)^2 \end{bmatrix}^{-1} \left(\mathbf{v}_P - \begin{bmatrix} \cos(\zeta_1) \zeta_2^2 \\ 4 \cos(\zeta_1) \sin(\zeta_1) \zeta_2^2 \end{bmatrix} \right) \quad (1.100)$$

berechnet werden kann. Mit der Wahl $\eta = \zeta_1$ ergibt sich die Zustandstransformation

$$\begin{bmatrix} \boldsymbol{\xi} \\ \eta \end{bmatrix} = \Phi(\boldsymbol{\chi}) = \begin{bmatrix} x - \cos(\zeta_1) \\ v_l \cos(\varphi) + \sin(\zeta_1)\zeta_2 \\ y - \cos(\zeta_1)\sin(\zeta_1) \\ v_l \sin(\varphi) + ((\sin \zeta_1)^2 - (\cos \zeta_1)^2)\zeta_2 \\ \zeta_1 \end{bmatrix} \quad (1.101)$$

und die transversale Normalform (1.82) zu

$$\dot{\xi}_1 = \xi_2 \quad (1.102a)$$

$$\dot{\xi}_2 = v_{P,1} \quad (1.102b)$$

$$\dot{\xi}_3 = \xi_4 \quad (1.102c)$$

$$\dot{\xi}_4 = v_{P,2} \quad (1.102d)$$

$$\dot{\eta} = f^0(\eta, \boldsymbol{\xi}) . \quad (1.102e)$$

Die interne Dynamik am Pfad lautet

$$\dot{\eta} = f^0(\eta, \mathbf{0}) = \pm \frac{|v_l|}{\sqrt{4(\cos \eta)^4 - 5(\cos \eta)^2 + 2}} . \quad (1.103)$$

Aus (1.103) sieht man unmittelbar, dass am Pfad für $v_l = \text{konst.}$ die zeitliche Ableitung des Pfadparameters $\dot{\eta} = \dot{\theta}$ nicht konstant sein kann. Umgekehrt erklärt dies auch, wieso bei konstantem $\dot{\theta}$ die Längsgeschwindigkeit v_l des Fahrzeuges schwanken muss, wie es in Beispiel 1.7 ersichtlich ist.

Mit Reglern analog zu (1.92) kann der Pfad stabilisiert werden. Die Trajektorien des Fahrzeuges im geschlossenen Kreis sehen qualitativ sehr ähnlich zu denen von Beispiel 1.7 aus.

1.4 Modellprädiktive Pfadfolgeregelung

In diesem Abschnitt soll ein Pfadfolgeregler basierend auf der *modellprädiktiven Regelung (MPC)* entwickelt werden. Grundsätzlich gibt es sehr viele Möglichkeiten, wie MPC für die Pfadfolgeregelung eingesetzt werden kann. Welche Lösung im konkreten Fall gewählt wird, hängt unter anderem von der Systemstruktur ab und ob der Pfad parametriert oder implizit definiert ist.

In diesem Abschnitt wird die Formulierung des modellprädiktiven Reglers für den folgenden Fall vorgestellt. Die Systemstruktur sei durch

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}) \quad (1.104a)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) \quad (1.104b)$$

mit $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{u} \in \mathbb{R}^m$ und $\mathbf{y} \in \mathbb{R}^p$ gegeben. Der Einfachheit halber gelte $m = p$. Weiters

sollen parametrierte Pfade \mathcal{P} der Form

$$\mathcal{P} = \{\bar{\mathbf{y}} \in \mathbb{R}^p \mid \bar{\mathbf{y}} = \boldsymbol{\sigma}(\theta), \theta \in \mathcal{T}\} \quad (1.105)$$

verfolgt werden. Die *Stellgrößen* unterliegen *Beschränkungen*, die durch

$$\mathbf{u}_{\min} \leq \mathbf{u} \leq \mathbf{u}_{\max} \quad (1.106)$$

gegeben sind. Alle Zustände \mathbf{x} seien direkt messbar bzw. bekannt. Bezuglich der Dynamik des Pfadparameters θ wird wieder angenommen, dass dieser über das Hilfssystem

$$\dot{\zeta} = \mathbf{F}\zeta + \mathbf{h}v_\theta \quad (1.107)$$

der Ordnung \hat{r} vollkommen analog zu (1.72) beschrieben wird. Das Ziel (Z3) der Bewegung auf dem Pfad sei so formulierbar, dass gewisse Komponenten ζ_i einen gewünschten Wert $\zeta_{d,i} = \text{konst.}$ annehmen sollen. Die Ordnung \hat{r} soll vorerst einfach passend zu dieser Anforderung gewählt sein. Es sei darauf hingewiesen, dass der Schwerpunkt in diesem Abschnitt auf der *Formulierung* des modellprädiktiven Reglers liegt. Es werden *keine Stabilitätsbetrachtungen* durchgeführt.

Als eine Möglichkeit, den Pfad zu stabilisieren, bietet sich folgende MPC-Formulierung an. In jedem Abtastschritt $t_k = kT_a$ wird das Optimalsteuerungsproblem

$$\min_{\mathbf{u}(\cdot), v_\theta(\cdot)} J(t_k, \mathbf{u}(\cdot), v_\theta(\cdot)) \quad (1.108a)$$

$$\text{u.B.v.} \quad \dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}) \quad \mathbf{x}(t_k) = \mathbf{x}_k \quad (1.108b)$$

$$\dot{\zeta} = \mathbf{F}\zeta + \mathbf{h}v_\theta \quad \zeta(t_k) = \zeta_k \quad (1.108c)$$

$$\mathbf{u}_{\min} \leq \mathbf{u}(t) \leq \mathbf{u}_{\max} \quad \forall t \in [t_k, t_k + T] \quad (1.108d)$$

$$\zeta_1(t) \in \mathcal{T} \quad \forall t \in [t_k, t_k + T] \quad (1.108e)$$

mit

$$J(t_k, \mathbf{u}(\cdot), v_\theta(\cdot)) = \int_{t_k}^{t_k+T} l(\mathbf{x}, \zeta, \mathbf{u}, v_\theta) dt + V(\mathbf{x}(t_k + T), \zeta(t_k + T)) \quad (1.108f)$$

und

$$l(\mathbf{x}, \zeta, \mathbf{u}, v_\theta) = \frac{1}{2} \left[\left(\mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}(\zeta_1) \right)_{\mathbf{Q}_e} + \left(\zeta - \zeta_d \right)_{\mathbf{Q}_\zeta} + \left(\mathbf{u} \right)_{\mathbf{R}_u} + r_\theta v_\theta^2 \right] \quad (1.108g)$$

$$V(\mathbf{x}, \zeta) = \frac{1}{2} \left[\left(\mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}(\zeta_1) \right)_{\mathbf{S}_e} + \left(\zeta - \zeta_d \right)_{\mathbf{S}_\zeta} \right] \quad (1.108h)$$

gelöst. Dabei ist $\left(\mathbf{x} \right)_\mathbf{Q}$ eine Kurzschreibweise für die quadratische Form $\left(\mathbf{x} \right)_\mathbf{Q} = \mathbf{x}^T \mathbf{Q} \mathbf{x}$. Die Matrizen \mathbf{Q}_e , \mathbf{R}_u und \mathbf{S}_e seien positiv definit. Da wie oben ausgeführt nur gewisse Komponenten ζ_i von ζ einen konstanten Sollwert $\zeta_{d,i}$ annehmen sollen, sind die Matrizen \mathbf{Q}_ζ und \mathbf{S}_ζ im Allgemeinen nur positiv semidefinit. Die restlichen Komponenten von ζ_d können beliebig festgesetzt werden. Weiters gelte $r_\theta > 0$. Der Optimierungshorizont ist durch $T \geq T_a$ gegeben. Wie schon in Kapitel ?? werden im Folgenden optimale Lösungen

eines Optimalsteuerungsproblems mit $*$ gekennzeichnet. Zusätzlich wird im Argument der entsprechenden Funktionen der Abtastzeitpunkt hinzugefügt, an dem sie ermittelt wurden. Dementsprechend lautet die optimale Lösung von (1.108) $\mathbf{u}^*(\cdot, t_k)$ und $v_\theta^*(\cdot, t_k)$ mit den zugehörigen Zustandstrajektorien $\mathbf{x}^*(\cdot, t_k)$ und $\zeta^*(\cdot, t_k)$. Der optimale Wert des Kostenfunktionalen J ergibt sich zu $J^*(t_k) = J(t_k, \mathbf{u}^*(\cdot, t_k), v_\theta^*(\cdot, t_k))$. Das Optimalsteuerungsproblem (1.108) wird mit den Anfangswerten \mathbf{x}_k und ζ_k gelöst. Der Anfangswert \mathbf{x}_k entspricht direkt dem Wert des Zustandes des Systems (1.104a) zum Zeitpunkt t_k . Für das Hilfssystem wird die optimale Lösung aus dem letzten Abtastschritt t_{k-1} , ausgewertet zum Zeitpunkt t_k , herangezogen, d. h. $\zeta_k = \zeta^*(t_k, t_{k-1})$. Der Anfangswert für ζ im ersten Zeitschritt t_0 kann so bestimmt werden wie am Ende von Abschnitt 1.3.3 beschrieben.

Intuitiv ist der modellprädiktive Regler basierend auf (1.108) durch die positiv definite Gewichtung des Fehlers $\mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}(\zeta_1)$ in der Lage, das Ziel (Z1) zu erfüllen. Darüber hinaus können durch die Gewichtung von $\zeta - \zeta_d$ die Ziele (Z3) erreicht werden. Das Problem bei der obigen MPC-Formulierung besteht darin, dass die Invarianzeigenschaft (Z2) im Allgemeinen nicht gegeben ist. Dies liegt an der absoluten Gewichtung der Stellgrößen \mathbf{u} . Eine Gewichtung der Stellgrößen ist für gewöhnlich sinnvoll, um unnötig große Amplituden der Stellgrößen zu vermeiden bzw. um ein gewünschtes Verhalten im geschlossenen Kreis einzustellen. Allerdings ist für eine exakte Verfolgung des Pfades, d. h. für $\mathbf{h}(\mathbf{x}) \equiv \boldsymbol{\sigma}(\zeta_1)$, üblicherweise eine Stellgröße $\mathbf{u} \neq \mathbf{0}$ erforderlich, die nicht konstant ist. Dies bewirkt, dass dadurch unvermeidliche Kosten in J entstehen, die auch bei stationärer Verfolgung des Pfades nicht verschwinden. Grundsätzlich kann es Fälle geben, für die eine nicht exakte Verfolgung des Pfades ($\mathbf{h}(\mathbf{x}) \neq \boldsymbol{\sigma}(\zeta_1)$) bewirkt, dass die absoluten Stellgrößenamplituden kleiner werden. Damit sinkt der Anteil der Kosten $(\mathbf{u})_{\mathbf{R}_u}$ während die Kosten $(\mathbf{h}(\mathbf{x}) - \boldsymbol{\sigma}(\zeta_1))_{\mathbf{Q}_e}$ steigen. Man kann sich nun vorstellen, dass es einen Kompromiss zwischen diesen beiden Anteilen geben kann, sodass der Wert des Kostenfunktionalen kleiner wird. Es kann also Fälle geben, für die es optimaler ist, eine Abweichung vom Pfad in Kauf zu nehmen. Dies bewirkt aber automatisch ein Verletzen der Invarianzeigenschaft (Z2).

Aus obiger Argumentation kann abgeleitet werden, dass es günstiger ist, die Stellgrößen \mathbf{u} nicht absolut zu gewichten, sondern nur deren Differenz zu den für exakte Verfolgung des Pfades benötigten Stellgrößen. Diese zu berechnen, stellt sich für allgemeine Systemklassen als herausfordernde Aufgabe dar. Daher soll im Folgenden eine *Einschränkung auf flache Systeme* der Form (1.104) vorgenommen werden, für die diese Aufgabe systematisch gelöst werden kann. Dies resultiert in einer MPC-Formulierung, die auch die Invarianzbedingung (Z2) sicherstellt.

Ein flaches System (1.104) zeichnet sich unter anderem dadurch aus, dass stets Parameterungen der Zustands- und Stellgrößen in der Form

$$\mathbf{x} = \boldsymbol{\psi}_x(\mathbf{y}, \dot{\mathbf{y}}, \ddot{\mathbf{y}}, \dots, \mathbf{y}^{(\beta-1)}) \quad (1.109a)$$

$$\mathbf{u} = \boldsymbol{\psi}_u(\mathbf{y}, \dot{\mathbf{y}}, \ddot{\mathbf{y}}, \dots, \mathbf{y}^{(\beta)}) \quad (1.109b)$$

mit endlichem $\beta \in \mathbb{N}$ gefunden werden können. Der nächste Schritt besteht darin, den Pfad $\boldsymbol{\sigma}(\theta)$ mit den zeitlichen Ableitungen

$$\mathbf{y} = \boldsymbol{\sigma}(\theta) \quad (1.110a)$$

$$\dot{\mathbf{y}} = \dot{\boldsymbol{\sigma}}(\theta, \dot{\theta}) = \boldsymbol{\sigma}'(\theta)\dot{\theta} \quad (1.110\text{b})$$

$$\ddot{\mathbf{y}} = \ddot{\boldsymbol{\sigma}}(\theta, \dot{\theta}, \ddot{\theta}) = \boldsymbol{\sigma}''(\theta)\dot{\theta}^2 + \boldsymbol{\sigma}'(\theta)\ddot{\theta} \quad (1.110\text{c})$$

$$\vdots \quad (1.110\text{d})$$

$$\mathbf{y}^{(\beta)} = \boldsymbol{\sigma}^{(\beta)}(\theta, \dot{\theta}, \ddot{\theta}, \dots, \theta^{(\beta)}) \quad (1.110\text{e})$$

in (1.109) einzusetzen. Dies ergibt

$$\mathbf{x} = \psi_x(\boldsymbol{\sigma}(\theta), \boldsymbol{\sigma}'(\theta)\dot{\theta}, \boldsymbol{\sigma}''(\theta)\dot{\theta}^2 + \boldsymbol{\sigma}'(\theta)\ddot{\theta}, \dots, \boldsymbol{\sigma}^{(\beta-1)}(\theta, \dot{\theta}, \ddot{\theta}, \dots, \theta^{(\beta-1)})) \quad (1.111\text{a})$$

$$\mathbf{u} = \psi_u(\boldsymbol{\sigma}(\theta), \boldsymbol{\sigma}'(\theta)\dot{\theta}, \boldsymbol{\sigma}''(\theta)\dot{\theta}^2 + \boldsymbol{\sigma}'(\theta)\ddot{\theta}, \dots, \boldsymbol{\sigma}^{(\beta)}(\theta, \dot{\theta}, \ddot{\theta}, \dots, \theta^{(\beta)})) \quad (1.111\text{b})$$

und motiviert die Wahl $\hat{r} = \beta + 1$. Daraus resultiert $\zeta = [\theta \ \dot{\theta} \ \ddot{\theta} \ \dots \ \theta^{(\beta)}]^T$ sowie gemäß (1.71)

$$\theta^{(\beta+1)} = v_\theta . \quad (1.112)$$

Durch Substituieren des Pfadparameters und seiner Ableitungen mit den Komponenten von ζ erhält man aus (1.111) die neuen Abbildungen

$$\mathbf{x} = \mathbf{p}_x(\zeta) \quad (1.113\text{a})$$

$$\mathbf{u} = \mathbf{p}_u(\zeta) . \quad (1.113\text{b})$$

Wenn die Beziehungen (1.113) erfüllt sind, befindet sich das System für beliebige Werte des Pfadparameters und seiner Ableitungen exakt am Pfad \mathcal{P} . Man beachte, dass (1.113a) die parametrisierte Form der Menge Γ^* darstellt, d. h. es gilt

$$\Gamma^* = \left\{ \bar{\mathbf{x}} \in \mathbb{R}^n \mid \bar{\mathbf{x}} = \mathbf{p}_x(\zeta), \zeta \in \mathcal{T} \times \mathbb{R}^\beta \right\} . \quad (1.114)$$

Die Abbildungen (1.113) liefern in Abhängigkeit von ζ die Werte der Zustände und der Stellgrößen für die exakte Verfolgung des Pfades im Ausgangsraum. Damit kann eine neue MPC-Formulierung basierend auf dem Optimalsteuerungsproblem

$$\min_{\mathbf{u}(\cdot), v_\theta(\cdot)} J(t_k, \mathbf{u}(\cdot), v_\theta(\cdot)) \quad (1.115\text{a})$$

$$\text{u.B.v.} \quad \dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}, \mathbf{u}) \quad \mathbf{x}(t_k) = \mathbf{x}_k \quad (1.115\text{b})$$

$$\dot{\zeta} = \mathbf{F}\zeta + \mathbf{h}v_\theta \quad \zeta(t_k) = \zeta_k \quad (1.115\text{c})$$

$$\mathbf{u}_{\min} \leq \mathbf{u}(t) \leq \mathbf{u}_{\max} \quad \forall t \in [t_k, t_k + T] \quad (1.115\text{d})$$

$$\zeta_1(t) \in \mathcal{T} \quad \forall t \in [t_k, t_k + T] \quad (1.115\text{e})$$

mit

$$J(t_k, \mathbf{u}(\cdot), v_\theta(\cdot)) = \int_{t_k}^{t_k+T} l(\mathbf{x}, \zeta, \mathbf{u}, v_\theta) dt + V(\mathbf{x}(t_k + T), \zeta(t_k + T)) \quad (1.115\text{f})$$

und

$$l(\mathbf{x}, \zeta, \mathbf{u}, v_\theta) = \frac{1}{2} \left[\left(\mathbf{x} - \mathbf{p}_x(\zeta) \right)_{\mathbf{Q}_e} + \left(\zeta - \zeta_d \right)_{\mathbf{Q}_\zeta} + \left(\mathbf{u} - \mathbf{p}_u(\zeta) \right)_{\mathbf{R}_u} + r_\theta v_\theta^2 \right] \quad (1.115\text{g})$$

$$V(\mathbf{x}, \boldsymbol{\zeta}) = \frac{1}{2} \left[(\mathbf{x} - \mathbf{p}_x(\boldsymbol{\zeta}))_{\mathbf{S}_e} + (\boldsymbol{\zeta} - \boldsymbol{\zeta}_d)_{\mathbf{S}_{\zeta}} \right] \quad (1.115h)$$

entwickelt werden.

Wenn der modellprädiktive Regler basierend auf (1.115) einen asymptotisch stabilen geschlossenen Kreis bewirkt, dann gilt $J^*(t_k) \rightarrow 0$ für $k \rightarrow \infty$. Das bedeutet, dass durch die positiv definite Gewichtung die Differenz $\mathbf{x} - \mathbf{p}_x(\boldsymbol{\zeta})$ verschwindet und zumindest die relevanten Komponenten von $\boldsymbol{\zeta} - \boldsymbol{\zeta}_d$ gegen Null konvergieren. Daher werden auch mit dem modellprädiktiven Regler basierend auf (1.115) die Ziele (Z1) und (Z3) berücksichtigt. Da nun keine absolute Gewichtung von \mathbf{u} sondern eine relative Gewichtung zu $\mathbf{p}_u(\boldsymbol{\zeta})$ verwendet wird, ist im Allgemeinen auch die Invarianzbedingung (Z2) erfüllt. Voraussetzung dafür ist allerdings, dass $\mathbf{u}_{\min} \leq \mathbf{p}_u(\boldsymbol{\zeta}^*(t, t_k)) \leq \mathbf{u}_{\max} \forall t \in [t_k, t_k + T]$ erfüllt ist, weil nur dann $\mathbf{u}^*(\cdot, t_k) \equiv \mathbf{p}_u(\boldsymbol{\zeta}^*(\cdot, t_k))$ gelten kann. Grob gesprochen hängt die Erfüllung der Bedingung $\mathbf{u}_{\min} \leq \mathbf{p}_u(\boldsymbol{\zeta}^*(\cdot, t_k)) \leq \mathbf{u}_{\max}$ von der Größe der Komponenten von $\boldsymbol{\zeta}^*(\cdot, t_k)$ ab. Daher müssen sowohl die Gewichtungen \mathbf{Q}_{ζ} , \mathbf{S}_{ζ} und r_{θ} als auch $\boldsymbol{\zeta}_d$ (d. h. das Ziel (Z3)) so gewählt werden, dass $\mathbf{u}^*(\cdot, t_k) \equiv \mathbf{p}_u(\boldsymbol{\zeta}^*(\cdot, t_k))$ erfüllbar ist. Für die Erfüllung der Invarianzbedingung muss weiters gelten, dass $\mathbf{x}(t_0) = \mathbf{p}_x(\boldsymbol{\zeta}(t_0))$ gilt, d. h. insbesondere ist die Wahl von $\boldsymbol{\zeta}(t_0)$ auch hier wieder entscheidend.

Example 1.9 (Fahrzeug mit variabler Vorwärtsgeschwindigkeit). Als Anwendungsbeispiel für die modellprädiktive Pfadfolgeregelung soll wieder das Fahrzeug mit variabler Vorwärtsgeschwindigkeit (1.48) dienen. Für dieses ist die Annahme $m = p$ erfüllt. Analog zu Beispiel 1.7 soll der Pfadfolgeregel sicherstellen, dass das Fahrzeug einer Lemniskate folgt und diese alle 20 Sekunden komplett abfährt (Ziel (Z3)).

Durch Anschreiben des Systemausgangs und Bilden der Ableitungen

$$\mathbf{y} = \begin{bmatrix} x \\ y \end{bmatrix} \quad (1.116a)$$

$$\dot{\mathbf{y}} = \begin{bmatrix} v_l \cos(\varphi) \\ v_l \sin(\varphi) \end{bmatrix} \quad (1.116b)$$

$$\ddot{\mathbf{y}} = \underbrace{\begin{bmatrix} -v_l^2 \sin(\varphi) & \cos(\varphi) \\ v_l^2 \cos(\varphi) & \sin(\varphi) \end{bmatrix}}_{\mathbf{D}(\mathbf{x})} \mathbf{u} \quad (1.116c)$$

ergibt sich, dass das mathematische Modell des Fahrzeugs (1.48) eingangs-zustands-linearisierbar und damit flach ist. Aus (1.116a) und (1.116b) erhält man zwei Lösungen für die Parametrierung der Zustände

$$\psi_x^+(\mathbf{y}, \dot{\mathbf{y}}) = \begin{bmatrix} y_1 \\ y_2 \\ \arctan\left(\frac{\dot{y}_2}{\dot{y}_1}\right) \\ \sqrt{\dot{y}_1^2 + \dot{y}_2^2} \end{bmatrix} \quad \psi_x^-(\mathbf{y}, \dot{\mathbf{y}}) = \begin{bmatrix} y_1 \\ y_2 \\ \arctan\left(\frac{-\dot{y}_2}{-\dot{y}_1}\right) \\ -\sqrt{\dot{y}_1^2 + \dot{y}_2^2} \end{bmatrix}. \quad (1.117)$$

Aus (1.116c) folgt die Parametrierung der Stellgrößen

$$\psi_u^\pm(\mathbf{y}, \dot{\mathbf{y}}, \ddot{\mathbf{y}}) = \mathbf{D}^{-1}(\mathbf{x})|_{\mathbf{x}=\psi_x^\pm(\mathbf{y}, \dot{\mathbf{y}})} \ddot{\mathbf{y}}, \quad (1.118)$$

für die nach Einsetzen von (1.117) ebenfalls zwei Lösungen

$$\psi_u^+(\mathbf{y}, \dot{\mathbf{y}}, \ddot{\mathbf{y}}) = \begin{bmatrix} \frac{\dot{y}_1 \ddot{y}_2 - \dot{y}_2 \ddot{y}_1}{(\dot{y}_1^2 + \dot{y}_2^2)^{\frac{3}{2}}} \\ \frac{(\dot{y}_1^2 + \dot{y}_2^2)^{\frac{3}{2}}}{\sqrt{\dot{y}_1^2 + \dot{y}_2^2}} \\ \frac{\dot{y}_1 \ddot{y}_1 + \dot{y}_2 \ddot{y}_2}{\sqrt{\dot{y}_1^2 + \dot{y}_2^2}} \end{bmatrix} \quad \psi_u^-(\mathbf{y}, \dot{\mathbf{y}}, \ddot{\mathbf{y}}) = \begin{bmatrix} -\frac{\dot{y}_1 \ddot{y}_2 - \dot{y}_2 \ddot{y}_1}{(\dot{y}_1^2 + \dot{y}_2^2)^{\frac{3}{2}}} \\ -\frac{(\dot{y}_1^2 + \dot{y}_2^2)^{\frac{3}{2}}}{\sqrt{\dot{y}_1^2 + \dot{y}_2^2}} \\ -\frac{\dot{y}_1 \ddot{y}_1 + \dot{y}_2 \ddot{y}_2}{\sqrt{\dot{y}_1^2 + \dot{y}_2^2}} \end{bmatrix} \quad (1.119)$$

existieren. Aus (1.119) folgt $\beta = 2$ und damit $\hat{r} = \beta + 1 = 3$ mit dem Hilfssystem $\theta^{(3)} = v_\theta$. Durch Einsetzen des Pfades mit den zugehörigen Ableitungen

$$\boldsymbol{\sigma}_l(\theta) = \begin{bmatrix} \cos(\theta) \\ \cos(\theta) \sin(\theta) \end{bmatrix} \quad \dot{\boldsymbol{\sigma}}_l(\theta, \dot{\theta}) = \dot{\theta} \begin{bmatrix} -\sin \theta \\ -(\sin \theta)^2 + (\cos \theta)^2 \end{bmatrix} \quad (1.120a)$$

$$\ddot{\boldsymbol{\sigma}}_l(\theta, \dot{\theta}, \ddot{\theta}) = \begin{bmatrix} -\dot{\theta}^2 \cos(\theta) - \ddot{\theta} \sin(\theta) \\ -4\dot{\theta}^2 \cos(\theta) \sin(\theta) + \ddot{\theta}(-(\sin \theta)^2 + (\cos \theta)^2) \end{bmatrix} \quad (1.120b)$$

und $\theta = \zeta_1$, $\dot{\theta} = \zeta_2$, $\ddot{\theta} = \zeta_3$ in (1.117) und (1.119) erhält man

$$\mathbf{p}_x^+(\zeta) = \begin{bmatrix} \cos \zeta_1 \\ \cos(\zeta_1) \sin(\zeta_1) \\ \arctan\left(\frac{\zeta_2(2(\cos \zeta_1)^2 - 1)}{-\sin(\zeta_1)\zeta_2}\right) \\ \sqrt{\zeta_2^2(4(\cos \zeta_1)^4 - 5(\cos \zeta_1)^2 + 2)} \end{bmatrix} \quad (1.121a)$$

$$\mathbf{p}_x^-(\zeta) = \begin{bmatrix} \cos \zeta_1 \\ \cos(\zeta_1) \sin(\zeta_1) \\ \arctan\left(\frac{-\zeta_2(2(\cos \zeta_1)^2 - 1)}{\sin(\zeta_1)\zeta_2}\right) \\ -\sqrt{\zeta_2^2(4(\cos \zeta_1)^4 - 5(\cos \zeta_1)^2 + 2)} \end{bmatrix} \quad (1.121b)$$

sowie $\mathbf{p}_u^+(\zeta)$ und $\mathbf{p}_u^-(\zeta)$, die auf Grund ihres Umfanges nicht explizit angegeben sind.

Damit das Ziel (Z3) erreicht werden kann, wird

$$\zeta_d = \begin{bmatrix} \text{beliebig} \\ \pm \frac{2\pi}{20} \\ 0 \end{bmatrix} \quad (1.122)$$

gesetzt. Die Lösungen (1.121) für $\mathbf{p}_x(\zeta)$ unterscheiden sich in der Orientierung des Fahrzeugs und dem Vorzeichen der Vorwärtsgeschwindigkeit. Das heißt, der Unterschied besteht darin, ob das Fahrzeug den Pfad vor- oder rückwärts abfährt.

Hinsichtlich einer Verwendung von $\mathbf{p}_x(\zeta)$ und $\mathbf{p}_u(\zeta)$ im Optimalsteuerungsproblem stellt sich die Frage, welche Lösung wann verwendet wird. Grundsätzlich soll gefordert werden, dass das Fahrzeug den Pfad in Vorwärtsrichtung, d. h. mit $v_l > 0$, abfahren soll. Wenn nun $\zeta_{d,2} > 0$ gilt, dann ergibt sich im geschlossenen Kreis im eingeschwungenen Zustand ebenfalls $\zeta_2 > 0$. Für diesen eingeschwungenen Zustand soll $v_l > 0$ gelten, womit für $\zeta_2 > 0$ die Lösungen $\mathbf{p}_x^+(\zeta)$ und $\mathbf{p}_u^+(\zeta)$ verwendet werden müssen. Für $\zeta_2 \leq 0$ kommen $\mathbf{p}_x^-(\zeta)$ und $\mathbf{p}_u^-(\zeta)$ zum Einsatz. Im Fall $\zeta_{d,2} < 0$ gilt im eingeschwungenen Zustand $\zeta_2 < 0$. Dafür soll aber ebenfalls $v_l > 0$ gelten. Daher sind die Verhältnisse für $\zeta_{d,2} < 0$ genau umgekehrt. Falls $\zeta_2 < 0$ wird $\mathbf{p}_x^+(\zeta)$ und $\mathbf{p}_u^+(\zeta)$ verwendet, ansonsten $\mathbf{p}_x^-(\zeta)$ und $\mathbf{p}_u^-(\zeta)$.

Den im Folgenden gezeigten Simulationsergebnissen liegen die Parameterwerte

$$\mathbf{Q}_e = \begin{bmatrix} 100 & 0 & 0 & 0 \\ 0 & 100 & 0 & 0 \\ 0 & 0 & 100 & 0 \\ 0 & 0 & 0 & 100 \end{bmatrix} \quad \mathbf{S}_e = \begin{bmatrix} 10 & 0 & 0 & 0 \\ 0 & 10 & 0 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 10 \end{bmatrix} \quad \mathbf{R}_u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (1.123a)$$

$$\mathbf{Q}_\zeta = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{S}_\zeta = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad r_\theta = 1 \quad (1.123b)$$

zugrunde. Man beachte, dass \mathbf{Q}_ζ und \mathbf{S}_ζ passend zu (Z3) gewählt sind. Die Beschränkungen der Stellgrößen werden zu

$$\mathbf{u}_{\min} = \begin{bmatrix} -5 \\ -3 \end{bmatrix} \quad \mathbf{u}_{\max} = \begin{bmatrix} 5 \\ 3 \end{bmatrix} \quad (1.124)$$

festgelegt. Da $\theta = \zeta_1 \in \mathbb{R}$, ist die Beschränkung (1.115e) hinfällig. Der Optimierungshorizont sei $T = 4\text{s}$ und die Abtastzeit wird zu $T_a = 50\text{ms}$ gewählt. Das Optimalsteuerungsproblem (1.115) wird durch ein direktes Verfahren und anschließende Volldiskretisierung [1.0] in ein statisches Optimierungsproblem überführt, vgl. Abschnitt ???. Dieses statische Optimierungsproblem wird in weiterer Folge mit dem SQP-Verfahren des Softwarepaketes SNOPT [1.0] gelöst. Die resultierenden optimalen Stellgrößen werden konstant über das Abtastintervall aufgeschaltet, d. h. $\mathbf{u}(t) = \mathbf{u}^*(t_k, t_k)$, $v_\theta(t) = v_\theta^*(t_k, t_k)$, $t \in [t_k, t_k + T_a]$. Die erste Komponente des Anfangswerts ζ_0 im Abtastschritt $t_0 = 0$ wird gemäß (1.84) bestimmt. Die zweite und dritte Komponente wird jeweils mit der zweiten und dritten Komponente von ζ_d gemäß (1.122) gleichgesetzt.

Abbildung 1.10 zeigt die Zeitverläufe im geschlossenen Kreis für den Anfangswert $\mathbf{x}_0 = \begin{bmatrix} 3 & 1 & \frac{3\pi}{2} & 0.45 \end{bmatrix}^\top$. Aus der anfänglichen Abweichung konvergieren sowohl die Stellgrößen als auch die Zustände schnell zu \mathbf{p}_u und \mathbf{p}_x . Die Stellgrößenbeschränkungen werden mehrmals aktiv. Der optimale Wert des Kostenfunktionalen $J^*(\cdot)$ wird von Abtastschritt zu Abtastschritt kleiner.

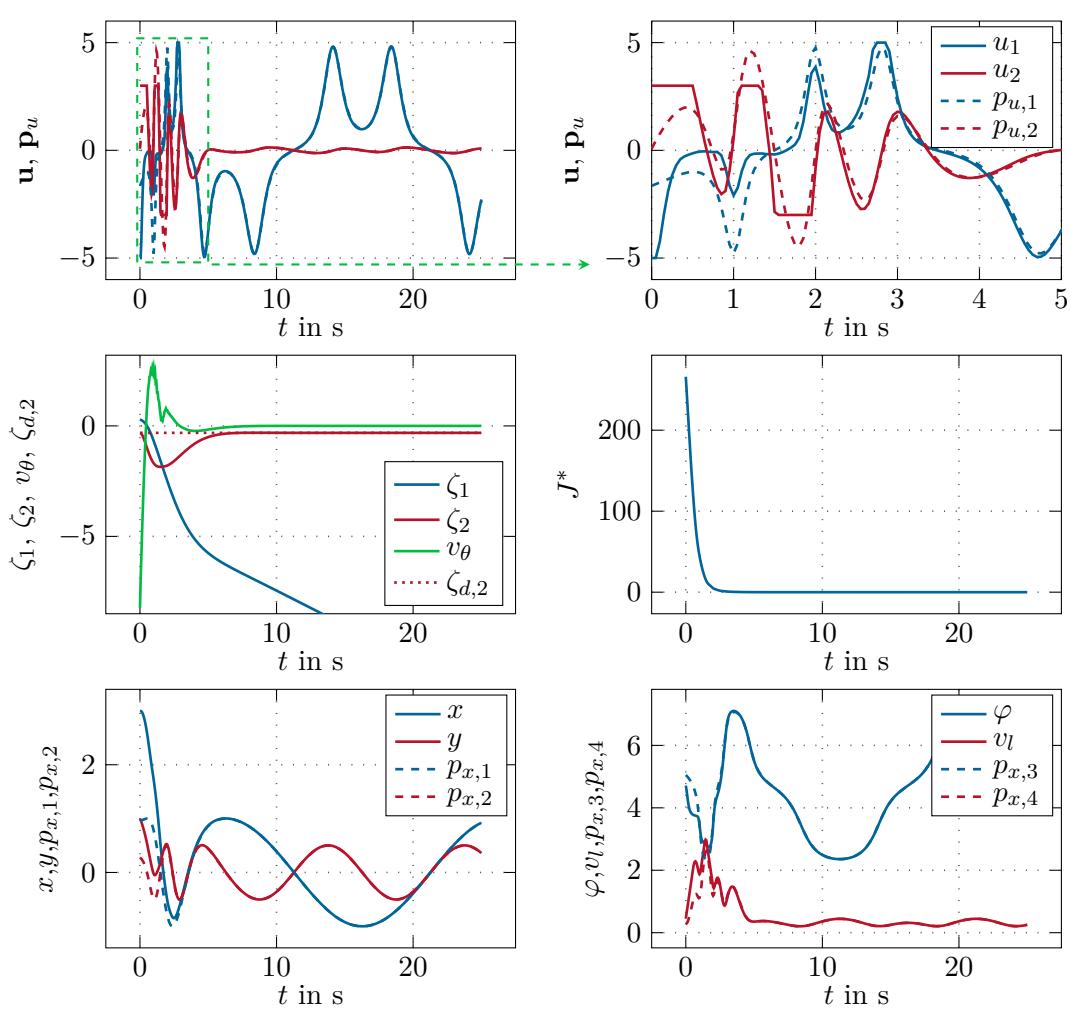


Figure 1.10: Zeitverläufe im geschlossenen Kreis bei der modellprädiktiven Pfadfolgeregelung einer Lemniskate für das Fahrzeug mit $v_l \neq \text{konst.}$

In Abbildung 1.11 ist die zugehörige Trajektorie des Fahrzeuges dargestellt.

Die Trajektorien des Fahrzeuges für mehrere verschiedene Anfangswerte sind in Abbildung 1.12 zu sehen.

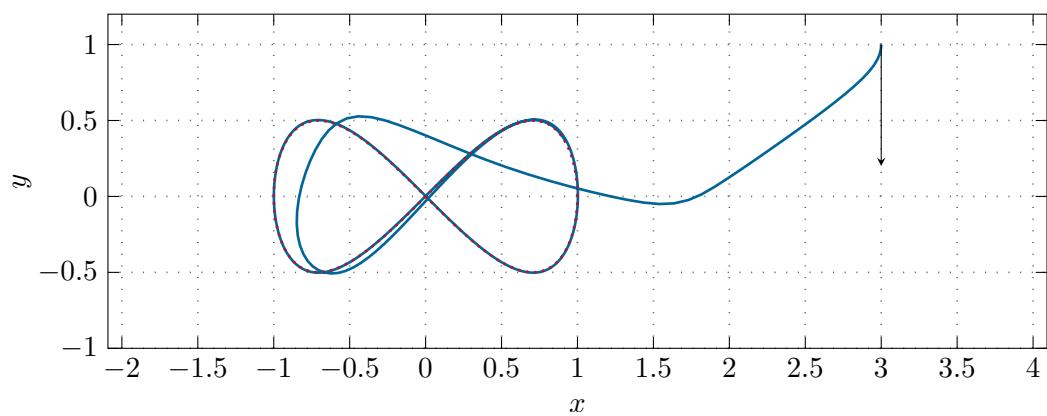


Figure 1.11: Trajektorie des Fahrzeuges mit $v_l \neq \text{konst.}$ für die modellprädiktive Pfadfolgeregelung einer Lemniskate.

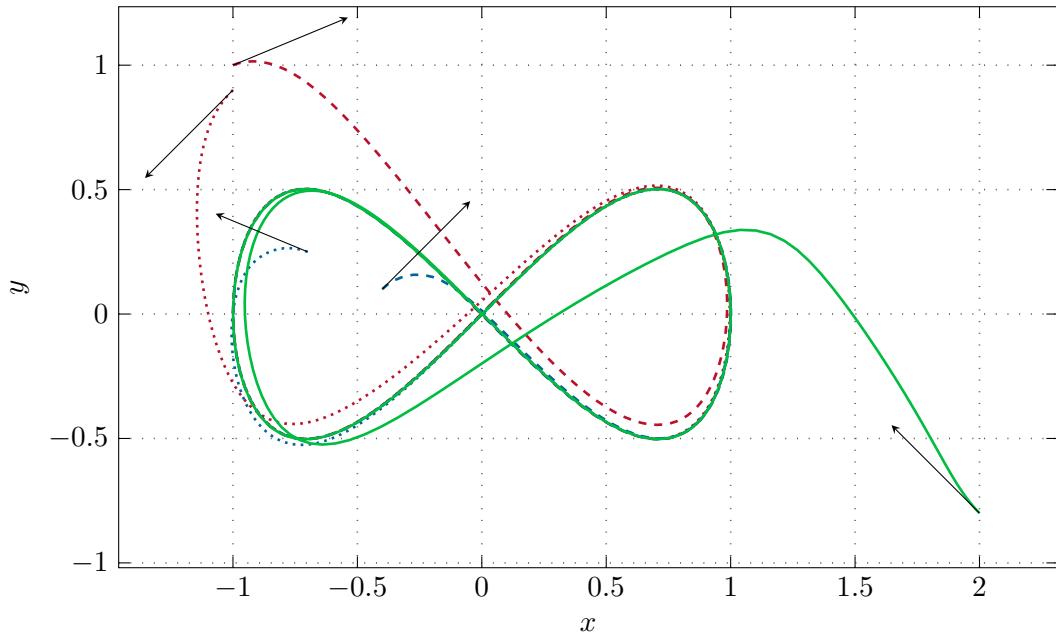


Figure 1.12: Trajektorien des Fahrzeuges mit $v_l \neq \text{konst.}$ für die modellprädiktive Pfadfolgeregelung einer Lemniskate und verschiedenen Anfangswerten.

1.5 Literatur

- [1.0] A. Kugi, *Skriptum zur VO Regelungssysteme 2 (SS 2019)*, Institut für Automatisierungs- und Regelungstechnik, TU Wien, 2019. [Online]. Available: <https://www.acin.tuwien.ac.at/master/regelungssysteme-2/>.
- [1.0] V. A. Toponogov, *Differential Geometry of Curves and Surfaces - A Concise Guide*. Boston: Birkhäuser, 2006.
- [1.0] A. Isidori, *Nonlinear Control Systems*, 3rd ed. London: Springer, 1995.
- [1.0] C. Nielsen, “Set stabilization using transverse feedback linearization,” Ph.D. dissertation, University of Toronto, 2009.
- [1.0] T. Faulwasser, “Optimization-based solutions to constrained trajectory-tracking and path-following problems,” Ph.D. dissertation, Otto-von-Guericke-Universität Magdeburg, 2012.
- [1.0] J. T. Betts, *Practical Methods for Optimal Control Using Nonlinear Programming* (Advances in Design and Control). Philadelphia, USA: SIAM - Society for Industrial and Applied Mathematics, 2001.
- [1.0] P. E. Gill, W. Murray, and M. A. Saunders, *User’s guide for SNOPT version 7: Software for large-scale nonlinear programming*, Department of Mathematics, University of California, San Diego, CA, USA, Jun. 2008.

2 Dissipativität und Passivität

Vereinfachend gesprochen, ist das Konzept der Dissipativität und Passivität die systemtheoretische Verallgemeinerung des Energieerhaltungsprinzips, welches besagt, dass in einem abgeschlossenen System Energie weder erzeugt noch vernichtet werden kann. Eine nähere Betrachtung des systemtheoretischen Konzeptes der Dissipativität wird jedoch zeigen, dass dies a priori mit dem Prinzip der Energieerhaltung nichts zu tun hat und lediglich bei gewissen physikalischen Systemen analoge Aussagen zulässt. Diese Analogie zu physikalischen Systemen trägt aber sicherlich zum Verständnis dieser Konzepte bei, weshalb im Folgenden zwei physikalische Systeme, ein Wärmeübertragungssystem und ein elektromechanisches System, diskutiert werden.

2.1 Glühsimulator

Abbildung 2.1 zeigt die schematische Darstellung eines so genannten Glühsimulators, der dazu verwendet wird, durch Ohmsches Erwärmen und freie bzw. erzwungene Konvektion (Pressluft oder Ventilator) für Metallproben vorgegebene Temperaturprofile abzufahren.

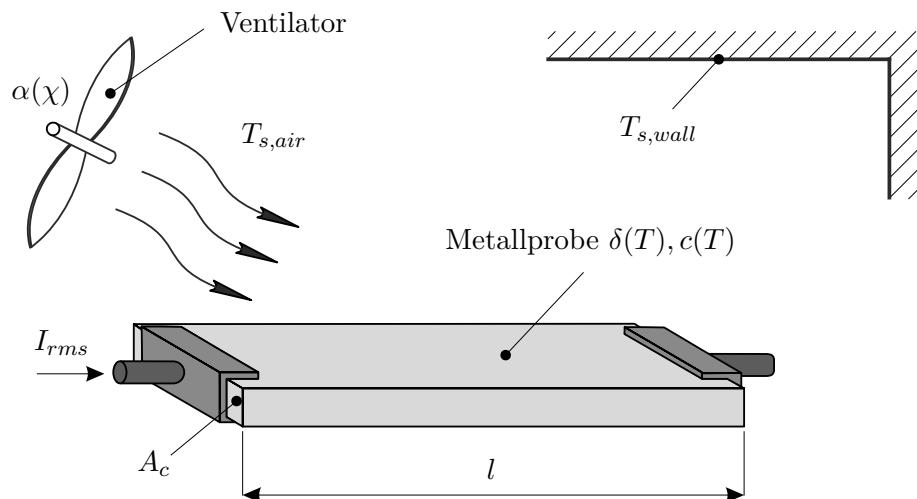


Figure 2.1: Glühsimulator.

Es ist naheliegend für dieses System die elektromechanischen Effekte zu vernachlässigen und die Änderung der im System gespeicherten Energie allein durch die Änderung der thermisch gespeicherten Energie zu erfassen. Das Energieerhaltungsprinzip besagt dann, dass die Änderung der thermisch gespeicherten Energie V der Beziehung

$$\frac{d}{dt}V = p_{in} - p_{out} \quad (2.1)$$

genügt, wobei p_{in} und p_{out} die Energieflüsse in das System und vom System beschreiben. Es wird angenommen, dass die Temperatur T in der Metallprobe zu jedem Zeitpunkt t gleichförmig verteilt ist, dass die Oberfläche der Probe sehr klein verglichen mit den umgebenden Wänden ist, und dass die Wärmeleitung vernachlässigt werden kann. Die in der Probe gespeicherte thermische Energie V lautet

$$V(T) = c(T)mT \quad (2.2)$$

mit der konstanten Probenmasse m und der spezifischen Wärmekapazität $c(T)$. Mit Hilfe des Ohmschen Gesetzes errechnet sich der Energiefluss in die Probe zu

$$p_{in} = I_{rms}^2 \delta(T) \frac{l}{A_c} \quad (2.3)$$

mit dem Effektivwert des durch die Probe fließenden Stromes I_{rms} , dem spezifischen Widerstand $\delta(T)$, der Länge der Probe l und der Probenquerschnittsfläche A_c . Die Energieflüsse von der Probe in die Umgebung werden einerseits durch die freie und erzwungene Konvektion

$$p_{out,1} = \alpha(\chi) A_s (T - T_{s,air}) \quad (2.4)$$

und andererseits durch die Wärmestrahlung

$$p_{out,2} = \varepsilon \sigma A_s (T^4 - T_{s,wall}^4) \quad (2.5)$$

verursacht. Dabei bezeichnen A_s die Oberfläche der Metallprobe, $T_{s,air}$ und $T_{s,wall}$ die Temperaturen der umgebenden Luft und Wände, ε ist der Emissionsgrad, $\sigma = 5.67 \cdot 10^{-8} \text{ Wm}^{-2}\text{K}^{-4}$ die Stefan-Boltzmann Konstante und $\alpha(\chi)$ ist der Konvektionskoeffizient, wobei χ im Falle eines Lüfters für die Drehwinkelgeschwindigkeit des Lüfters und im Falle von Druckluft für den Druck steht. Bei freier Konvektion ist $\alpha(\chi)$ konstant und liegt im Bereich von $2 - 25 \text{ Wm}^{-2}\text{K}^{-1}$. Das mathematische Modell des Glühsimulators erhält man einfach durch Einsetzen von (2.2) - (2.5) in (2.1) mit der Zustandsgröße T und den Eingangsgrößen $\mathbf{u}^T = [I_{rms}, \chi, T_{s,air}, T_{s,wall}]$. Integriert man (2.1) entlang einer Lösungskurve vom Zeitpunkt $t_0 = 0$ zum Zeitpunkt t für gegebene Eingangsgrößen $\mathbf{u}(\tau)$, $0 \leq \tau \leq t$, dann erhält man

$$V(T(t)) - V(T(0)) = \int_0^t s(I_{rms}, \chi, T_{s,air}, T_{s,wall}, T) d\tau \quad (2.6)$$

mit

$$s(I_{rms}, \chi, T_{s,air}, T_{s,wall}, T) = I_{rms}^2 \delta(T) \frac{l}{A_c} - \alpha(\chi) A_s (T - T_{s,air}) - \varepsilon \sigma A_s (T^4 - T_{s,wall}^4). \quad (2.7)$$

Gleichung (2.6) besagt, dass die zum Zeitpunkt t im System gespeicherte thermische Energie V gleich der zum Zeitpunkt $t_0 = 0$ gespeicherten Energie plus oder minus der in dieser Zeit mit der so genannten Versorgungsrate $s(I_{rms}, \chi, T_{s,air}, T_{s,wall}, T)$ dem System zu- oder abgeführten Energie ist.

2.2 Einfaches Elektromagnetventil

Abbildung 2.2 zeigt das Elektromagnetventil mit einem zylindrischen Gehäuse und einem zylindrischen Stössel mit der Masse m und dem Durchmesser D . Die aus N Windungen bestehende Spule mit einem gesamten Innenwiderstand R wird mit einer Spannung U_0 versorgt. Es wird angenommen, dass der magnetische Widerstand des Gehäuses und des Stössels Null ist, dass die Gleithülse die gleiche Permeabilität wie Luft besitzt und dass für die geometrischen Abmessungen gilt $h \ll D$ und $\delta \ll b$ (keine Streuflüsse).

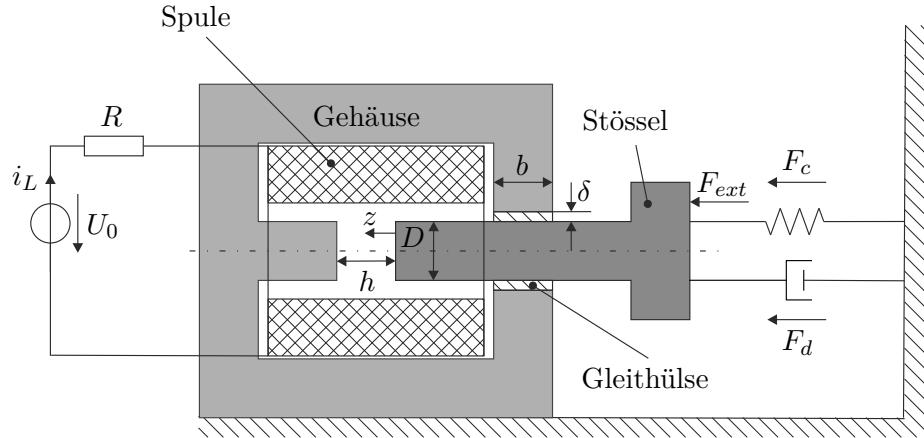


Figure 2.2: Einfaches Elektromagnetventil.

Auf analoge Art und Weise zu (2.1) gilt für die Änderung der im System gespeicherten Energie V die Beziehung

$$\frac{d}{dt}V = p_{in} - p_{out} - p_{diss} \quad (2.8)$$

mit den Energieflossen p_{in} und p_{out} , die über die Systemgrenzen in das System bzw. vom System fließen und mit der in Wärme dissipierten Leistung p_{diss} .

Unter den obigen Voraussetzungen errechnet sich die im Magnetkreis gespeicherte Koenergie in der Form

$$\check{w}_L = \frac{1}{2}L(z)i_L^2 \quad (2.9)$$

mit der Ersatzinduktivität des magnetischen Kreises

$$L(z) = \frac{\mu_0 N^2 D^2 \pi (D + \delta) \pi b}{4(h - z)(D + \delta) \pi b + \delta D^2 \pi} \quad (2.10)$$

und der Permeabilität von Luft $\mu_0 = 4\pi \times 10^{-7} \text{ Vs/(A m)}$.

Exercise 2.1. Rechnen Sie die Beziehung für die Induktivität $L(z)$ von (2.10) nach.

Da das betrachtete Elektromagnetventil magnetisch linear ist, sind die Ausdrücke für Energie \hat{w}_L und Koenergie \check{w}_L identisch. Die auf den Stössel wirkende Magnetkraft errechnet sich zu

$$F_{mag} = \frac{\partial}{\partial z} \check{w}_L = \frac{1}{2} \frac{\partial L(z)}{\partial z} i_L^2. \quad (2.11)$$

Wie in Abbildung 2.2 gezeichnet, wirkt der Stössel gegen ein lineares Feder-Dämpfer System mit der Dämpfungskraft $F_d = dv$, $v = \dot{z}$, $d > 0$, der Federkraft $F_c = cz(t)$, $c > 0$ und einer externen Kraft F_{ext} . Das mathematische Modell des Elektromagnetventils lautet dann

$$\frac{d}{dt}z = v \quad (2.12)$$

$$\frac{d}{dt}v = \frac{1}{m} \left(\frac{1}{2} \frac{\partial L(z)}{\partial z} i_L^2 - cz - dv + F_{ext} \right) \quad (2.13)$$

$$\frac{d}{dt}i_L = \frac{1}{L(z)} \left(U_0 - Ri_L - \frac{\partial L(z)}{\partial z} i_L v \right) \quad (2.14)$$

mit den Zustandsgrößen $\mathbf{x}^T = [z, v, i_L]$ und den Eingangsgrößen $\mathbf{u}^T = [U_0, F_{ext}]$.

Die im System gespeicherte Energie setzt sich nun aus der magnetischen Energie (2.9), der kinetischen Energie des Stössels und der potenziellen Energie der Feder

$$V = \frac{1}{2} \left(L(z) i_L^2 + mv^2 + cz^2 \right) \quad (2.15)$$

zusammen. Die Änderung der gespeicherten Energie V entlang einer Lösungskurve ergibt sich in der Form

$$\frac{d}{dt}V = \underbrace{U_0 i_L + F_{ext} v}_{p_{in} - p_{out}} - \underbrace{(dv^2 + Ri_L^2)}_{p_{diss}}. \quad (2.16)$$

Integriert man nun wieder (2.16) entlang einer Lösungskurve vom Zeitpunkt $t_0 = 0$ zum Zeitpunkt t für gegebene Eingangsgrößen $\mathbf{u}(\tau)$, $0 \leq \tau \leq t$, dann erhält man wegen $p_{diss} \geq 0$

$$V(\mathbf{x}(t)) - V(\mathbf{x}(t_0)) \leq \int_{t_0}^t s(U_0, F_{ext}, i_L, v) d\tau \quad (2.17)$$

mit der Versorgungsrate

$$s(U_0, F_{ext}, i_L, v) = U_0 i_L + F_{ext} v. \quad (2.18)$$

2.3 Systemtheoretisches Konzept

2.3.1 Dissipativität

Den nachfolgenden Betrachtungen liege ein nichtlineares dynamisches System der Form

$$\begin{aligned} \frac{d}{dt}\mathbf{x} &= \mathbf{f}(\mathbf{x}, \mathbf{u}) \\ \mathbf{y} &= \mathbf{h}(\mathbf{x}, \mathbf{u}) \end{aligned} \quad (2.19)$$

mit dem Zustand $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$, dem Stelleingang $\mathbf{u} \in \mathcal{U} \subset \mathbb{R}^m$ und dem Ausgang $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^p$ zu Grunde. Es sei angenommen, dass der Zustand $\mathbf{x}(t)$ zu jedem Zeitpunkt t eindeutig durch die Wahl der Eingangsgröße $\mathbf{u}(t)$ und des Anfangszustandes $\mathbf{x}(0) = \mathbf{x}_0$, bestimmt ist. Dies erlaubt es, die so genannte Versorgungsrate $s(\mathbf{u}, \mathbf{y}) : \mathcal{U} \times \mathcal{Y} \rightarrow \mathbb{R}$, eine

reellwertige Funktion, die für alle Anfangswerte $\mathbf{x}_0 \in \mathcal{X}$ und alle Eingangsgrößen \mathbf{u} die Bedingung

$$\int_0^t |s(\mathbf{u}, \mathbf{y})| d\tau < \infty \quad (2.20)$$

für alle Zeiten $t \geq 0$ erfüllt, einzuführen.

Definition 2.1. Das System (2.19) heißt *dissipativ bezüglich der Versorgungsrate s*, wenn eine nichtnegative Funktion $V(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$ so existiert, dass die so genannte *integrale Dissipativitätsungleichung*

$$V(\mathbf{x}(t)) - V(\mathbf{x}(0)) \leq \int_0^t s(\mathbf{u}(\tau), \mathbf{y}(\tau)) d\tau \quad (2.21)$$

für alle Anfangswerte $\mathbf{x}(0) \in \mathcal{X}$ und alle Eingangsgrößen $\mathbf{u}(t)$ für alle Zeiten $t \geq 0$ erfüllt ist. Die Funktion $V(\mathbf{x})$ wird als *Speicherfunktion* bezeichnet. Falls in (2.21) das Gleichheitszeichen gilt, nennt man das System (2.19) *verlustlos bezüglich der Versorgungsrate s*.

Im Sinne dieser Definition ist der Glühsimulator von Abbildung 2.1 verlustlos bezüglich der Versorgungsrate (2.7) und das Elektromagnetventil von Abbildung 2.2 ist dissipativ bezüglich der Versorgungsrate (2.18). Wenn die Speicherfunktion $V(\mathbf{x})$ bezüglich \mathbf{x} stetig differenzierbar ist, dann kann man die Änderung von $V(\mathbf{x})$ entlang einer Lösungskurve von (2.19) berechnen und man erhält die so genannte differenzielle Dissipativitätsungleichung

$$\frac{d}{dt} V(\mathbf{x}) \leq s(\mathbf{u}(t), \mathbf{y}(t)) \quad (2.22)$$

für alle Zeiten $t \geq 0$.

2.3.2 Passivität

Die Passivität kann als Spezialfall der Dissipativität aufgefasst werden. Zur Definition betrachte man wiederum das System (2.19), wobei nun die Dimension des Systemeingangs m gleich der Dimension des Ausgangs p ist.

Definition 2.2. Das System (2.19) mit $m = p$ nennt man *passiv*, wenn eine Konstante δ so existiert, dass die Ungleichung

$$\int_0^t \mathbf{y}^T \mathbf{u} d\tau \geq \delta \quad (2.23)$$

für alle zulässigen Eingangsgrößen $\mathbf{u}(t)$ und alle $t \geq 0$ erfüllt ist.

Wenn darüberhinaus für geeignete reelle Konstanten α, β die Ungleichung

$$\int_0^t \mathbf{y}^T \mathbf{u} d\tau \geq \delta + \alpha \int_0^t \mathbf{u}^T \mathbf{u} d\tau \quad \text{bzw.} \quad \int_0^t \mathbf{y}^T \mathbf{u} d\tau \geq \delta + \beta \int_0^t \mathbf{y}^T \mathbf{y} d\tau \quad (2.24)$$

für alle zulässigen Eingangsgrößen $\mathbf{u}(t)$ und alle $t \geq 0$ erfüllt ist, dann nennt man das System α -eingangspassiv bzw. β -ausgangspassiv.

Offensichtlich muss $\delta \leq 0$ gelten, denn die Ungleichung (2.23) muss auch für die Eingangsgröße $\mathbf{u}(t) = \mathbf{0}$ gültig sein.

Theorem 2.1 (Verbindung Passivität und Dissipativität). *Existiert nun für das System (2.19) mit $m = p$ eine nichtnegative Funktion $V(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}$ so, dass gilt (integrale Passivitätsungleichung)*

$$V(\mathbf{x}(t)) - V(\mathbf{x}(0)) \leq \int_0^t \mathbf{y}^T \mathbf{u} d\tau \quad (2.25)$$

für alle zulässigen Eingangsgrößen $\mathbf{u}(t)$, alle $\mathbf{x}(0)$ und alle $t \geq 0$, dann ist das System (2.19) vom Eingang \mathbf{u} zum Ausgang \mathbf{y} passiv. Offensichtlich ist dies gemäß Definition 2.1 äquivalent dazu, dass das System (2.19) bezüglich der speziellen bilinearen Versorgungsrate $s(\mathbf{u}, \mathbf{y}) = \langle \mathbf{y}, \mathbf{u} \rangle = \mathbf{y}^T \mathbf{u}$ dissipativ ist. Ist darüberhinaus das System (2.19) bezüglich der Versorgungsrate $s(\mathbf{u}, \mathbf{y}) = \mathbf{y}^T \mathbf{u} - \alpha \|\mathbf{u}\|^2$ bzw. $s(\mathbf{u}, \mathbf{y}) = \mathbf{y}^T \mathbf{u} - \beta \|\mathbf{y}\|^2$ für geeignete reelle Konstanten α, β dissipativ, so ist (2.19) α -eingangspassiv bzw. β -ausgangspassiv. Ein verlustloses passives System nennt man in diesem Zusammenhang auch ein konservatives System.

Proof. Der Beweis des Satzes ist trivial, da wegen $V(\mathbf{x}) \geq 0$ aus (2.25) unmittelbar folgt

$$\int_0^t \mathbf{y}^T \mathbf{u} d\tau \geq -V(\mathbf{x}(0)) = \delta. \quad (2.26)$$

□

Mit dieser Definition erkennt man unmittelbar, dass das Elektromagnetventil von Abbildung 2.2 mit dem Eingang $\mathbf{u}^T = [U_0, F_{ext}]$ und dem Ausgang $\mathbf{y}^T = [i_L, v]$ passiv, ja sogar β -ausgangspassiv mit $0 < \beta < \min(d, R)$ ist, da für die dissipierte Leistung von (2.16) gilt $p_{diss} = dv^2 + Ri_L^2 \geq \beta \|\mathbf{y}\|^2$.

Die physikalische Interpretation der Passivitätsungleichung (2.25) lautet nun wie folgt: Gibt der Ausdruck $\mathbf{y}^T \mathbf{u}$ eine Leistung an (z.B. geeignete Paare von Strömen und Spannungen bei elektrischen Systemen oder kollokierte Geschwindigkeiten und Kräfte bei mechanischen Systemen) und ist $V(\mathbf{x})$ die im System gespeicherte Energie, so besagt die Passivitätsungleichung (2.25), dass die Zunahme der im System gespeicherten Energie kleiner oder gleich der dem System zugeführten Energie ist.

Exercise 2.2. Zeigen Sie, dass der Integrator mit der Zustandsdarstellung

$$\begin{aligned} \frac{dx}{dt} &= u \\ y &= x \end{aligned} \quad (2.27)$$

passiv ist.

Exercise 2.3. Unter welchen Voraussetzungen an die Parameter $\sigma_0, \sigma_1, \sigma_2, r_C, r_H$ und v_0 beschreibt das *LuGre-Reibmodell* (siehe z. B. Skriptum zur VO Regelungssysteme 2 [0]) ein passives System vom Eingang Δv zum Ausgang F_R . Zur Wiederholung soll das LuGre-Reibmodell nochmals in der Form

$$\begin{aligned}\frac{d}{dt}z &= \Delta v - \frac{\text{abs}(\Delta v)}{\chi(\Delta v)}\sigma_0 z \\ F_R &= \sigma_0 z + \sigma_1 \frac{d}{dt}z + \sigma_2 \Delta v\end{aligned}\quad (2.28)$$

mit

$$\chi(\Delta v) = r_C + (r_H - r_C) \exp\left(-\left(\frac{\Delta v}{v_0}\right)^2\right) \quad (2.29)$$

angeschrieben werden.

Exercise 2.4. Zeigen Sie, dass eine nichtlineare Kennlinie $y = \psi(u)$, die die Sektorbedingung $k_1 u^2 \leq \psi(u)u \leq k_2 u^2$ erfüllt, k_1 -eingangspassiv und $(\frac{1}{k_2})$ -ausgangspassiv gemäß Definition 2.2 ist.

2.3.3 Eigenschaften Passiver Systeme

Passive Systeme haben nun die bemerkenswerte Eigenschaft, dass die Parallelschaltung und die Rückkopplung passiver Systeme, wie in Abbildung 2.3 dargestellt, wiederum passiv ist.

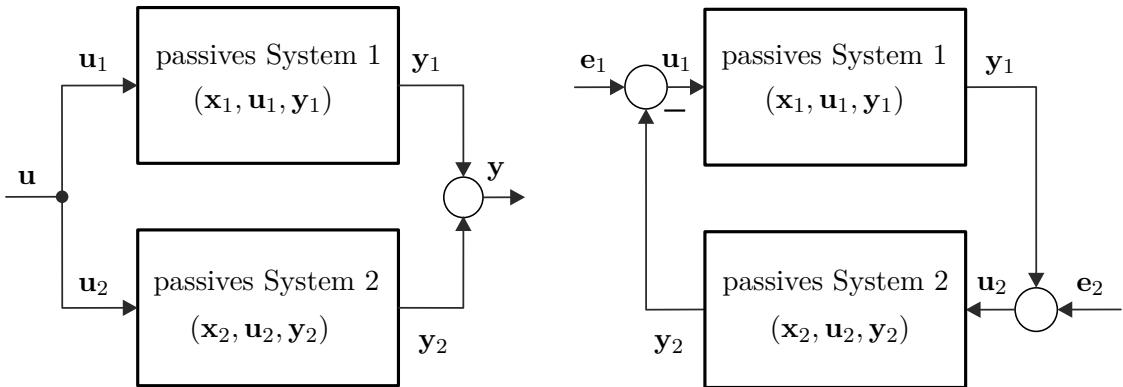


Figure 2.3: Parallelschaltung und Rückkopplung zweier passiver Systeme.

Proof. Um dies zu zeigen, nimmt man zwei passive Systeme der Form (2.19) mit $m = p$ an. Für diese existieren dann zwei nichtnegative Speicherfunktionen $V_1(\mathbf{x}_1)$

und $V_2(\mathbf{x}_2)$, die den Passivitätsungleichungen

$$\begin{aligned} V_1(\mathbf{x}_1(t)) - V_1(\mathbf{x}_1(0)) &\leq \int_0^t \mathbf{y}_1^T \mathbf{u}_1 d\tau \\ V_2(\mathbf{x}_2(t)) - V_2(\mathbf{x}_2(0)) &\leq \int_0^t \mathbf{y}_2^T \mathbf{u}_2 d\tau \end{aligned} \quad (2.30)$$

genügen. Für die Parallelschaltung nach Abbildung 2.3 gilt $\mathbf{u}_1 = \mathbf{u}_2 = \mathbf{u}$, $\mathbf{y} = \mathbf{y}_1 + \mathbf{y}_2$ und damit

$$V_1(\mathbf{x}_1(t)) + V_2(\mathbf{x}_2(t)) - V_1(\mathbf{x}_1(0)) - V_2(\mathbf{x}_2(0)) \leq \int_0^t (\mathbf{y}_1^T + \mathbf{y}_2^T) \mathbf{u} d\tau \quad (2.31)$$

bzw.

$$V(\mathbf{x}(t)) - V(\mathbf{x}(0)) \leq \int_0^t \mathbf{y}^T \mathbf{u} d\tau \quad (2.32)$$

mit der nichtnegativen Speicherfunktion $V(\mathbf{x}) = V_1(\mathbf{x}_1) + V_2(\mathbf{x}_2)$ und dem Zustand $\mathbf{x}^T = [\mathbf{x}_1^T, \mathbf{x}_2^T]$. \square

Exercise 2.5. Zeigen Sie, dass der geschlossene Kreis der Rückkopplung zweier passiver Systeme (siehe Abbildung 2.3, rechtes Bild) vom Eingang $(\mathbf{e}_1, \mathbf{e}_2)$ zum Ausgang $(\mathbf{y}_1, \mathbf{y}_2)$ passiv ist.

Darüberhinaus ist auch die Hintereinanderschaltung zweier passiver Systeme gemäß Abbildung 2.4 passiv, sofern das Verbindungssystem energieerhaltend ist, d.h. folgende Zusammenschaltungsbedingung

$$\int_0^t (\mathbf{y}_1^T \mathbf{u}_I + \mathbf{y}_2^T \mathbf{y}_I) d\tau = 0 \quad (2.33)$$

erfüllt ist.

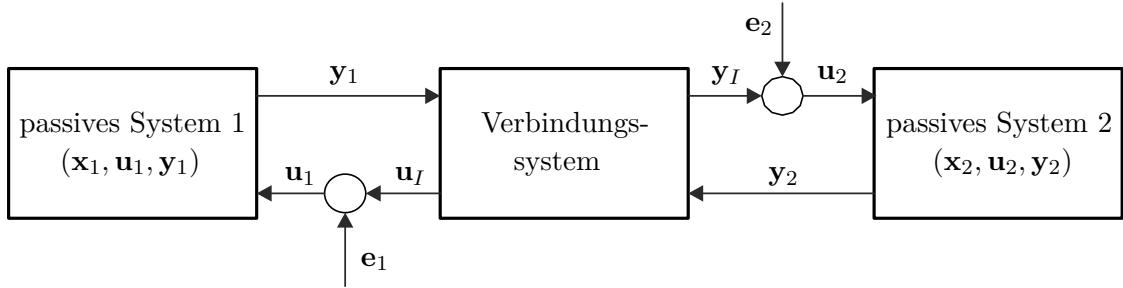


Figure 2.4: Hintereinanderschaltung passiver Systeme.

Man überzeugt sich leicht, dass dies der Fall ist, da die nachfolgende Passivitätsungleichung

$$V(\mathbf{x}(t)) - V(\mathbf{x}(0)) \leq \int_0^t (\mathbf{y}_1^T \mathbf{e}_1 + \mathbf{y}_2^T \mathbf{e}_2) d\tau \quad (2.34)$$

mit $V(\mathbf{x}) = V_1(\mathbf{x}_1) + V_2(\mathbf{x}_2)$ und $\mathbf{x}^T = [\mathbf{x}_1^T, \mathbf{x}_2^T]$ gilt. Gerade diese Eigenschaft wird bei gewissen passivitätsbasierten Reglerentwurfsverfahren genutzt, wobei das System 1 einer passiven Strecke und das System 2 einem passiven Regler entspricht. Für das Verbindungssystem wird in diesem Fall ein System der Form

$$\begin{bmatrix} \mathbf{u}_I \\ \mathbf{y}_I \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{U}_I(\mathbf{x}) \\ -\mathbf{U}_I^T(\mathbf{x}) & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \quad (2.35)$$

mit einer vorerst beliebigen quadratischen Matrix $\mathbf{U}_I(\mathbf{x})$ gewählt.

Exercise 2.6. Zeigen Sie, dass (2.35) die Zusammenschaltungsbedingung (2.33) erfüllt.

2.3.4 Passivität und Lyapunov-Stabilität

Es sei angenommen, dass das System (2.19) passiv mit einer stetig differenzierbaren, positiv definiten Speicherfunktion $V(\mathbf{x})$ ist. Dann folgt unmittelbar aus der Passivitätsungleichung (2.25) in ihrer differenziellen Form

$$\frac{d}{dt}V(\mathbf{x}) \leq \mathbf{y}^T \mathbf{u}, \quad (2.36)$$

dass die Ruhelage $\mathbf{x} = \mathbf{0}$ des freien Systems (2.19), also für $\mathbf{u} = \mathbf{0}$, stabil im Sinne von Lyapunov ist mit der Lyapunovfunktion $V(\mathbf{x})$. Ob die Ruhelage asymptotisch stabil ist, muss von Fall zu Fall mit Hilfe des Invarianzprinzips von Krassovskii-LaSalle untersucht werden.

Für die Rückkopplung zweier passiver Systeme, wie sie im rechten Teil von Abbildung 2.3 gezeigt ist, kann die asymptotische Stabilität der Ruhelage des freien geschlossenen Kreises, also für $\mathbf{e}_1 = \mathbf{e}_2 = \mathbf{0}$, auf Eigenschaften der Teilsysteme zurückgeführt werden.

Theorem 2.2. Angenommen, die Ruhelage $\mathbf{x}_1 = \mathbf{0}$ des Teilsystems 1 ist asymptotisch stabil und α -eingangspassiv gemäß Definition 2.2 mit einer stetig differenzierbaren, positiv definiten Speicherfunktion $V_1(\mathbf{x}_1)$. Weiters sei das Teilsystem 2 nullzustandsermittelbar und β -ausgangspassiv gemäß Definition 2.2 mit einer stetig differenzierbaren, positiv definiten Speicherfunktion $V_2(\mathbf{x}_2)$. Die Ruhelage des geschlossenen Kreises $(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{0}, \mathbf{0})$ ist dann asymptotisch stabil, wenn $\alpha + \beta > 0$ gilt.

Bevor dieser Satz gezeigt wird, sollen noch die Begriffe der Nullzustandsermittelbarkeit und Nullzustandsbeobachtbarkeit definiert werden.

Definition 2.3. Das System (2.19) heißt nullzustandsermittelbar (nullzustandsbeobachtbar), wenn aus $\mathbf{u}(t) = \mathbf{0}$ und $\mathbf{y}(t) = \mathbf{0}$ für alle Zeiten $t \geq 0$ folgt $\lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$ ($\mathbf{x}(t) = \mathbf{0}$ für alle Zeiten $t \geq 0$).

Proof. Zum Beweis von Satz 2.2 wähle man als Lyapunovfunktion des geschlossenen Kreises $V(\mathbf{x}) = V_1(\mathbf{x}_1) + V_2(\mathbf{x}_2)$ und bilde deren zeitliche Ableitung

$$\frac{d}{dt}V(\mathbf{x}) \leq -(\alpha + \beta)\|\mathbf{y}_2\|^2. \quad (2.37)$$

Da aber nach Satz 2.2 $\alpha + \beta > 0$ ist, folgt unmittelbar, dass die Ruhelage des geschlossenen Kreises $(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{0}, \mathbf{0})$ stabil im Sinne von Lyapunov ist. Aufgrund

der Nullzustandsermittelbarkeit des Teilsystems 2 und der asymptotischen Stabilität der Ruhelage $\mathbf{x}_1 = \mathbf{0}$ des Teilsystems 1 kann man zeigen, dass die größte positiv invariante Menge, die in $\mathcal{H} = \left\{ \mathbf{x} \in \mathcal{X} \mid \frac{d}{dt} V(\mathbf{x}) = 0 \right\}$ enthalten ist, der Ursprung $(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{0}, \mathbf{0})$ ist. Damit ist aber nach dem Invarianzprinzip von Krassovskii-LaSalle die Ruhelage des geschlossenen Kreises $(\mathbf{x}_1, \mathbf{x}_2) = (\mathbf{0}, \mathbf{0})$ asymptotisch stabil. \square

Satz 2.2 wird im Zusammenhang mit dem Begriff der *absoluten Stabilität* benötigt, insbesondere zur Herleitung des *Kreis- und Popov-Kriteriums*.

2.4 Lineare passive Systeme

Für ein lineares zeitinvariantes System der Form

$$\begin{aligned} \frac{d}{dt} \mathbf{x} &= \mathbf{A}\mathbf{x} + \mathbf{b}u \\ y &= \mathbf{c}^T \mathbf{x} + du \end{aligned} \quad (2.38)$$

lässt sich die Eigenschaft der Passivität auch an Hand der zugehörigen Übertragungsfunktion

$$G(s) = \frac{\hat{y}(s)}{\hat{u}(s)} = \mathbf{c}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{b} + d \quad (2.39)$$

beurteilen. Ohne Einschränkung der Allgemeinheit werden hier nur Eingrößensysteme behandelt, für Mehrgrößensysteme sei auf die am Ende angeführte Literatur verwiesen. Gemäß Definition 2.2 ist das System (2.38) genau dann passiv, wenn folgende Ungleichung

$$\int_0^t y u d\tau \geq 0 \quad (2.40)$$

erfüllt ist. Damit lässt sich folgender Satz für die Passivität linearer zeitinvarianter Eingrößensysteme angeben:

Theorem 2.3. Das lineare zeitinvariante System (2.38) mit der Übertragungsfunktion $G(s)$ von (2.39) ist

(1) genau dann passiv, wenn gilt

$$\operatorname{Re}(G(i\omega)) \geq 0 \quad \text{für alle } \omega, \quad (2.41)$$

(2) genau dann α -eingangspassiv mit $\alpha > 0$, wenn gilt

$$\operatorname{Re}(G(i\omega)) \geq \alpha > 0 \quad \text{für alle } \omega \quad (2.42)$$

(3) und genau dann β -ausgangspassiv mit $\beta > 0$, wenn gilt

$$\operatorname{Re}(G(i\omega)) \geq \beta |G(i\omega)|^2 > 0 \quad \text{für alle } \omega. \quad (2.43)$$

Man beachte, dass die Überprüfung der Bedingungen (2.41) - (2.43) sehr einfach an Hand der Nyquist-Ortskurve von $G(s)$ möglich ist.

Proof. Zum Beweis dieses Satzes benötigt man das so genannte *Theorem von Parseval*. Bezeichnen $x(t)$ und $y(t)$ zwei quadratisch integrierbare Zeitfunktionen, also $x(t), y(t) \in L_2(-\infty, \infty)$, und

$$\hat{x}(\omega) = \int_{-\infty}^{\infty} x(t) \exp(-I\omega t) dt \quad \text{bzw.} \quad \hat{y}(\omega) = \int_{-\infty}^{\infty} y(t) \exp(-I\omega t) dt \quad (2.44)$$

seien die zugehörigen Fouriertransformierten, dann gilt für das innere Produkt

$$\int_{-\infty}^{\infty} x(t)y(t) dt = \langle x, y \rangle = \langle \hat{x}, \hat{y} \rangle = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{x}(\omega) \hat{y}^*(\omega) d\omega . \quad (2.45)$$

Aus (2.45) folgt dann unmittelbar die Beziehung

$$\|x\|_2 = \|\hat{x}\|_2 . \quad (2.46)$$

Um das Theorem von Parseval für den Beweis von Satz 2.3 anwenden zu können, wird der Abschneideoperator $(\cdot)_T$ in der Form

$$u_T(t) = \begin{cases} u(t) & \text{für } t \leq T \\ 0 & \text{für } t > T \end{cases} \quad (2.47)$$

eingeführt. Weiters wird angenommen, dass die Zeitfunktionen $u(t)$ und $y(t)$ kausal sind, d.h. $u(t) = 0$ und $y(t) = 0$ für $t < 0$. Damit erhält man

$$\int_0^T u(t)y(t) dt = \int_{-\infty}^{\infty} u_T(t)y(t) dt = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{u}_T(\omega) \hat{y}^*(\omega) d\omega \quad (2.48)$$

bzw. mit $\hat{y}(\omega) = G(I\omega)\hat{u}_T(\omega)$ ergibt sich

$$\begin{aligned} \int_0^T u(t)y(t) dt &= \frac{1}{2\pi} \int_{-\infty}^{\infty} G^*(I\omega)\hat{u}_T(\omega)\hat{u}_T^*(\omega) d\omega \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} (\text{Re}(G(I\omega)) - I\text{Im}(G(I\omega)))|\hat{u}_T(\omega)|^2 d\omega . \end{aligned} \quad (2.49)$$

Da die linke Seite von (2.49) rein reell ist, muss der Imaginärteil auf der rechten Seite verschwinden, und es gilt

$$\int_0^T u(t)y(t) dt = \frac{1}{2\pi} \int_{-\infty}^{\infty} \text{Re}(G(I\omega))|\hat{u}_T(\omega)|^2 d\omega . \quad (2.50)$$

” \Leftarrow ”: Setzt man nun voraus, dass (2.42) gilt, dann folgt

$$\int_0^T u(t)y(t) dt \geq \frac{\alpha}{2\pi} \int_{-\infty}^{\infty} |\hat{u}_T(\omega)|^2 d\omega = \alpha \int_0^T u^2(t) dt \quad (2.51)$$

und damit nach Definition 2.2 die α -Eingangspassivität von (2.38).

” \Rightarrow ”: Umgekehrt, wenn das System (2.38) α -eingangspassiv ist, dann existiert ein $\alpha > 0$ so, dass die Ungleichung

$$\int_0^T u(t)y(t)dt \geq \alpha \int_0^T u^2(t)dt \quad (2.52)$$

erfüllt ist, bzw. mit Hilfe des Theorems von Parseval erhält man

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} (\operatorname{Re}(G(\mathrm{i}\omega)) - \alpha) |\hat{u}_T(\omega)|^2 d\omega \geq 0 . \quad (2.53)$$

Die Ungleichung (2.53) ist aber nur dann für alle Eingangsgrößen $u(t)$ gültig, wenn für alle ω gilt $\operatorname{Re}(G(\mathrm{i}\omega)) \geq \alpha$. Angenommen, es existiert ein ω_0 so, dass $\operatorname{Re}(G(\mathrm{i}\omega_0)) < \alpha$ ist, dann sieht man, dass für die Eingangsgröße $u(t) = U \sin(\omega_0 t)$ und hinreichend großes T die Ungleichung (2.53) nicht erfüllt ist. Damit ist aber Punkt (2) und für $\alpha = 0$ auch Punkt (1) von Satz 2.3 bewiesen.

Exercise 2.7. Beweisen Sie Punkt (3) von Satz 2.3.

□

Als einfaches Anwendungsbeispiel soll gezeigt werden, dass der PID-Regler

$$R(s) = V \frac{1 + T_I s}{s} \frac{1 + T_D s}{1 + \alpha T_D s} \quad (2.54)$$

mit den positiven Parametern V , T_I , T_D und $0 < \alpha < 1$ passiv ist. Dazu berechne man einfach

$$\operatorname{Re}(R(\mathrm{i}\omega)) = \frac{V(T_I + T_D(1 - \alpha) + \alpha T_D^2 T_I w^2)}{1 + \alpha^2 T_D^2 w^2} > 0 . \quad (2.55)$$

Exercise 2.8. Zeigen Sie, dass ein PI-Regler passiv ist.

Exercise 2.9. Zeigen Sie, dass das lineare zeitinvariante System (2.38) mit der Übertragungsfunktion $G(s)$ von (2.39) passiv ist, wenn

$$|\arg(G(\mathrm{i}\omega))| \leq \frac{\pi}{2} . \quad (2.56)$$

Exercise 2.10. Betrachten Sie einen einschleifigen Standardregelkreis mit einer passiven Strecke $G(s)$ und einem α -eingangspassiven Regler $R(s)$ mit $\alpha > 0$. Zeigen Sie, dass der geschlossene Kreis BIBO-stabil ist.

Tip: Verwenden Sie dazu das Nyquistkriterium.

Exercise 2.11. Der Zusammenhang zwischen Strom $\hat{i}(x, s)$ und Spannung $\hat{u}(x, s)$ an der Stelle $x = 0$ und an der Stelle $x = l$ einer langen elektrischen Leitung mit dem Kapazitätsbelag c , dem Induktivitätsbelag l , dem Widerstandsbelag r und dem

Leitwertsbelag g lautet

$$\begin{bmatrix} \hat{u}(0, s) \\ \hat{i}(0, s) \end{bmatrix} = \begin{bmatrix} \cosh(\gamma(s)l) & Z_0(s) \sinh(\gamma(s)l) \\ \frac{1}{Z_0(s)} \sinh(\gamma(s)l) & \cosh(\gamma(s)l) \end{bmatrix} \begin{bmatrix} \hat{u}(l, s) \\ \hat{i}(l, s) \end{bmatrix}, \quad (2.57)$$

wobei $Z_0(s)$ den Wellenwiderstand und $\gamma(s)$ den Ausbreitungskoeffizienten

$$Z_0(s) = \sqrt{\frac{r + sl}{g + sc}} \quad \text{und} \quad \gamma(s) = \sqrt{(r + sl)(g + sc)} \quad (2.58)$$

bezeichnen. überprüfen Sie für verschiedene Lastimpedanzen $Z_L(s)$ mit

$$\hat{u}(l, s) = Z_L(s) \hat{i}(l, s) \quad (2.59)$$

die Passivität der Übertragungsfunktion $G(s) = \frac{\hat{u}(0, s)}{\hat{i}(0, s)}$.

2.5 Positive Reellheit

Bei linearen zeitinvarianten Systemen (2.38) wird an Stelle der Passivität sehr oft der Begriff der positiven Reellheit der zugehörigen Übertragungsfunktion (2.39) verwendet. Ohne Beweis sei angemerkt, dass das System (2.38) genau dann passiv ist, wenn (2.39) positiv reell ist.

Theorem 2.4. Eine Übertragungsfunktion $G(s)$ ist genau dann positiv reell, wenn

- (1) $G(s)$ keine Pole in der rechten offenen s -Halbebene besitzt,
- (2) $\operatorname{Re}(G(\mathrm{j}\omega)) \geq 0$ ist für alle ω , für die gilt, $\mathrm{j}\omega$ ist kein Pol von $G(s)$ und
- (3) wenn $s = \mathrm{j}\omega_0$ ein Pol von $G(s)$ ist, dann ist dieser einfach und für endliches ω_0 muss das Residuum

$$\lim_{s \rightarrow \mathrm{j}\omega_0} (s - \mathrm{j}\omega_0) G(s) \quad (2.60)$$

positiv und reell sein. Ist ω_0 unendlich, dann muss der Grenzwert

$$\lim_{\omega \rightarrow \infty} \frac{G(\mathrm{j}\omega)}{\mathrm{j}\omega} \quad (2.61)$$

positiv und reell sein.

Man nennt $G(s)$ streng positiv reell, wenn $G(s - \delta)$ für ein geeignetes $\delta > 0$ positiv reell ist.

Exercise 2.12. Zeigen Sie, dass die Bedingungen

- (1) die Graddifferenz zwischen Zähler- und Nennerpolynom von $G(s)$ sind $-1, 0$ oder 1 und
 - (2) $G(s)$ hat keine Nullstellen in der rechten offenen s -Halbebene
- notwendig dafür sind, dass $G(s)$ positiv reell ist.

Exercise 2.13. Sind die nachfolgenden Übertragungsfunktionen

$$G_1(s) = -(s - 3), G_2(s) = \frac{1}{s^2 + 2s + 1}, G_3(s) = \frac{s + 1}{s^2 + 1}, G_4(s) = \frac{s + 10}{(s + 1)(s + 2)} \quad (2.62)$$

positiv reell?

Wie im nachfolgenden Satz gezeigt wird, hängt die positive Reellheit einer Übertragungsfunktion $G(s)$ eng mit der Lösbarkeit eines Gleichungssystems zusammen. Für den Beweis dieses Satzes sei auf die am Ende angeführte Literatur verwiesen.

Theorem 2.5 (Kalman-Yakubovich-Popov (KYP)-Lemma). Gegeben ist das System (2.38), wobei angenommen wird, dass das Paar (\mathbf{A}, \mathbf{b}) erreichbar und das Paar $(\mathbf{c}^T, \mathbf{A})$ beobachtbar ist. Die Übertragungsfunktion (2.39) ist genau dann positiv reell (passiv), wenn ein Skalar w , ein Vektor \mathbf{m} und eine positiv definite Matrix \mathbf{P} so existieren, dass nachfolgende Bedingungen

$$\begin{aligned} \mathbf{PA} + \mathbf{A}^T \mathbf{P} &= -\mathbf{m} \mathbf{m}^T \\ \mathbf{Pb} &= \mathbf{c} - \mathbf{m}w \\ w^2 &= 2d \end{aligned} \quad (2.63)$$

erfüllt sind. Die Übertragungsfunktion (2.39) ist darüberhinaus genau dann streng positiv reell nach Satz 2.4, wenn Skalare w und $\varepsilon > 0$, ein Vektor \mathbf{m} und eine positiv definite Matrix \mathbf{P} so existieren, dass nachfolgende Bedingungen

$$\begin{aligned} \mathbf{PA} + \mathbf{A}^T \mathbf{P} &= -\mathbf{m} \mathbf{m}^T - \varepsilon \mathbf{P} \\ \mathbf{Pb} &= \mathbf{c} - \mathbf{m}w \\ w^2 &= 2d \end{aligned} \quad (2.64)$$

erfüllt sind.

Exercise 2.14. Angenommen $w, \mathbf{m}, \mathbf{P} > \mathbf{0}$ und $\varepsilon > 0$ seien Lösungen von (2.64). Zeigen Sie, dass dann im Falle $d \neq 0$ die Riccati-Gleichung

$$\mathbf{P} \left(\frac{\varepsilon}{2} \mathbf{E} + \mathbf{A} \right) + \left(\frac{\varepsilon}{2} \mathbf{E} + \mathbf{A}^T \right) \mathbf{P} + (\mathbf{c} - \mathbf{Pb}) \frac{1}{2d} (\mathbf{c}^T - \mathbf{b}^T \mathbf{P}) = \mathbf{0} \quad (2.65)$$

erfüllt ist.

Als Anwendung des KYP Lemmas betrachte man den geschlossenen Regelkreis von

Abbildung 2.5 mit der nichtlinearen passiven Strecke im Vorwärtszweig und dem streng positiv reellen Regler im Rückwärtszweig.

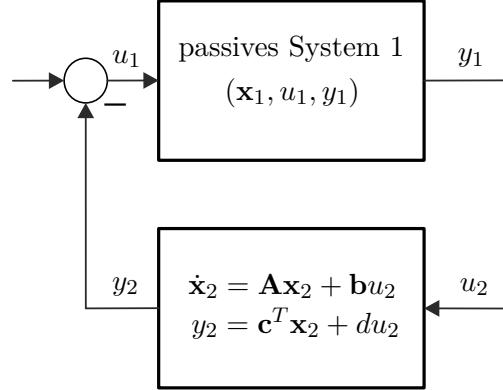


Figure 2.5: Passives System mit linearem Regler.

Angenommen das passive nichtlineare System habe eine stetig differenzierbare, positiv definite Speicherfunktion $V_1(\mathbf{x}_1)$, die der differenziellen Passivitätsungleichung (siehe (2.36))

$$\frac{d}{dt}V_1(\mathbf{x}_1) = -W_1(\mathbf{x}_1) + y_1 u_1 \leq y_1 u_1 , \quad (2.66)$$

mit der positiv semidefiniten Funktion $W_1(\mathbf{x}_1)$ genügt. Für das Weitere sei der streng positiv reelle Regler durch folgende Minimalrealisierung

$$\begin{aligned} \frac{d}{dt}\mathbf{x}_2 &= \mathbf{A}\mathbf{x}_2 + \mathbf{b}u_2 \\ y_2 &= \mathbf{c}^T \mathbf{x}_2 + du_2 \end{aligned} \quad (2.67)$$

beschrieben. Aufgrund des KYP Lemmas Satz 2.5 findet man für das System (2.67) Skalare w und $\varepsilon > 0$, einen Vektor \mathbf{m} und eine positiv definite Matrix \mathbf{P} so, dass (2.64) erfüllt ist. Damit ergibt sich die *Lyapunov-Funktion des geschlossenen Kreises* von Abbildung 2.5 zu

$$V_e(\mathbf{x}_1, \mathbf{x}_2) = V_1(\mathbf{x}_1) + \frac{1}{2}\mathbf{x}_2^T \mathbf{P} \mathbf{x}_2 . \quad (2.68)$$

Um dies zu zeigen, berechnet man die zeitliche Änderung von (2.68) entlang der Lösungskurve und berücksichtigt die Zusammenschaltungsbedingung $u_1 = -y_2$ und $u_2 = y_1$ gemeinsam mit (2.64) und (2.66)

$$\begin{aligned}
\frac{d}{dt} V_e(\mathbf{x}_1, \mathbf{x}_2) &= -W_1(\mathbf{x}_1) + y_1 u_1 + \frac{1}{2} \underbrace{\dot{\mathbf{x}}_2^T \mathbf{P} \mathbf{x}_2}_{(\mathbf{x}_2^T \mathbf{A}^T + u_2 \mathbf{b}^T) \mathbf{P} \mathbf{x}_2} + \frac{1}{2} \underbrace{\mathbf{x}_2^T \mathbf{P} \dot{\mathbf{x}}_2}_{\mathbf{x}_2^T \mathbf{P} (\mathbf{A} \mathbf{x}_2 + \mathbf{b} u_2)} \\
&= -W_1(\mathbf{x}_1) + y_1 u_1 + \frac{1}{2} \mathbf{x}_2^T \underbrace{(\mathbf{A}^T \mathbf{P} + \mathbf{P} \mathbf{A})}_{-\mathbf{m} \mathbf{m}^T - \varepsilon \mathbf{P}} \mathbf{x}_2 + \mathbf{x}_2^T \underbrace{\mathbf{P} \mathbf{b}}_{\mathbf{c} - \mathbf{m} w} u_2 \\
&= -W_1(\mathbf{x}_1) - \underbrace{y_1 \mathbf{c}^T \mathbf{x}_2}_{=} - d y_1^2 - \frac{1}{2} \mathbf{x}_2^T \mathbf{m} \mathbf{m}^T \mathbf{x}_2 - \frac{1}{2} \varepsilon \mathbf{x}_2^T \mathbf{P} \mathbf{x}_2 + \underbrace{\mathbf{x}_2^T \mathbf{c} y_1}_{=} - \mathbf{x}_2^T \mathbf{m} w y_1 \quad (2.69) \\
&= -W_1(\mathbf{x}_1) - \frac{1}{2} \varepsilon \mathbf{x}_2^T \mathbf{P} \mathbf{x}_2 - \frac{1}{2} y_1^2 \underbrace{(2d)}_{w^2} - \frac{1}{2} \mathbf{x}_2^T \mathbf{m} \mathbf{m}^T \mathbf{x}_2 - \mathbf{x}_2^T \mathbf{m} w y_1 \\
&= -W_1(\mathbf{x}_1) - \frac{1}{2} \varepsilon \mathbf{x}_2^T \mathbf{P} \mathbf{x}_2 - \frac{1}{2} (\mathbf{m}^T \mathbf{x}_2 + w y_1)^T (\mathbf{m}^T \mathbf{x}_2 + w y_1) \leq 0.
\end{aligned}$$

Dies zeigt unmittelbar die Stabilität des geschlossenen Kreises von Abbildung 2.5.

2.6 Kanonische Form Passiver Systeme

Bevor eine kanonische Form für passive Systeme vorgestellt wird, soll gezeigt werden, dass die wohlbekannten Euler-Lagrange Gleichungen passiv sind.

2.6.1 Hamiltonsche Systeme

Betrachtet man ein endlich-dimensionales Lagrangesches System mit n Freiheitsgraden und den generalisierten Koordinaten $\mathbf{q} \in \mathbb{R}^n$, dann folgen bekannterweise die Bewegungsgleichungen aus den Euler-Lagrange Gleichungen in der Form

$$\frac{d}{dt} \left(\frac{\partial L}{\partial v_k} \right) - \frac{\partial L}{\partial q_k} = \tau_k, \quad k = 1, \dots, n \quad (2.70)$$

mit der Lagrangefunktion $L(\mathbf{q}, \mathbf{v})$, den generalisierten Geschwindigkeiten $\frac{d}{dt} \mathbf{q} = \mathbf{v}$ und den generalisierten Kräften τ_k , $k = 1, \dots, n$. Bei einfachen Lagrageschen Systemen entspricht die Langrangepunktionsfunktion der Differenz aus kinetischer und potenzieller Energie

$$L(\mathbf{q}, \mathbf{v}) = T(\mathbf{q}, \mathbf{v}) - V(\mathbf{q}). \quad (2.71)$$

Es sei angenommen, dass sich die generalisierten Kräfte $\boldsymbol{\tau}$ aus externen Kräften $\boldsymbol{\tau}_e$ (Stell- und Störeingänge im Regelungstechnischen Sinne) und dissipativen Kräften $\boldsymbol{\tau}_d^T = -\left(\frac{\partial}{\partial \mathbf{v}} R\right)(\mathbf{v})$ mit der Rayleighsche Dissipationsfunktion $R(\mathbf{v})$ und

$$\left(\frac{\partial}{\partial \mathbf{v}} R \right)(\mathbf{v}) \mathbf{v} \geq 0 \quad (2.72)$$

zusammensetzen. Damit ergibt sich (2.70) zu

$$\frac{d}{dt} \left(\frac{\partial L}{\partial v_k} \right) - \frac{\partial L}{\partial q_k} + \frac{\partial}{\partial v_k} R = \tau_{e,k}, \quad k = 1, \dots, n. \quad (2.73)$$

Definition 2.4. Man bezeichnet das Lagrangesche System (2.73) *voll gedämpft*, wenn die Rayleighsche Dissipationsfunktion $R(\mathbf{v})$ folgender Ungleichung

$$\left(\frac{\partial}{\partial \mathbf{v}} R \right)(\mathbf{v}) \mathbf{v} \geq \sum_{k=1}^n \beta_k v_k^2, \quad \beta_k > 0, \quad k = 1, \dots, n \quad (2.74)$$

genügt. Ist ein $\beta_k = 0$, dann spricht man auch von einem *nicht voll gedämpften* Lagrangeschen System.

Mit Hilfe der generalisierten Impulskoordinaten

$$p_k = \frac{\partial L}{\partial v_k}, \quad k = 1, \dots, n \quad (2.75)$$

und der *Legendre-Transformation* $(q_k, v_k) \rightarrow (q_k, p_k)$ erhält man direkt aus den Euler-Lagrange Gleichungen (2.70) die äquivalenten *Hamiltonschen Gleichungen*

$$\begin{aligned} \frac{d}{dt} q_k &= \frac{\partial H}{\partial p_k} \\ \frac{d}{dt} p_k &= -\frac{\partial H}{\partial q_k} + \tau_k, \quad k = 1, \dots, n \end{aligned} \quad (2.76)$$

mit der Hamiltonfunktion

$$H(\mathbf{q}, \mathbf{p}) = \sum_{k=1}^n p_k v_k - L(\mathbf{q}, \mathbf{v}). \quad (2.77)$$

Der Satz über implizite Funktionen besagt, dass die generalisierten Geschwindigkeiten v_k aus (2.75) genau dann lokal berechnet werden können, wenn die Matrix $\left[\frac{\partial^2}{\partial v_i \partial v_j} L \right]$ regulär ist. Man spricht dann auch von einer *nichtdegenerierten Lagrangefunktion* L .

Proof. Zum Beweis betrachte man die kurzen Ableitungen

$$\frac{\partial H}{\partial p_k} = v_k + \sum_{j=1}^n \left(p_j \frac{\partial v_j}{\partial p_k} - \underbrace{\frac{\partial L}{\partial v_j} \frac{\partial v_j}{\partial p_k}}_{=p_j} \right) = v_k = \frac{d}{dt} q_k \quad (2.78)$$

und

$$\frac{\partial H}{\partial q_k} = \sum_{j=1}^n \left(p_j \frac{\partial v_j}{\partial q_k} - \underbrace{\frac{\partial L}{\partial v_j} \frac{\partial v_j}{\partial q_k}}_{=p_j} \right) - \frac{\partial L}{\partial q_k} = \tau_k - \frac{d}{dt} \left(\frac{\partial L}{\partial v_k} \right) = \tau_k - \frac{d}{dt} p_k. \quad (2.79)$$

□

Wenn die kinetische Energie $T(\mathbf{q}, \mathbf{v})$ in (2.71) die Form

$$T(\mathbf{q}, \mathbf{v}) = \frac{1}{2} \mathbf{v}^T \mathbf{D}(\mathbf{q}) \mathbf{v} \quad (2.80)$$

mit der positiv definiten Massenmatrix $\mathbf{D}(\mathbf{q})$ hat, dann entspricht die Hamiltonfunktion

$$H(\mathbf{q}, \mathbf{p}) = \sum_{k=1}^n p_k v_k - \frac{1}{2} \mathbf{v}^T \mathbf{D}(\mathbf{q}) \mathbf{v} + V(\mathbf{q}) = \frac{1}{2} \mathbf{v}^T \mathbf{D}(\mathbf{q}) \mathbf{v} + V(\mathbf{q}) \quad (2.81)$$

der im System gespeicherten Energie. Berechnet man die zeitliche Änderung der Hamiltonfunktion (2.81)

$$\frac{d}{dt} H(\mathbf{q}, \mathbf{p}) = \sum_{k=1}^n \left[\frac{\partial H}{\partial q_k} \frac{\partial H}{\partial p_k} + \underbrace{\frac{\partial H}{\partial p_k} \left(-\frac{\partial H}{\partial q_k} - \frac{\partial}{\partial v_k} R + \tau_{e,k} \right)}_{v_k} \right] \leq \sum_{k=1}^n v_k \tau_{e,k}, \quad (2.82)$$

dann sieht man, dass das Lagrangesche System gemäß Definition 2.2 passiv ist mit der Eingangsgröße $\boldsymbol{\tau}_e$, der Ausgangsgröße $\mathbf{v} = \frac{d}{dt} \mathbf{q}$ und der Speicherfunktion $H(\mathbf{q}, \mathbf{p})$. Ist darüberhinaus das Lagrangesche System gemäß Definition 2.4 voll gedämpft, dann ist das Lagrangesche System wegen (2.74) sogar β -ausgangspassiv mit $\beta = \min_k(\beta_k)$, $k = 1, \dots, n$, da gilt

$$\frac{d}{dt} H(\mathbf{q}, \mathbf{p}) \leq \sum_{k=1}^n v_k \tau_{e,k} - \sum_{k=1}^n \beta_k v_k^2 \leq \sum_{k=1}^n v_k \tau_{e,k} - \min_k(\beta_k) \|\mathbf{v}\|_2^2. \quad (2.83)$$

Man sagt dann auch, dass v_k der zur generalisierten Kraft $\tau_{e,k}$ *kollokierte Ausgang* ist. D.h., die Paarung $(\tau_{e,k}, v_k)$ beschreibt einen Energieeingang in das System, wie z.B. zusammengehörende Strom und Spannungen, Kräfte und Geschwindigkeiten oder Momente und Drehwinkelgeschwindigkeiten. Im Rahmen der Netzwerkstheorie werden solche Paarungen von Strom und Spannung, die einen Energieeingang bilden, auch als *Tor* (im Englischen *port*) bezeichnet. Die Generalisierung der Hamiltonschen Gleichungen (2.76) in Kombination mit dem Torkonzept führt direkt zur Klasse der *Port-Hamiltonschen Systeme*.

2.6.2 Port-Hamiltonsche Systeme

Ein finit-dimensionales Port-Hamiltonsches System lässt sich in der Form

$$\frac{d}{dt} \mathbf{x} = (\mathbf{J}(\mathbf{x}) - \mathbf{S}(\mathbf{x})) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \mathbf{G}_u(\mathbf{x}) \mathbf{u} \quad (2.84)$$

mit dem Zustand $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ und dem Stelleingang $\mathbf{u} \in \mathcal{U} \subset \mathbb{R}^m$ formulieren. Dabei bezeichnet $V(\mathbf{x})$, $V(\mathbf{0}) = 0$, eine stetig differenzierbare positiv definite Speicherfunktion und die Einträge der Matrizen $\mathbf{G}_u(\mathbf{x})$, $\mathbf{J}(\mathbf{x}) = -\mathbf{J}^T(\mathbf{x})$ und $\mathbf{S}(\mathbf{x}) = \mathbf{S}^T(\mathbf{x}) \geq \mathbf{0}$ seien glatte Funktionen in \mathbf{x} . Wählt man als Ausgang $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^m$ den *kollokierten Ausgang*

$$\mathbf{y} = \mathbf{G}_u^T(\mathbf{x}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T, \quad (2.85)$$

dann erkennt man unmittelbar aus der differenziellen Passivitätsungleichung

$$\frac{d}{dt} V = \mathbf{y}^T \mathbf{u} - \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T \mathbf{S}(\mathbf{x}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T \leq \mathbf{y}^T \mathbf{u}, \quad (2.86)$$

dass das System (2.84) passiv ist mit der Speicherfunktion $V(\mathbf{x})$. Die Darstellung in der Form von (2.84) erlaubt mehr als nur die einfache Feststellung der Passivität – sie ermöglicht, falls die Speicherfunktion $V(\mathbf{x})$ gleich der im System gespeicherten Gesamtentnergie ist, einen tieferen Einblick in die Energieflüsse des Systems im Inneren und mit der Systemumgebung: Die schiefsymmetrische Matrix $\mathbf{J}(\mathbf{x})$ ist nämlich mit den Energieflüssen im Systeminneren verbunden, die symmetrische, positiv semidefinite Matrix $\mathbf{S}(\mathbf{x})$ umfasst das Verhalten der dissipativen Effekte und $\mathbf{G}_u(\mathbf{x})$ beschreibt den Energieaustausch des Systems mit der Systemumgebung über die Systemtore. Wenn (2.84) keine dissipativen Elemente enthält, also $\mathbf{S}(\mathbf{x}) = \mathbf{0}$ ist, dann ist das System verlustlos bezüglich der Versorgungsrate $\mathbf{y}^T \mathbf{u}$.

Eine perfekte Aktuator/Sensor Kollokation bringt den Vorteil mit sich, dass eine einfache (zustandsabhängige) Rückführung des kollokierten Ausgangs (2.85) der Form

$$\mathbf{u} = -\mathbf{K}(\mathbf{x})\mathbf{y} = -\mathbf{K}(\mathbf{x})\mathbf{G}_u^T(\mathbf{x})\left(\frac{\partial V}{\partial \mathbf{x}}\right)^T, \quad (2.87)$$

mit der positiv definiten Matrix $\mathbf{K}(\mathbf{x}) > 0$ für alle $\mathbf{x} \in \mathcal{X}$ bei stabilen Strecken die Stabilität im geschlossenen Kreis erhält, da gilt

$$\frac{d}{dt}V = -\left(\frac{\partial V}{\partial \mathbf{x}}\right)\left(\mathbf{S}(\mathbf{x}) + \mathbf{G}_u(\mathbf{x})\mathbf{K}(\mathbf{x})\mathbf{G}_u^T(\mathbf{x})\right)\left(\frac{\partial V}{\partial \mathbf{x}}\right)^T \leq 0. \quad (2.88)$$

In der Literatur wird diese Art der Rückführung (2.87) im Zusammenhang mit Port-Hamiltonschen Systemen als *damping injection* bezeichnet oder bei allgemeinen nichtlinearen Systemen mit affinem Eingang als *Jurdjevic-Quinn Rückführung*.

Example 2.1 (Port-Hamiltonsche Darstellung des Elektromagnetventils (2.14)). Um das mathematische Modell des Elektromagnetventils (2.14) in Port-Hamiltonsche Darstellung (2.84) zu bringen, führt man die neuen Zustandsgrößen $\mathbf{x}^T = [z, p = mv, \psi_L = L(z)i_L]$ ein. Die im Magnetventil gespeicherte Energie gemäß (2.15) formuliert im neuen Zustand $[z, p, \psi_L]$

$$V = \frac{1}{2}\left(\frac{1}{L(z)}\psi_L^2 + \frac{1}{m}p^2 + cz^2\right) \quad (2.89)$$

wird in weiterer Folge als Speicherfunktion verwendet. Mit

$$\frac{\partial V}{\partial \mathbf{x}} = \begin{bmatrix} cz - \frac{1}{2}\frac{\partial L(z)}{\partial z}\frac{\psi_L^2}{L^2(z)} & \frac{p}{m} & \frac{\psi_L}{L(z)} \end{bmatrix} \quad (2.90)$$

und den Systemgleichungen (2.14) im transformierten Zustand

$$\begin{aligned} \frac{d}{dt}z &= \frac{p}{m} \\ \frac{d}{dt}p &= \left(\frac{1}{2}\frac{\partial L(z)}{\partial z}\frac{\psi_L^2}{L^2(z)} - cz - d\frac{p}{m} + F_{ext}\right) \\ \frac{d}{dt}\psi_L &= U_0 - R\frac{\psi_L}{L(z)} \end{aligned} \quad (2.91)$$

ergibt sich unmittelbar die Port-Hamiltonsche Darstellung (2.84) zu

$$\frac{d}{dt} \begin{bmatrix} z \\ p \\ \psi_L \end{bmatrix} = \left(\underbrace{\begin{bmatrix} 0 & 1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}}_{\mathbf{J}(\mathbf{x})} - \underbrace{\begin{bmatrix} 0 & 0 & 0 \\ 0 & d & 0 \\ 0 & 0 & R \end{bmatrix}}_{\mathbf{S}(\mathbf{x})} \right) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \underbrace{\begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}}_{\mathbf{G}_u(\mathbf{x})} \underbrace{\begin{bmatrix} U_0 \\ F_{ext} \end{bmatrix}}_{\mathbf{u}} . \quad (2.92)$$

Der zugehörige kollokierte Ausgang gemäß (2.85) lautet

$$\mathbf{y} = \mathbf{G}_u^T(\mathbf{x}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T = \begin{bmatrix} \psi_L \\ \frac{p}{L(z)} \\ \frac{v}{m} \end{bmatrix} = \begin{bmatrix} i_L \\ v \end{bmatrix} . \quad (2.93)$$

Exercise 2.15. Stellen Sie die mathematischen Modelle des Balkens mit rollender Kugel und des Krans mit einem Schwenkarm aus dem Skriptum zur VO Regelungssysteme 2 [0] als Port-Hamiltonsche Systeme dar.

Exercise 2.16. Stellen Sie die unterschiedlichen Gleichstrommaschinen vom Abschnitt 1.7 aus dem Skriptum zur VO Regelungssysteme 2 [0] als Port-Hamiltonsche Systeme dar.

2.7 Passivitätsbasierter Reglerentwurf

Ein mit der Port-Hamiltonschen Struktur (2.84) unmittelbar verbundenes Reglerentwurfsverfahren ist die so genannte *IDA-PBC* (*Interconnection and Damping Assignment Passivity-Based Control*). Dazu sei folgender Satz formuliert:

Theorem 2.6 (IDA-PBC). *Gegeben ist das nichtlineare System*

$$\frac{d}{dt} \mathbf{x} = \mathbf{f}(\mathbf{x}) + \mathbf{G}_u(\mathbf{x}) \mathbf{u} \quad (2.94)$$

mit dem Zustand $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n$ und dem Stelleingang $\mathbf{u} \in \mathcal{U} \subset \mathbb{R}^m$ mit $m < n$. Von der Matrix $\mathbf{G}_u(\mathbf{x})$ wird vorausgesetzt, dass diese für alle $\mathbf{x} \in \mathcal{X}$ spaltenregulär ist, d.h. $\text{rang}(\mathbf{G}_u(\mathbf{x})) = m$. Im Weiteren bezeichne $\mathbf{G}_u^\perp(\mathbf{x})$ den Linksannihilator von $\mathbf{G}_u(\mathbf{x})$, d.h. $\mathbf{G}_u^\perp(\mathbf{x}) \mathbf{G}_u(\mathbf{x}) = \mathbf{0}$, und $V_d(\mathbf{x})$ sei die Speicherfunktion des geschlossenen Kreises und habe an der gewünschten Ruhelage $\mathbf{x} = \mathbf{x}_d$ ein striktes Minimum, d.h.

$$V_d(\mathbf{x}) > V_d(\mathbf{x}_d) \quad \text{für alle } \mathbf{x} \neq \mathbf{x}_d, \quad \left(\frac{\partial V_d}{\partial \mathbf{x}} \right)(\mathbf{x}_d) = \mathbf{0} \quad \text{und} \quad \left(\frac{\partial^2 V_d}{\partial \mathbf{x}^2} \right)(\mathbf{x}_d) > 0 . \quad (2.95)$$

Damit ist $V_d(\mathbf{x}) - V_d(\mathbf{x}_d)$ positiv definit und eignet sich als Lyapunovfunktion für den geschlossenen Kreis. Angenommen die Matrizen $\mathbf{J}_d(\mathbf{x}) = -\mathbf{J}_d^T(\mathbf{x})$, $\mathbf{S}_d(\mathbf{x}) = \mathbf{S}_d^T(\mathbf{x}) \geq 0$, der Linksannihilator $\mathbf{G}_u^\perp(\mathbf{x})$ und die Speicherfunktion $V_d(\mathbf{x})$ genügen der Bedingung (PBC matching equation)

$$\mathbf{G}_u^\perp(\mathbf{x})\mathbf{f}(\mathbf{x}) = \mathbf{G}_u^\perp(\mathbf{x})(\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T, \quad (2.96)$$

dann ergibt sich mit der Zustandsrückführung

$$\mathbf{u} = \boldsymbol{\beta}(\mathbf{x}) = (\mathbf{G}_u^T(\mathbf{x})\mathbf{G}_u(\mathbf{x}))^{-1}\mathbf{G}_u^T(\mathbf{x})\left\{(\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T - \mathbf{f}(\mathbf{x})\right\} \quad (2.97)$$

eingesetzt in (2.94) ein geschlossener Kreis in Port-Hamiltonscher Form

$$\frac{d}{dt}\mathbf{x} = (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T \quad (2.98)$$

mit der stabilen gewünschten Ruhelage des geschlossenen Kreises $\mathbf{x} = \mathbf{x}_d$. Wenn die Menge $\{\mathbf{x}_d\}$ die größte positive invariante Menge von

$$\left\{ \mathbf{x} \in \mathbb{R}^n \mid \left(\frac{\partial V_d}{\partial \mathbf{x}}\right)\mathbf{S}_d(\mathbf{x})\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T = 0 \right\} \quad (2.99)$$

ist, dann ist $\mathbf{x} = \mathbf{x}_d$ sogar asymptotisch stabil.

Proof. Setzt man die rechten Seiten von (2.98) und (2.94) mit (2.97) gleich, d.h.

$$\mathbf{f}(\mathbf{x}) + \mathbf{G}_u(\mathbf{x})\boldsymbol{\beta}(\mathbf{x}) = (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T, \quad (2.100)$$

und multipliziert man mit $\mathbf{G}_u^\perp(\mathbf{x})$ von links, so erhält man unmittelbar die PBC matching equation (2.96). Die Zustandsrückführung (2.97) folgt direkt aus (2.100) durch Multiplikation mit der Pseudoinversen $(\mathbf{G}_u^T(\mathbf{x})\mathbf{G}_u(\mathbf{x}))^{-1}\mathbf{G}_u^T(\mathbf{x})$ von links. Man beachte, dass die zuvor angenommene Spaltenregularität von $\mathbf{G}_u(\mathbf{x})$ die Regularität der Pseudoinversen garantiert. Im nächsten Schritt soll gezeigt werden, dass (2.94) mit (2.97) tatsächlich (2.98) entspricht, also gilt

$$\begin{aligned} \Psi &= \mathbf{f}(\mathbf{x}) + \mathbf{G}_u(\mathbf{x})\left(\mathbf{G}_u^T(\mathbf{x})\mathbf{G}_u(\mathbf{x})\right)^{-1}\mathbf{G}_u^T(\mathbf{x})\left\{(\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T - \mathbf{f}(\mathbf{x})\right\} \\ &\quad - (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d(\mathbf{x}))\left(\frac{\partial V_d}{\partial \mathbf{x}}\right)^T = \mathbf{0}. \end{aligned} \quad (2.101)$$

Dazu multipliziere man (2.101) mit der regulären Matrix

$$\mathbf{T}(\mathbf{x}) = \begin{bmatrix} \mathbf{G}_u^T(\mathbf{x}) \\ \mathbf{G}_u^\perp(\mathbf{x}) \end{bmatrix} \quad (2.102)$$

und zufolge der PBC matching condition (2.96) folgt $\mathbf{T}(\mathbf{x})\Psi = \mathbf{0}$ und damit unmittelbar $\Psi = \mathbf{0}$. \square

Die Schwierigkeit dieser Reglerentwurfsmethode besteht offensichtlich darin, die PBC matching equation (2.96), welche ein *System partieller Differentialgleichungen* darstellt, zu lösen.

Dazu sei erwähnt, dass

- die Matrizen $\mathbf{J}_d(\mathbf{x}) = -\mathbf{J}_d^T(\mathbf{x})$ und $\mathbf{S}_d(\mathbf{x}) = \mathbf{S}_d^T(\mathbf{x}) \geq 0$ frei zu wählen sind,
- die Speicherfunktion des geschlossenen Kreises $V_d(\mathbf{x})$ abgesehen von der Bedingung (2.95) ebenfalls frei gewählt werden kann,
- und der Linksannihilator $\mathbf{G}_u^\perp(\mathbf{x})$ mit jeder regulären $(n-m) \times (n-m)$ Matrix $\Lambda(\mathbf{x})$ von links multipliziert werden kann, d.h. $\tilde{\mathbf{G}}_u^\perp(\mathbf{x}) = \Lambda(\mathbf{x})\mathbf{G}_u^\perp(\mathbf{x})$, ohne die PBC matching equation (2.96) zu ändern. Die Matrix $\Lambda(\mathbf{x})$ stellt somit einen weiteren Entwurfsfreiheitsgrad dar.

In den letzten Jahren haben sich im Wesentlichen folgende Varianten des IDA-PBC Entwurfsverfahrens durchgesetzt:

- *Non-Parametrized IDA-PBC*: In diesem Fall wird die Struktur der Zusammenschaltung in Form der Matrizen $\mathbf{J}_d(\mathbf{x}) = -\mathbf{J}_d^T(\mathbf{x})$ und $\mathbf{S}_d(\mathbf{x}) = \mathbf{S}_d^T(\mathbf{x}) \geq 0$ vorgegeben. Mit bekanntem $\mathbf{G}_u^\perp(\mathbf{x})$ resultiert die PBC matching equation (2.96) zu einer partiellen Differentialgleichung für die Speicherfunktion $V_d(\mathbf{x})$. Aus der Familie aller Lösungen müssen dann jene extrahiert werden, die die Bedingung (2.95) erfüllen. In der Literatur, siehe beispielsweise [0], findet man auch Bedingungen für die Existenz einer Lösung der zugrundeliegenden partiellen Differentialgleichung (2.96).
- *Algebraic IDA-PBC*: In diesem Fall wird die Speicherfunktion $V_d(\mathbf{x})$ unter der Bedingung (2.95) festgelegt und die PBC matching equation (2.96) degeneriert zu einer algebraischen Gleichung für die Bestimmung der Matrizen $\mathbf{J}_d(\mathbf{x}) = -\mathbf{J}_d^T(\mathbf{x})$ und $\mathbf{S}_d(\mathbf{x}) = \mathbf{S}_d^T(\mathbf{x}) \geq 0$.
- *Parametrized IDA-PBC*: Hier wird die Speicherfunktion $V_d(\mathbf{x})$ auf eine bestimmte Klasse eingeschränkt, beispielsweise bei mechanischen Systemen, dass die gewünschte potenzielle Energie nur von den generalisierten Lagekoordinaten abhängt und die gewünschte kinetische Energie eine quadratische Form in den generalisierten Geschwindigkeiten ist. Diese spezielle Form von $V_d(\mathbf{x})$ impliziert eine neue PBC matching equation mit Einschränkungen bezüglich der Wahl von $\mathbf{J}_d(\mathbf{x}) = -\mathbf{J}_d^T(\mathbf{x})$ und $\mathbf{S}_d(\mathbf{x}) = \mathbf{S}_d^T(\mathbf{x}) \geq 0$.

Example 2.2. Als Anwendungsbeispiel betrachte man eine *permanentmagnetisch erregte Synchronmaschine* in *dq*-Darstellung

$$\begin{aligned} L_d \frac{d}{dt} i_d &= -R_s i_d + \omega L_q i_q + u_d \\ L_q \frac{d}{dt} i_q &= -R_s i_q - \omega(L_d i_d + \Phi) + u_q \\ J \frac{d}{dt} \omega &= p((L_d - L_q)i_d i_q + \Phi i_q) - \tau_l \end{aligned} \quad (2.103)$$

mit den Statorströmen i_d und i_q sowie der Drehwinkelgeschwindigkeit des Rotors ω als Zustandsgrößen, den Statorspannungen u_d und u_q als Stellgrößen und dem Lastmoment τ_l . Im Weiteren bezeichnet J das Trägheitsmoment, R_s den Statorwicklungswiderstand, L_d und L_q die Statorinduktivitäten, p die Polpaarzahl und Φ den Fluss des Permanentmagneten im Rotor. Es sei an dieser Stelle erwähnt, dass für den Fall eines *gleichförmigen Luftspaltes* gilt $L_d = L_q = L$ und sich damit das mathematische Modell (2.103) entsprechend vereinfacht.

Wählt man nun als Zustandsgrößen $\mathbf{x}^T = [x_1, x_2, x_3] = [L_d i_d, L_q i_q, J\omega/p]$, dann lässt sich (2.103) in Form eines Port-Hamiltonschen Systems

$$\frac{d}{dt} \mathbf{x} = (\mathbf{J}(\mathbf{x}) - \mathbf{S}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \mathbf{G}_u \mathbf{u} + \mathbf{g}_d \tau_l \quad (2.104)$$

mit der Energiefunktion als Speicherfunktion

$$V(\mathbf{x}) = \frac{1}{2L_d} x_1^2 + \frac{1}{2L_q} x_2^2 + \frac{p}{2J} x_3^2 \quad (2.105)$$

und

$$\mathbf{J}(\mathbf{x}) = \begin{bmatrix} 0 & 0 & x_2 \\ 0 & 0 & -(x_1 + \Phi) \\ -x_2 & x_1 + \Phi & 0 \end{bmatrix}, \quad \mathbf{S} = \begin{bmatrix} R_s & 0 & 0 \\ 0 & R_s & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (2.106)$$

sowie

$$\mathbf{G}_u = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad \mathbf{g}_d = \begin{bmatrix} 0 \\ 0 \\ -1/p \end{bmatrix} \quad \text{und} \quad \mathbf{u} = \begin{bmatrix} u_d \\ u_q \end{bmatrix} \quad (2.107)$$

schreiben.

Exercise 2.17. Zeigen Sie die Gültigkeit von (2.104).

Es soll nun mit Hilfe der *Non-Parametrized IDA-PBC* eine Zustandsrückführung gemäß

Satz 2.6 so entworfen werden, dass der stationäre Arbeitspunkt

$$\mathbf{x}_d^T = [0, x_{2,d}, x_{3,d}] \quad \text{mit} \quad x_{2,d} = \frac{\bar{\tau}_l L_q}{\Phi p} \quad (2.108)$$

für ein konstantes Moment $\bar{\tau}_l$ und eine gewünschte Drehinkelgeschwindigkeit $\omega_d = x_{3,d}p/J$ stabilisiert wird. Die Struktur des geschlossenen Kreises $\mathbf{J}_d(\mathbf{x})$ und \mathbf{S}_d wird nun entsprechend einer Maschine mit gleichförmigem Luftspalt gewählt, d.h., es gilt $L_d = L_q = L$.

Exercise 2.18. Zeigen Sie, dass für $L_d = L_q = L$ die Matrizen $\mathbf{J}_d(\mathbf{x})$ und \mathbf{S}_d des zu (2.103) zugehörigen Port-Hamiltonschen Systems folgende Struktur aufweisen

$$\mathbf{J}_d(\mathbf{x}) = \begin{bmatrix} 0 & \frac{Lp}{J}x_3 & 0 \\ -\frac{Lp}{J}x_3 & 0 & -\Phi \\ 0 & \Phi & 0 \end{bmatrix} \quad \text{und} \quad \mathbf{S}_d = \mathbf{S}. \quad (2.109)$$

Die PBC matching equation (2.96) lautet dann

$$(\mathbf{J}(\mathbf{x}) - \mathbf{S}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \mathbf{G}_u \beta(\mathbf{x}) + \mathbf{g}_d \bar{\tau}_l = (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d) \left(\frac{\partial V_d}{\partial \mathbf{x}} \right)^T \quad (2.110)$$

bzw. mit dem Linksannihilator von \mathbf{G}_u

$$\mathbf{G}_u^\perp = [0, 0, 1] \quad (2.111)$$

und den Größen $V_a(\mathbf{x}) = V_d(\mathbf{x}) - V(\mathbf{x})$ sowie

$$\mathbf{J}_a(\mathbf{x}) = \mathbf{J}_d(\mathbf{x}) - \mathbf{J}(\mathbf{x}) = \begin{bmatrix} 0 & \frac{Lp}{J}x_3 & -x_2 \\ -\frac{Lp}{J}x_3 & 0 & x_1 \\ x_2 & -x_1 & 0 \end{bmatrix} \quad (2.112)$$

ergibt sich

$$\mathbf{G}_u^\perp (\mathbf{J}(\mathbf{x}) - \mathbf{S}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \mathbf{G}_u^\perp \mathbf{g}_d \bar{\tau}_l = \mathbf{G}_u^\perp (\mathbf{J}(\mathbf{x}) + \mathbf{J}_a(\mathbf{x}) - \mathbf{S}) \left(\left(\frac{\partial V_a}{\partial \mathbf{x}} \right)^T + \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T \right) \quad (2.113)$$

bzw.

$$-\mathbf{G}_u^\perp \mathbf{J}_a(\mathbf{x}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T + \mathbf{G}_u^\perp \mathbf{g}_d \bar{\tau}_l = \mathbf{G}_u^\perp (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}) \left(\frac{\partial V_a}{\partial \mathbf{x}} \right)^T. \quad (2.114)$$

Die Auswertung von (2.114) resultiert in folgender partieller Differentialgleichung

$$-\frac{x_2 x_1}{L_d} + \frac{x_2 x_1}{L_q} - \frac{1}{p} \bar{\tau}_l = \Phi \frac{\partial V_a}{\partial x_2}, \quad (2.115)$$

deren allgemeine Lösung sich wie folgt

$$V_a(x_1, x_2, x_3) = \alpha_1 \left(\frac{1}{2} x_2^2 x_1 \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) - \frac{x_2}{\Phi p} \bar{\tau}_l \right) + \psi(x_1, x_3) \quad (2.116)$$

mit dem positiven Parameter α_1 und einer noch zu wählenden Funktion $\psi(x_1, x_3)$ darstellen lässt. Damit besitzt die Speicherfunktion des geschlossenen Kreises $V_d = V + V_a$ folgende Struktur

$$V_d = \frac{1}{2L_d}x_1^2 + \frac{1}{2L_q}x_2^2 + \frac{p}{2J}x_3^2 + \frac{1}{2}\alpha_1 x_2^2 x_1 \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) - \alpha_1 \frac{x_2}{\Phi p} \bar{\tau}_l + \psi(x_1, x_3). \quad (2.117)$$

Die Aufgabe besteht nun darin, die Funktion $\psi(x_1, x_3)$ so festzulegen, dass die Bedingungen (2.95) erfüllt werden. Man kann sich nun einfach überzeugen, dass der Ansatz

$$\psi(x_1, x_3) = -\frac{1}{2}\alpha_1 \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) x_1 x_{2,d}^2 + \frac{\alpha_2}{2} x_1^2 - \frac{p}{2J} x_3^2 + \frac{\alpha_3}{2} (x_3 - x_{3,d})^2 - \frac{1}{2L_q} x_{2,d}^2 \quad (2.118)$$

mit den positiven Entwurfsparametern α_1 , α_2 und α_3 diese Bedingungen erfüllt. Dazu berechnet man für

$$V_d = \left(\frac{1}{2L_d} + \frac{\alpha_2}{2} \right) x_1^2 + \left(\frac{1}{2L_q} + \frac{\alpha_1}{2} x_1 \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) \right) (x_2^2 - x_{2,d}^2) - \frac{\alpha_1}{L_q} x_2 x_{2,d} + \frac{\alpha_3}{2} (x_3 - x_{3,d})^2 \quad (2.119)$$

vorerst den Gradienten und wertet diesen an der Stelle $\mathbf{x} = \mathbf{x}_d$ (siehe (2.108)) aus

$$\left(\frac{\partial}{\partial \mathbf{x}} V_d \right)^T (\mathbf{x}_d) = \begin{bmatrix} \left(\frac{1}{L_d} + \alpha_2 \right) x_{1,d} \\ \left(\frac{1}{L_q} + \alpha_1 x_{1,d} \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) \right) x_{2,d} - \frac{\alpha_1}{L_q} x_{2,d} \\ 0 \end{bmatrix}. \quad (2.120)$$

Offensichtlich ist für $\alpha_1 = 1$ die Forderung $\left(\frac{\partial}{\partial \mathbf{x}} V_d \right)^T (\mathbf{x}_d) = \mathbf{0}$ erfüllt. Um zu gewährleisten, dass \mathbf{x}_d ein striktes lokales Minimum von V_d ist, muss im Weiteren

$$\left(\frac{\partial^2}{\partial \mathbf{x}^2} V_d \right) (\mathbf{x}_d) = \begin{bmatrix} \frac{1}{L_d} + \alpha_2 & \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) x_{2,d} & 0 \\ \left(\frac{L_d - L_q}{L_d L_q \Phi} \right) x_{2,d} & \frac{1}{L_q} & 0 \\ 0 & 0 & \alpha_3 \end{bmatrix} \quad (2.121)$$

positiv definit sein, was durch geeignete Wahl der Parameter $\alpha_2 > 0$ und $\alpha_3 > 0$ mit

$$\frac{1}{L_d} + \alpha_2 > 0 \quad \text{und} \quad \left(\frac{1}{L_d} + \alpha_2 \right) \frac{1}{L_q} - \left(\frac{L_d - L_q}{L_d L_q \Phi} \right)^2 x_{2,d}^2 > 0 \quad (2.122)$$

sichergestellt wird. Die Zustandsrückführung errechnet sich dann gemäß (2.97) in der Form

$$\boldsymbol{\beta}(\mathbf{x}) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \left\{ (\mathbf{J}_d(\mathbf{x}) - \mathbf{S}_d) \left(\frac{\partial V_d}{\partial \mathbf{x}} \right)^T - (\mathbf{J}(\mathbf{x}) - \mathbf{S}) \left(\frac{\partial V}{\partial \mathbf{x}} \right)^T - \mathbf{g}_d \Phi p x_{2,d} \right\}. \quad (2.123)$$

Exercise 2.19. Bestimmen Sie die expliziten Ausdrücke des Zustandsregelgesetzes (2.123).

3 Iterative learning control

A large number of physical or industrial processes operate in an repetitive fashion, whereby the same task is performed over and over again under identical or at least very similar conditions. *Iterative learning control* (ILC) is based on the notion that the performance of such repetitive processes can be improved by learning from previous trials (i.e., iterations). Typical application examples include various robotic pick-and-place tasks, switched operation of electronic or optical systems or batched manufacturing processes.

Learning on observed information from previous trials is a very general idea and holds true for many every-day activities. For example, a basketball player throwing a ball repeatedly towards the hoop from a fixed position can improve his or her success rate over time. By observing the trajectories of the ball, the player obtains information on how to subsequently modify the throwing motion such that the future outcome improves. Iterative learning strategies can thus be seen as an *adaptive open-loop control scheme* that refines its input signals during operation through repetition and learning. This way, iterative learning strategies achieve high control performance in the presence of large uncertainty - either due to any inevitable *model-plant mismatch* or under the effect of *external disturbances* - as long as its effect is (almost) identical in each trial. Conversely, a non-learning controller is not able to leverage this additional information and thus reproduces the same residual control error in each subsequent iteration.

While learning can be applied to a large variety of problems, ILC specifically considers *tracking problems*, i.e., a system \mathbf{G} is operated in a repetitive fashion whereby the input $\mathbf{u}_j(t) \in \mathcal{U}$ during the j -th iteration yields the corresponding output $\mathbf{y}_j(t) \in \mathcal{Y}$ as

$$\mathbf{y}_j(t) = \mathbf{G}[\mathbf{x}_j(0), \mathbf{u}_j(t)]. \quad (3.1)$$

Within such a setting, ILC methods typically assume that:

- (i.) every iteration ends within an identical time interval $t \in [0, t_f]$,
- (ii.) every iteration starts from an (almost) identical initial state $\mathbf{x}_j(0) \approx \mathbf{x}(0)$,
- (iii.) there exists a unique input $\mathbf{u}_d(t)$ that yields the desired output $\mathbf{y}_d(t)$.

The main goal is now to design an operator $\Psi : \mathcal{U} \times \mathcal{Y} \rightarrow \mathcal{U}$ that produces an updated input

$$\mathbf{u}_{j+1}(t) = \Psi(\mathbf{u}_j(t), \mathbf{e}_j(t)) \quad (3.2)$$

based on the previous input and the resulting tracking error $\mathbf{e}_j(t) = \mathbf{y}_d(t) - \mathbf{y}_j(t)$ such that the output sequence asymptotically converges to the desired output $\mathbf{y}_d(t)$, i.e.,

$$\lim_{j \rightarrow \infty} \mathbf{G}[\mathbf{x}(0), \mathbf{u}_j(t)] = \mathbf{y}_d(t). \quad (3.3)$$

This iterative learning process is illustrated in Fig. 3.1. In some sense, ILC can be seen as a method to increase the robustness of feedforward control against model uncertainties and repetitive disturbances by determining it online through a fixed-point iteration (3.2), whose properties depend on the chosen operator Ψ . For simplicity, we will restrict ourselves to linear operators Ψ , wherefore the learning law (3.2) can be written as

$$\mathbf{u}_{j+1}(t) = \mathbf{L}_u \mathbf{u}_j(t) + \mathbf{L}_e \mathbf{e}_j(t) = \mathbf{Q}(\mathbf{u}_j(t) + \mathbf{L} \mathbf{e}_j(t)) \quad (3.4)$$

with the linear operators \mathbf{L}_u , \mathbf{L}_e , \mathbf{Q} , and \mathbf{L} , respectively. Correspondingly, we will mainly restrict ourselves to linear, time-invariant systems of the form

$$\dot{\mathbf{x}}_j(t) = \mathbf{A}\mathbf{x}_j(t) + \mathbf{B}\mathbf{u}_j(t) \quad (3.5)$$

$$\mathbf{y}_j(t) = \mathbf{C}\mathbf{x}_j(t) + \mathbf{D}\mathbf{u}_j(t) \quad (3.6)$$

with the state vector $\mathbf{x}_j^T(t) \in \mathbb{R}^n$, the input vector $\mathbf{u}_j^T(t) \in \mathbb{R}^l$ and the output vector $\mathbf{y}_j^T(t) \in \mathbb{R}^m$.

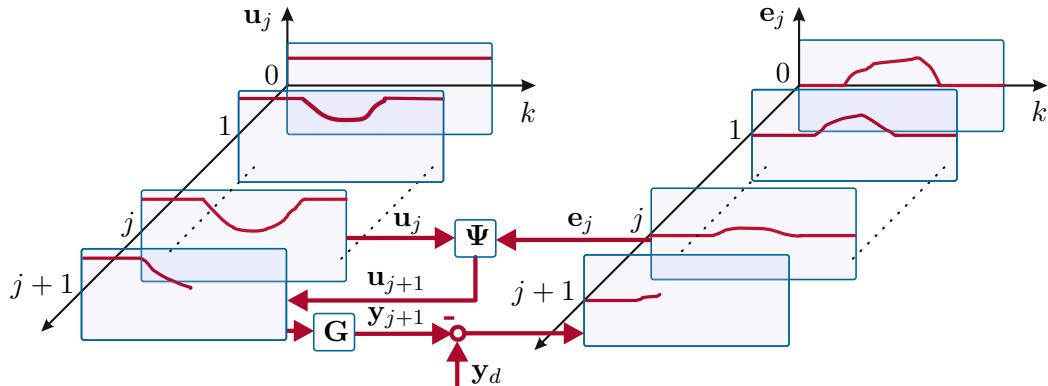


Figure 3.1: Graphical illustration of iterative learning control (ILC).

There are a number of different concepts in control theory that aim to include some element of learning. Most notably, *adaptive control* is using information from the past to increase the performance of the closed-loop system by continuously modifying model or (feedback) control parameters. All uncertainty is thus represented in a parametric form. Conversely, ILC modifies the feedforward inputs applied to the system and is thus not limited to parametric uncertainties. However, this entails that the learned input signal is specific to the desired output $\mathbf{y}_d(t)$. One property that is used to distinguish ILC [3.0] from other learning concepts [3.0] is the so-called *identical initialization condition*, i.e., that every iteration starts from the exact same initial state $\mathbf{x}_j(0) = \mathbf{x}(0)$. In a more relaxed formulation, $\mathbf{x}_j(0)$ may be a stochastic quantity, but at least is independent of previous iterations. *Repetitive control* on the other hand assumes that the initial state of the current iteration is given by the terminal state of the previous one such that the sequence of all iterations can be seen as a continuously operated system.

3.1 Fixed-point iterations

Calculating zeros of a function

$$\gamma(\mathbf{z}) = \mathbf{0} \quad (3.7)$$

with $\gamma(\mathbf{z}) : \mathbb{R}^N \rightarrow \mathbb{R}^N$ and $\mathbf{z}^T = [z_1 \dots z_N]$ can always be recast as a fixed-point problem

$$\psi(\mathbf{z}) = \mathbf{z}. \quad (3.8)$$

A *fixed-point iteration* is a sequence $\mathbf{z}_0, \mathbf{z}_1, \dots$ given by

$$\mathbf{z}_{j+1} = \psi(\mathbf{z}_j), \quad j = 0, 1, 2, \dots . \quad (3.9)$$

that converges to a solution \mathbf{z}_∞ of (3.8) depending on the chosen function ψ . Note that there is no unique fixed-point formulation for a given zero-point problem, i.e., the choices

- $\psi(\mathbf{z}) = \mathbf{z} - \gamma(\mathbf{z})$
- $\psi(\mathbf{z}) = \mathbf{z} + 2\gamma(\mathbf{z})$
- $\psi(\mathbf{z}) = \mathbf{z} - (\frac{\partial}{\partial \mathbf{z}} \gamma)^{-1}(\mathbf{z})\gamma(\mathbf{z})$

are equally valid. However, the choice of ψ crucially determines the convergence properties of the fixed-point iteration. Therefore, we define:

Definition 3.1 (Convergence). The iteration (3.9) is

- *locally convergent (LC)* to \mathbf{z}_∞ if there exists a $\delta > 0$ such that the iteration (3.9) exists and converges to \mathbf{z}_∞ for all starting points $\|\mathbf{z}_0 - \mathbf{z}_\infty\| < \delta$,
- *globally convergent (GC)* if the iteration (3.9) converges to \mathbf{z}_∞ for all \mathbf{z}_0 .

Definition 3.2 (Stability). A fixed-point \mathbf{z}_∞ is

- *stable (S)* (in a Lyapunov sense), if for every $\varepsilon > 0$ there exists a $\delta > 0$ such that if $\|\mathbf{z}_0 - \mathbf{z}_\infty\| < \delta$, the sequence $\{\mathbf{z}_j\}$ exists and $\|\mathbf{z}_j - \mathbf{z}_\infty\| < \varepsilon$ for all $j \geq 1$,
- *attractive (A)*, if there exists a $\delta > 0$ such that $\|\mathbf{z}_0 - \mathbf{z}_\infty\| < \delta$ implies that $\{\mathbf{z}_j\}$ exists and $\lim_{j \rightarrow \infty} \mathbf{z}_j = \mathbf{z}_\infty$,
- *globally attractive (GA)* if $\delta = \infty$ above,
- *asymptotically stable (AS)*, if stable and attractive and *globally asymptotically stable (GAS)*, if stable and globally attractive.

Attractivity and convergence are equivalent concepts, whereby the following relations [3.0] hold true

$$GAS \implies GA \iff GC \implies A \iff AS. \quad (3.10)$$

In the following we will consider the special case of linear fixed-point iterations

$$\mathbf{z}_{j+1} = \Psi \mathbf{z}_j , \quad j = 0, 1, 2, \dots \quad (3.11)$$

Definition 3.3 (Stability and asymptotic stability of linear fixed-point iterations). A linear iteration (3.11) is stable, if

$$\sup_{j \geq 1} \|\Psi^j\| < \infty , \quad (3.12)$$

and asymptotically stable if

$$\lim_{j \rightarrow \infty} \|\Psi^j\| = 0 . \quad (3.13)$$

Definition 3.4 (Spectral radius). The spectrum of a matrix Γ denotes the set of all its eigenvalues, i.e.,

$$\sigma(\Gamma) = \{\lambda \in \mathbb{C} \mid \det(\lambda \mathbf{I} - \Gamma) = 0\} \quad (3.14)$$

and

$$\rho(\Gamma) = \max_{\lambda \in \sigma(\Gamma)} |\lambda| \quad (3.15)$$

denotes the spectral radius of Γ .

Theorem 3.1. A linear iteration (3.11) is stable iff $\rho(\Psi) \leq 1$ and all eigenvalues at 1 are distinct eigenvalues, i.e., their algebraic multiplicity equals 1. The iteration is further asymptotically stable iff $\rho(\Psi) < 1$.

Definition 3.5 (BIBO stability). A linear iteration

$$\mathbf{z}_{j+1} = \Psi \mathbf{z}_j + \Lambda \mathbf{v}_j , \quad \mathbf{z}_0 = \mathbf{0} \quad (3.16)$$

is called BIBO stable if every bounded input sequence $\{\mathbf{v}_j\}$ results in a bounded output sequence $\{\mathbf{z}_j\}$.

Theorem 3.2. A linear iteration $\mathbf{z}_{j+1} = \Psi \mathbf{z}_j + \Lambda \mathbf{v}_j , \mathbf{z}_0 = \mathbf{0}$ is BIBO stable iff $\rho(\Psi) < 1$.

Lemma 3.1. For a bounded sequence $\{\mathbf{z}_j\}$ and $\varepsilon > 0 \in \mathbb{R}$ with

$$\|\mathbf{z}_{j+1}\| \leq \rho \|\mathbf{z}_j\| + \varepsilon , \quad 0 \leq \rho < 1 \quad (3.17)$$

it follows that

$$\limsup_{j \rightarrow \infty} \|\mathbf{z}_j\| \leq \frac{1}{1-\rho} \varepsilon . \quad (3.18)$$

Proof. Iterative application of the iteration (3.17) yields

$$\begin{aligned} \|\mathbf{z}_1\| &\leq \rho \|\mathbf{z}_0\| + \varepsilon \\ \|\mathbf{z}_2\| &\leq \rho^2 \|\mathbf{z}_0\| + (1+\rho)\varepsilon \\ &\vdots \\ \|\mathbf{z}_j\| &\leq \rho^j \|\mathbf{z}_0\| + \sum_{j=0}^{j-1} \rho^j \varepsilon = \rho^j \|\mathbf{z}_0\| + \frac{1-\rho^j}{1-\rho} \varepsilon . \end{aligned} \quad (3.19)$$

Since $0 \leq \rho < 1$, one directly obtains (3.18) for $j \rightarrow \infty$. \square

While asymptotic stability ensures that an iteration (3.11) eventually converges, this is not necessarily happening in a monotonic way. Since the notion of iterative learning somewhat suggests that one is successively progressing towards a desired solution, i.e., that the control performance improves with every iteration, monotonicity is an important property of ILC algorithms.

Definition 3.6 (Largest singular value). The largest singular value of a matrix Ψ is given by

$$\bar{\sigma}(\Psi) = \sqrt{\rho(\Psi^T \Psi)} . \quad (3.20)$$

Our main interest in the largest singular value of a matrix stems from the fact that it is the induced norm of a matrix mapping two spaces equipped with the Euclidean norm $\|z\| = \sqrt{\mathbf{z}^T \mathbf{z}}$, i.e.,

$$\|\Psi \mathbf{z}_j\| \leq \|\Psi\|_{i,2} \|\mathbf{z}_j\| = \bar{\sigma}(\Psi) \|\mathbf{z}_j\| . \quad (3.21)$$

The largest singular value can thus be seen as an upper bound of the gain or amplification of the mapping given by the matrix Ψ .

Theorem 3.3 (Monotone convergence of linear iterations). *A linear iteration $\mathbf{z}_{j+1} = \Psi \mathbf{z}_j$ converges monotonically towards $\mathbf{0}$ in the l_2 -norm, i.e., it holds true that*

$$\|\mathbf{z}_{j+1}\| \leq \beta \|\mathbf{z}_j\| \quad \text{and thus} \quad \|\mathbf{z}_{j+1}\| \leq \beta^j \|\mathbf{z}_0\| \quad (3.22)$$

for $0 \leq \beta < 1$, if

$$\bar{\sigma}(\Psi) < 1 . \quad (3.23)$$

Note that since $\rho(\Psi) \leq \bar{\sigma}(\Psi)$, monotonic convergence unsurprisingly implies convergence.

3.2 Signals, systems, and the frequency domain

Integral transforms are essential tools to analyze temporal signals and their interaction with dynamical systems. Here, we consider a *signal* as a real-valued (measureable) function that maps the real numbers \mathbb{R} to \mathbb{R}^n . The set of all signals

$$\mathcal{S} = \{\mathbf{f} : \mathbb{R} \rightarrow \mathbb{R}^n\} \quad (3.24)$$

can intuitively be conceived as a set that contains all signals that can possibly occur in an engineering system - and many more. Since the sum of two signals \mathbf{f} and \mathbf{g} and the product of a signal with a scalar are again contained in this set, signals form a natural *linear vector space*.

Definition 3.7. A *linear vector space* over a field \mathbb{K} is a set \mathcal{S} with a binary operation $+ : \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{S}$ (addition) and a binary operation $\cdot : \mathbb{K} \times \mathcal{S} \rightarrow \mathcal{S}$ (multiplication) that fulfills

1. The set \mathcal{S} and the operation $+$ are a commutative group.
2. Multiplication with scalars $a, b \in \mathbb{K}$ for $\mathbf{f}, \mathbf{g} \in \mathcal{S}$ satisfies
 - $a(\mathbf{f} + \mathbf{g}) = a\mathbf{f} + a\mathbf{g}$
 - $(a + b)\mathbf{f} = a\mathbf{f} + b\mathbf{f}$
 - $(ab)\mathbf{f} = a(b\mathbf{f})$
 - $0\mathbf{f} = \mathbf{0}$ and $1\mathbf{f} = \mathbf{f}$.

To measure the size of a signal, one typically equips vector spaces with a suitable norm.

Definition 3.8. A *normed linear vector space* \mathcal{S} is a linear vector space equipped with real-valued function $\|\cdot\|_p : \mathcal{S} \rightarrow \mathbb{R}$ that adheres to the following properties:

1. $\|\mathbf{f}\|_p \geq 0$
2. $\|\mathbf{f}\|_p = 0 \iff \mathbf{f} = \mathbf{0}$
3. $\|a\mathbf{f}\|_p = |a|\|\mathbf{f}\|_p$
4. $\|\mathbf{f} + \mathbf{g}\|_p \leq \|\mathbf{f}\|_p + \|\mathbf{g}\|_p$

Note that a vector space can be equipped with different norms. The resulting normed vector spaces are different and may contain different elements.

We will now introduce typical signal spaces that we will require to analyze ILC algorithms. Since we consider signals with values in some \mathbb{R}^n , this implicitly requires a notion of size in this underlying vector space, for which we simply use the Euclidean norm. One of the most intuitive candidates for an ILC setting is the finite-horizon 2-norm defined by

$$\|\mathbf{f}\|_{2,[0,t_f]} = \left[\int_0^{t_f} \|\mathbf{f}(t)\| dt \right]^{\frac{1}{2}} = \left\{ \int_0^{t_f} \mathbf{f}(t)^T \mathbf{f}(t) dt \right\}^{\frac{1}{2}}. \quad (3.25)$$

The set of signals for which this norm is finite is known as the finite-horizon Lebesgue 2-space

$$\mathcal{L}_2([0, t_f]; \mathbb{R}^n) = \left\{ \mathbf{f} \in \mathcal{S} : \|\mathbf{f}\|_{2,[0,t_f]} < \infty \right\}. \quad (3.26)$$

For simplicity, we will simply write $\mathcal{L}_2([0, t_f])$ in cases where the dimension of the vector-valued signals is not important. Any signal that is continuous on $[0, t_f]$ is bounded and thus contained in $\mathcal{L}_2([0, t_f])$ but signals like $\frac{1}{|2t-t_f|}$ are not. In order to address stability issues or to apply frequency-domain methods, one must consider signals over infinite time intervals. Extending the considered horizon to infinity on both sides yields the usual infinite-horizon Lebesgue 2-space

$$\mathcal{L}_2(\mathbb{R}) = \left\{ \mathbf{f} \in \mathcal{S} : \|\mathbf{f}\|_2 < \infty \right\} \quad (3.27)$$

with the corresponding norm

$$\|\mathbf{f}\|_2 = \left[\int_{-\infty}^{\infty} \|\mathbf{f}(t)\|^2 dt \right]^{\frac{1}{2}}. \quad (3.28)$$

The spaces $\mathcal{L}_2([0, \infty))$ and $\mathcal{L}_2((-\infty, 0])$ can be defined analogously.

For convenience, we will restrict ourselves to *Hilbert spaces*, i.e., *complete* spaces whose norm is generated by an *inner product* $\|\mathbf{f}\|_2 = \sqrt{\langle \mathbf{f}, \mathbf{f} \rangle}$.

Definition 3.9 (Inner product). A mapping $\mathcal{S} \times \mathcal{S} \rightarrow \mathbb{K}$ that assigns a scalar to each two elements of a vector space is called an *inner product* if the conditions

1. $\langle \mathbf{f} + \mathbf{g}, \mathbf{h} \rangle = \langle \mathbf{f}, \mathbf{h} \rangle + \langle \mathbf{g}, \mathbf{h} \rangle$
2. $\langle \mathbf{f}, \mathbf{g} \rangle = \langle \mathbf{g}, \mathbf{f} \rangle^*$
3. $\langle a\mathbf{f}, \mathbf{g} \rangle = a\langle \mathbf{f}, \mathbf{g} \rangle$
4. $\langle \mathbf{f}, \mathbf{g} \rangle > 0$ and $\langle \mathbf{f}, \mathbf{f} \rangle = 0 \iff \mathbf{f} = 0$

hold true for $\mathbf{f}, \mathbf{g}, \mathbf{h} \in \mathcal{S}$ and $a \in \mathbb{K}$.

The space $\mathcal{L}_2(\mathbb{R})$ is a Hilbert space with inner product defined by

$$\langle \mathbf{f}, \mathbf{g} \rangle = \int_{-\infty}^{\infty} \mathbf{g}(t)^T \mathbf{f}(t) dt. \quad (3.29)$$

Two signals \mathbf{f} and \mathbf{g} are called orthogonal if $\langle \mathbf{f}, \mathbf{g} \rangle = 0$ analogous to orthogonality in \mathbb{R}^n . The spaces $\mathcal{L}_2([0, \infty))$, $\mathcal{L}_2((-\infty, 0])$, and $\mathcal{L}_2([0, t_f])$ are all Hilbert spaces in their own right, with the inner product integral taken over the appropriate time interval, e.g., for $\mathcal{L}_2([0, t_f])$ we have

$$\langle \mathbf{f}, \mathbf{g} \rangle_{[0,t_f]} = \int_0^{t_f} \mathbf{g}(t)^T \mathbf{f}(t) dt. \quad (3.30)$$

The Fourier transform

The Fourier transform is a unique mapping of certain (real-valued) functions of time to complex-valued functions of a single real variable ω using

$$\hat{\mathbf{f}}(j\omega) = \mathcal{F}\{\mathbf{f}\} = \int_{-\infty}^{\infty} \mathbf{f}(t) e^{-j\omega t} dt. \quad (3.31)$$

Note that the Fourier transform in (3.31) is performed for each component individually and a vector-valued signal thus simply yields a vector of all transformed components. To ensure that such a function $\hat{\mathbf{f}}(j\omega)$ exists and is reasonably well behaved, classical Fourier analysis assumes that each component of $\mathbf{f}(t)$ is absolutely integrable over the real line

$$\int_{-\infty}^{\infty} |f_i(t)| dt < \infty, \quad (3.32)$$

i.e., each component is a $L_1(-\infty, \infty)$ function. With some mathematical effort, the Fourier transform can be extended to $L_2(\mathbb{R})$ (and functions beyond that such as the Dirac delta function). To motivate this, one can introduce the inner product of the transformed signals

$$\langle \hat{\mathbf{f}}, \hat{\mathbf{g}} \rangle = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{\mathbf{g}}(j\omega)^H \hat{\mathbf{f}}(j\omega) d\omega, \quad (3.33)$$

which again introduces an Lebesgue 2-space $L_2(\mathbb{R})$ in the frequency domain. Here, $(\cdot)^H$ denotes the conjugate transpose or Hermitian transpose.

Theorem 3.4 (Parseval's theorem). *For $\mathbf{f}, \mathbf{g} \in L_2(\mathbb{R})$ and $\hat{\mathbf{f}}(j\omega) = \mathcal{F}\{\mathbf{f}\}, \hat{\mathbf{g}}(j\omega) = \mathcal{F}\{\mathbf{g}\} \in L_2(\mathbb{R})$ it follows that*

$$\langle \mathbf{f}, \mathbf{g} \rangle = \langle \hat{\mathbf{f}}, \hat{\mathbf{g}} \rangle \quad (3.34)$$

and thus

$$\|\mathbf{f}\|_2 = \|\hat{\mathbf{f}}\|_2. \quad (3.35)$$

The Fourier transform is thus a mapping between two Lebesgue 2-spaces that preserves the inner product and the norm, i.e., a Hilbert space *isomorphism*. Finally, the inverse transform is given by

$$\mathbf{f}(t) = \mathcal{F}^{-1}\{\hat{\mathbf{f}}\} = \frac{1}{2\pi} \int_{-\infty}^{\infty} \hat{\mathbf{f}}(j\omega) e^{j\omega t} d\omega. \quad (3.36)$$

The bilateral Laplace transform

For reasons that are quite obvious to a control engineer, one is interested to extend the Fourier transform from the imaginary axis $j\omega$ to the complex plane $s = \sigma + j\omega$, which directly yields the bilateral Laplace transform

$$\hat{\mathbf{f}}(s) = \mathcal{B}\{\mathbf{f}\} = \int_{-\infty}^{\infty} \mathbf{f}(t) e^{-st} dt \quad (3.37)$$

that is equivalent to the Fourier transform for $\sigma = 0$. In general, the Laplace transform will not converge for arbitrary values of s . From

$$\hat{\mathbf{f}}(\sigma + j\omega) = \int_{-\infty}^{\infty} \mathbf{f}(t) e^{-\sigma t} e^{-j\omega t} dt \quad (3.38)$$

we see that the region of convergence (ROC) of the Laplace transform is directly linked to the convergence of $\mathbf{f}(t)e^{-\sigma t}$.

Theorem 3.5. *The ROC of a signal has the following properties:*

1. *The ROC consists of strips parallel to the $j\omega$ -axis in the s -plane.*
2. *If $\mathbf{f}(t)$ is of finite duration and is absolutely integrable, then the ROC is the entire s -plane.*
3. *If $\mathbf{f}(t)$ is right-sided (i.e., $\mathbf{f}(t) = 0$ for $t < \bar{t}$) and if the line $\text{Re}\{s\} = \sigma_0$ is in the ROC, then all values of s with $\text{Re}\{s\} > \sigma_0$ are also included in the ROC.*
4. *If $\mathbf{f}(t)$ is left-sided (i.e., $\mathbf{f}(t) = 0$ for $t > \bar{t}$) and if the line $\text{Re}\{s\} = \sigma_0$ is in the ROC, then all values of s with $\text{Re}\{s\} < \sigma_0$ are also included in the ROC.*
5. *If the Laplace transform $\hat{\mathbf{f}}(s)$ is rational, then its ROC is bounded by poles or extends to infinity. No poles are contained in the ROC.*

In particular, a signal that is contained in $\mathcal{L}_2([0, \infty))$ (and assuming that it is absolutely integrable for simplicity) is thus bounded and analytic (i.e., holomorphic) for $\text{Re}\{s\} > 0$, which is the definition of the Hardy 2-space \mathcal{H}_2 . Conversely, a signal that is contained in $\mathcal{L}_2((-\infty, 0])$ is bounded and analytic for $\text{Re}\{s\} < 0$, which is the definition of the complementary Hardy 2-space \mathcal{H}_2^\perp . Since for every $\hat{\mathbf{f}} \in \mathcal{H}_2$ it follows that $\lim_{\sigma \rightarrow +0} \hat{\mathbf{f}} \in \mathcal{L}_2(\mathbb{R})$ and analogous $\hat{\mathbf{f}} \in \mathcal{H}_2^\perp$ implies that $\lim_{\sigma \rightarrow -0} \hat{\mathbf{f}} \in \mathcal{L}_2(\mathbb{R})$, we regard \mathcal{H}_2 and \mathcal{H}_2^\perp closed subspaces of the frequency-domain $\mathcal{L}_2(\mathbb{R})$. Since any time signal can be decomposed as

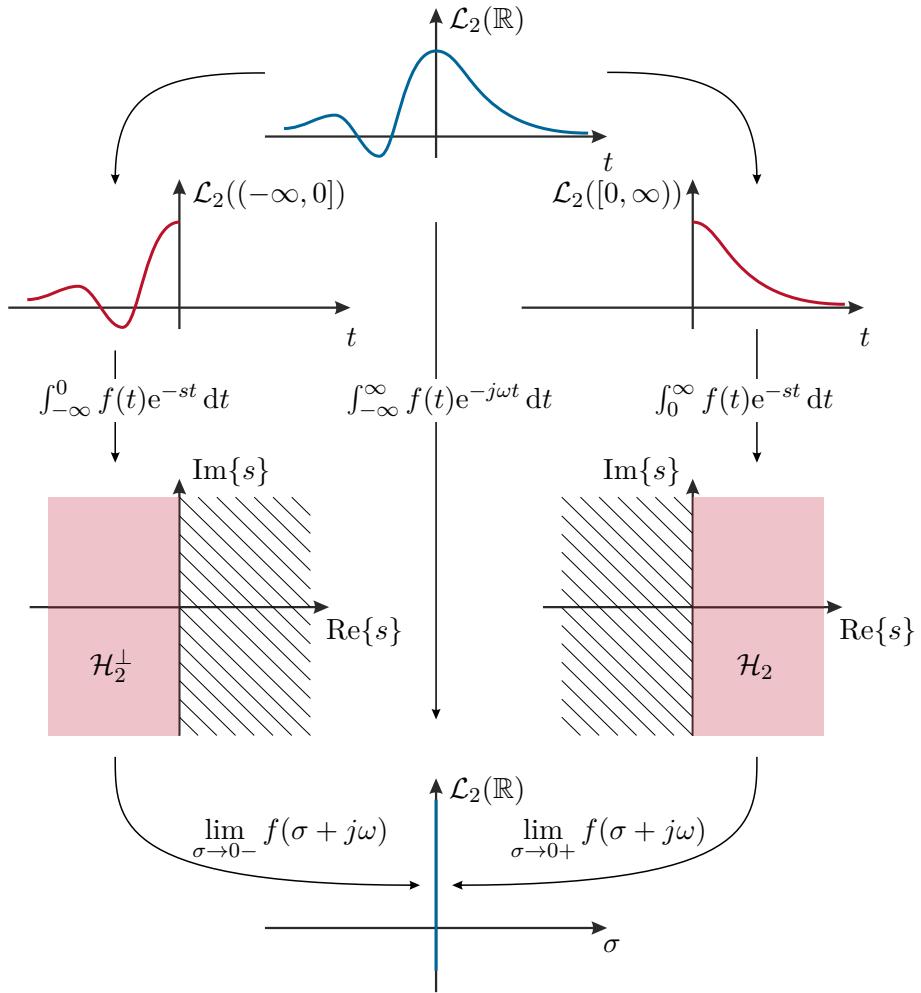
$$\mathbf{f}(t) = \mathbf{f}_1(t) + \mathbf{f}_2(t) \quad \text{with} \quad \mathbf{f}_1 \in \mathcal{L}_2([0, \infty)), \mathbf{f}_2 \in \mathcal{L}_2((-\infty, 0]) \quad (3.39)$$

and $\langle \mathbf{f}_1, \mathbf{f}_2 \rangle = 0$, one can see from $\mathcal{L}_2(\mathbb{R}) = \mathcal{H}_2 \cup \mathcal{H}_2^\perp$ and Parseval's theorem that

$$\langle \hat{\mathbf{f}}_1, \hat{\mathbf{f}}_2 \rangle = 0 \quad \text{and} \quad \mathcal{H}_2 \cap \mathcal{H}_2^\perp = 0, \quad (3.40)$$

which justifies the already used nomenclature \mathcal{H}_2 and \mathcal{H}_2^\perp , respectively. A graphical illustrations of these connections is given in Fig. 3.2. Finally, the inverse transformation to 3.37 is given by the contour integral

$$\mathbf{f}(t) = \mathcal{B}^{-1}\{\hat{\mathbf{f}}\} = \frac{1}{2\pi j} \int_{\sigma-j\infty}^{\sigma+j\infty} \hat{\mathbf{f}}(s) e^{st} ds \quad (3.41)$$

Figure 3.2: On $\mathcal{L}_2((-\infty, 0])$, $\mathcal{L}_2([0, \infty))$, \mathcal{H}_2 , and \mathcal{H}_2^\perp .

where the line of integration $\text{Re}\{s\} = \sigma$ is within the corresponding ROC of $\hat{\mathbf{f}}(s)$.

In the following, we will not make a notational distinction between time-domain signals and their frequency-domain representations and drop the hat notation unless it is necessary to avoid confusion. In general, the argument or the context sufficiently determines whether one is dealing with a time-domain or frequency-domain signal, which are anyhow isomorphic under the used integral transformations.

Input-output representation of linear systems

The input-output behavior of a dynamical system can be seen as a mapping from one signal space, the input space \mathcal{U} , to another space, the output space \mathcal{Y} :

$$\mathbf{G} : \mathbf{u} \in \mathcal{U} \mapsto \mathbf{y} = \mathbf{Gu} \in \mathcal{Y} \quad (3.42)$$

For the case $\mathcal{U} = \mathcal{L}_2(\mathbb{R}; \mathbb{R}^l)$ and $\mathcal{Y} = \mathcal{L}_2(\mathbb{R}; \mathbb{R}^m)$, the input-output behavior of any linear system can be represented by the integral operator

$$\mathbf{y}(t) = \int_{-\infty}^{\infty} \mathbf{G}(t, \tau) \mathbf{u}(\tau) d\tau \quad (3.43)$$

with the kernel function $\mathbf{G}(t, \tau) \in \mathbb{R}^{m \times l}$ that is usually referred to as the systems impulse response matrix. The system is said to be causal iff $\mathbf{G}(t, \tau) = 0$ for all $\tau > t$ and time-invariant if $\mathbf{G}(t, \tau) = \mathbf{G}(t - \tau, 0)$ for all t, τ . Linear time-invariant (LTI) systems are thus represented by a convolution-type integral

$$\mathbf{y}(t) = \int_{-\infty}^{\infty} \mathbf{G}(t - \tau) \mathbf{u}(\tau) d\tau \quad (3.44)$$

where $\mathbf{G}(t - \tau, 0)$ is denoted as $\mathbf{G}(t - \tau)$ for simplicity. A system is called *causal* if $\mathbf{G}(t) = 0$ for $t < 0$ and *anticausal* if $\mathbf{G}(t) = 0$ for $t > 0$. If a system is neither causal nor anticausal, it is called *non-causal* or *acausal*. By applying the Laplace transform (3.37) to (3.44), one obtains

$$\mathbf{y}(s) = \mathbf{G}(s) \mathbf{u}(s) \quad (3.45)$$

with the transfer function matrix

$$\mathbf{G}(s) = \int_{-\infty}^{\infty} \mathbf{G}(t) e^{-st} dt. \quad (3.46)$$

The following properties of a dynamic system are directly related to the ROC of the transfer function matrix:

Theorem 3.6 (Causality and stability). *For an LTI system $\mathbf{G}(s)$ with corresponding ROC it holds true that:*

1. *If it is causal then its ROC is a right-half plane. If $\mathbf{G}(s)$ is rational, then causality is equivalent to the ROC being the right-half plane to the right of the rightmost pole.*
2. *If it is anticausal then its ROC is a left-half plane. If $\mathbf{G}(s)$ is rational, then anticausality is equivalent to the ROC being the left-half plane to the left of the leftmost pole.*
3. *It is stable iff the ROC includes $\text{Re}\{s\} = 0$.*

Exercise 3.1. Show that the system

$$G(s) = \frac{\exp(s)}{s + 1} \quad \text{with } \text{Re}\{s\} > -1 \quad (3.47)$$

is not causal although the ROC is right to the rightmost pole.

Exercise 3.2. Specify all possible ROCs for the system

$$G(s) = \frac{s-1}{(s+1)(s-2)} \quad (3.48)$$

and determine the corresponding impulse response functions.

Note that if one assumes a system to be rational and causal, then this is equivalent to the usual stability criterion that all poles have negative real part, i.e., they lie in the left-half of the s -plane. Unlike feedback control, ILC methods do not need to be causal in time. In fact, strictly causal ILC algorithms are known to be inferior to noncausal algorithms since they lack the ability to *plan ahead*. This ability is closely linked to so-called *adjoint systems* that regularly appear in (optimization-based) feedforward control methods.

The \mathcal{L}_∞ norm

That a linear time-invariant system indeed maps some $\mathcal{L}_2(\mathbb{R}; \mathbb{R}^l)$ to $\mathcal{L}_2(\mathbb{R}; \mathbb{R}^m)$ is only true if $\mathbf{Gu} \in \mathcal{L}_2(\mathbb{R}; \mathbb{R}^m)$ for any $\mathbf{u} \in \mathcal{L}_2(\mathbb{R}; \mathbb{R}^l)$. Due to Parseval's theorem 3.4, one can evaluate the norm of a signal either in the time or frequency domain. The 2-norm of the output is thus

$$\begin{aligned} \|\mathbf{Gu}\|_2^2 &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathbf{G}(j\omega)\mathbf{u}(j\omega)\| d\omega \\ &\leq \frac{1}{2\pi} \int_{-\infty}^{\infty} \bar{\sigma}(\mathbf{G}(j\omega))^2 \|\mathbf{u}(j\omega)\| d\omega \\ &\leq \sup_{\omega} \bar{\sigma}(\mathbf{G}(j\omega))^2 \frac{1}{2\pi} \int_{-\infty}^{\infty} \|\mathbf{u}(j\omega)\| d\omega. \end{aligned} \quad (3.49)$$

Using the definition of the supremum norm or \mathcal{L}_∞ -norm

$$\|\mathbf{G}\|_\infty = \sup_{\omega} \bar{\sigma}(\mathbf{G}(j\omega)) \quad (3.50)$$

yields

$$\|\mathbf{Gu}\|_2 \leq \|\mathbf{G}\|_\infty \|\mathbf{u}\|_2, \quad (3.51)$$

i.e., the \mathcal{L}_∞ -norm is the induced operator norm of a mapping between two \mathcal{L}_2 -spaces. Consequently, a sufficient condition for $\mathbf{Gu} \in \mathcal{L}_2(\mathbb{R}; \mathbb{R}^m)$ is $\sup_{\omega} \bar{\sigma}(\mathbf{G}(j\omega)) < \infty$. Note that if $\mathbf{G}(s)$ is a rational transfer matrix, this is true if and only if $\mathbf{G}(s)$ has no poles on the imaginary axis.

Adjoint systems

For a given system $\mathbf{G} : \mathcal{L}_2(\mathbb{R}; \mathbb{R}^l) \rightarrow \mathcal{L}_2(\mathbb{R}; \mathbb{R}^m)$, the *adjoint system* is defined as the linear system $\mathbf{G}^\dagger : \mathcal{L}_2(\mathbb{R}; \mathbb{R}^m) \rightarrow \mathcal{L}_2(\mathbb{R}; \mathbb{R}^l)$ such that

$$\langle \mathbf{Gw}, \mathbf{y} \rangle = \langle \mathbf{w}, \mathbf{G}^\dagger \mathbf{y} \rangle. \quad (3.52)$$

It is easy to show that this uniquely defines \mathbf{G}^\dagger and that $(\mathbf{G}^\dagger)^\dagger = \mathbf{G}$. Furthermore, if a real-valued \mathbf{G} has a state-space representation $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D})$, then \mathbf{G}^\dagger has a realization $(-\mathbf{A}^T, -\mathbf{C}^T, \mathbf{B}^T, \mathbf{D}^T)$ and the transfer matrix

$$\mathbf{G}^\dagger(s) = \mathbf{G}^T(-s). \quad (3.53)$$

Exercise 3.3. Determine the adjoint \mathbf{G}^\dagger in the Hilbert space $\mathcal{L}_2([0, t_f])$ using (3.52).

Systems in state-space representation

For a given state-space representation of a LTI system with measurement noise

$$\dot{\mathbf{x}}(t) = \mathbf{Ax}(t) + \mathbf{Bu}(t), \quad \mathbf{x}(0) = \mathbf{x}_0 \quad (3.54a)$$

$$\mathbf{y}(t) = \mathbf{Cx}(t) + \mathbf{Du}(t) + \mathbf{w}(t), \quad (3.54b)$$

which is the main type of system we will investigate in this chapter, one obtains a corresponding input-output representation (3.1) as

$$\mathbf{y}(t) = \mathbf{C}\Phi(t)\mathbf{x}_0 + \mathbf{Du}(t) + \int_0^t \mathbf{C}\Phi(t-\tau)\mathbf{Bu}(\tau) d\tau + \mathbf{w}(t) \quad (3.55)$$

with $\Phi(t) = \exp(\mathbf{At})$. If $\mathbf{x}_0 = 0$, the system is said to be relaxed at $t = 0$. It is often convenient to assume that the system is relaxed in the infinitely remote past, wherefore $\mathbf{C}\Phi(t)\mathbf{x}_0$ vanishes and (3.55) can be extended to an infinite time horizon

$$\mathbf{y} = \mathbf{Gu} + \mathbf{w} \quad (3.56)$$

and $\mathbf{G}(t) = \mathcal{B}^{-1}\{\mathbf{G}(s)\} = \mathbf{C}\Phi(t-\tau)\mathbf{B} + \mathbf{D}\delta(t)$ using the Dirac-delta function $\delta(t)$. Note that (3.56) is more general than (3.54), since it can describe dynamic systems that do not have a state-space representation or whose state-space is infinite-dimensional, i.e., distributed-parameter systems described by partial differential equations. Further, the noise term \mathbf{v} can be arbitrary and may thus be used to also include process noise which is not present in (3.54).

3.3 Frequency-domain ILC methods on infinite time horizons

Following the previous section, we consider a linear LTI system that is relaxed in the infinitely remote past and performs the same task over and over again. Each iteration is then described by

$$\mathbf{y}_j = \mathbf{Gu}_j + \mathbf{w}_j \quad (3.57)$$

with the input and output quantities $\mathbf{y}_j(t) \in \mathbb{R}^l$ and $\mathbf{u}_j(t) \in \mathbb{R}^m$, respectively, as well as a exogenous disturbance $\mathbf{w}_j(t) \in \mathbb{R}^l$. We assume that the system \mathbf{G} is inherently BIBO-stable and thus $\mathbf{G}(j\omega) < \infty$ for all ω . Since ILC can be seen as an adaptive

open-loop control strategy, this assumption is quite natural. Unstable plants thus have to be stabilized before applying ILC methods. By using a linear ILC law (3.4) on an infinite time horizon, i.e.,

$$\mathbf{u}_{j+1} = \mathbf{Q}(\mathbf{u}_j + \mathbf{L}\mathbf{e}_j) \quad (3.58)$$

with

$$\mathbf{Qu} = \int_{-\infty}^{\infty} \mathbf{Q}(t - \tau) \mathbf{u}(\tau) d\tau \quad \text{and} \quad \mathbf{Lu} = \int_{-\infty}^{\infty} \mathbf{L}(t - \tau) \mathbf{u}(\tau) d\tau, \quad (3.59)$$

we want to iteratively track a desired output \mathbf{y}_d that exists and is uniquely defined by $\mathbf{y}_d = \mathbf{Gu}_d$. Following a signal processing nomenclature, \mathbf{Q} and \mathbf{L} are usually referred to as *Q-filter* and *L-filter* or *learning filter*, respectively.

3.3.1 Analysis of ILC laws

To analyze stability and convergence of the learning law (3.58), we assume a deterministic input-output behavior with $\mathbf{w}_j = 0$. Using the output error $\mathbf{e}_j = \mathbf{y}_d - \mathbf{y}_j$ and the learning law (3.58) yields

$$\mathbf{u}_{j+1} = \Psi \mathbf{u}_j + \Lambda \mathbf{y}_d, \quad (3.60)$$

with $\Psi = \mathbf{Q}(\mathbf{I} - \mathbf{LG})$ and $\Lambda = \mathbf{QL}$. The asymptotic input $\mathbf{u}_{j+1} = \mathbf{u}_j = \mathbf{u}_\infty$ of the input iteration (3.60) is given by

$$\mathbf{u}_\infty = (\mathbf{I} - \Psi)^{-1} \Lambda \mathbf{y}_d. \quad (3.61)$$

By introducing the input error $\bar{\mathbf{u}}_j = \mathbf{u}_j - \mathbf{u}_\infty$, one obtains the input error iteration

$$\bar{\mathbf{u}}_{j+1} = \Psi \bar{\mathbf{u}}_j, \quad (3.62)$$

for which the follow theorem assures stability:

Theorem 3.7 (Asymptotic stability of the ILC law). *The input error iteration (3.62) of the ILC law (3.58) is asymptotically stable if*

$$\sup_{\omega} \rho(\mathbf{Q}(j\omega)(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega))) < 1 \quad (3.63)$$

and \mathbf{u}_j converges to \mathbf{u}_∞ .

For practical reasons, we are usually much more interested in making stability assertions for the output error \mathbf{e}_j . Assuming that there exists a formal inverse \mathbf{G}^{-1} , we can rewrite the input-output behavior as $\mathbf{u}_j = \mathbf{G}^{-1}(\mathbf{y}_d - \mathbf{e}_j)$ which yields the output iteration

$$\mathbf{e}_{j+1} = \mathbf{G}\Psi\mathbf{G}^{-1}\mathbf{e}_j + (\mathbf{I} - \mathbf{G}\mathbf{Q}\mathbf{G}^{-1})\mathbf{y}_d. \quad (3.64)$$

It is immediately clear from this equation that perfect tracking, i.e., a vanishing asymptotic output error $\mathbf{e}_\infty = \mathbf{0}$, is only possible for $\mathbf{Q} = \mathbf{I}$ and one may be inclined to question other

choices for \mathbf{Q} . We will get back to this issue later. The asymptotic output error \mathbf{e}_∞ is given by

$$\begin{aligned}\mathbf{e}_\infty &= \mathbf{y}_d - \mathbf{G}\mathbf{u}_\infty \\ &= (\mathbf{I} - \mathbf{G}(\mathbf{I} - \mathbf{\Psi})^{-1}\mathbf{\Lambda})\mathbf{y}_d.\end{aligned}\quad (3.65)$$

Using $\bar{\mathbf{e}}_j = \mathbf{e}_j - \mathbf{e}_\infty$ yields

$$\bar{\mathbf{e}}_{j+1} = \mathbf{G}\mathbf{\Psi}\mathbf{G}^{-1}\bar{\mathbf{e}}_j. \quad (3.66)$$

Since

$$\rho(\mathbf{G}\mathbf{\Psi}\mathbf{G}^{-1}) = \rho(\mathbf{\Psi}), \quad (3.67)$$

it follows that:

Theorem 3.8 (Asymptotic stability of the output iteration). *The output iteration (3.64) of the ILC law (3.58) is asymptotically stable iff the input iteration is stable, i.e.,*

$$\sup_\omega \rho(\mathbf{Q}(j\omega)(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega))) < 1 \quad (3.68)$$

and \mathbf{e}_j then converges to the asymptotic tracking error

$$\mathbf{e}_\infty = (\mathbf{I} - \mathbf{G}(\mathbf{I} - \mathbf{\Psi})^{-1}\mathbf{\Lambda})\mathbf{y}_d = (\mathbf{I} - \mathbf{G}\mathbf{\Psi}\mathbf{G}^{-1})^{-1}(\mathbf{I} - \mathbf{G}\mathbf{Q}\mathbf{G}^{-1})\mathbf{y}_d. \quad (3.69)$$

Exercise 3.4. Show the equivalence of both expressions in (3.69).

If one wants to avoid (potentially) large transient errors during the learning process, monotonic convergence of the learning law has to be ensured.

Theorem 3.9 (Monotonic convergence of the input iteration). *The input iteration (3.60) of the ILC law (3.58) converges monotonically to \mathbf{u}_∞ , i.e., it holds that*

$$\|\mathbf{u}_{j+1} - \mathbf{u}_\infty\|_2 \leq \alpha \|\mathbf{u}_j - \mathbf{u}_\infty\|_2 \quad (3.70)$$

for $0 \leq \alpha < 1$ if

$$\|\mathbf{\Psi}\|_\infty = \sup_\omega \bar{\sigma}(\mathbf{Q}(j\omega)(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega))) = \alpha < 1. \quad (3.71)$$

Theorem 3.10 (Monotonic convergence of the output iteration). *The output iteration (3.64) of the ILC law (3.58) converges monotonically to \mathbf{e}_∞ , i.e., it holds that*

$$\|\mathbf{e}_{j+1} - \mathbf{e}_\infty\|_2 \leq \beta \|\mathbf{e}_j - \mathbf{e}_\infty\|_2 \quad (3.72)$$

for $0 \leq \alpha < 1$ if

$$\|\mathbf{G}\Psi\mathbf{G}^{-1}\|_\infty = \sup_{\omega} \bar{\sigma}(\mathbf{G}(j\omega)\mathbf{Q}(j\omega)(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega))\mathbf{G}^{-1}(j\omega)) = \beta < 1 . \quad (3.73)$$

Note that monotonic convergence of the input iteration *does not* imply monotonic convergence of the output iteration and vice versa.

3.3.2 Deterministic design methods

All types of ILC algorithms try to perform some kind of approximate inversion of the input-output behavior of the system. Since this inversion has to be performed online based on measurements, it is inherently constrained by the presence of disturbances and measurement noise. Ideally, an ILC algorithm is able to iteratively separate repeating disturbances and effects of a model-plant mismatch from non-repeating disturbances and measurement noise. Nevertheless, most ILC designs are developed within a deterministic framework with $\mathbf{w} = \mathbf{0}$ which we will adopt in this section. Note that we introduce ILC as a pure open-loop control strategy here, henceforth ILC algorithms are not able to respond to non-repeating (and unanticipated) disturbances. For the performance of the complete control concept, it is thus advisable to incorporate feedback control methods. For simplicity, we thus assume that the system \mathbf{G} already incorporates a suitable feedback control strategy and thus \mathbf{G} is BIBO-stable.

We know from the previous section that $\mathbf{Q} = \mathbf{I}$ is a necessary condition to achieve perfect tracking. One thus only wants to deviate from this ideal if necessary to achieve a stable ILC algorithm that is sufficiently *robust* to model variations. We will thus start with three typical design strategies to obtain suitable learning operators before considering the relation between Q-filtering and robustness. As we will see later, good learning filters usually avoid to learn at high frequencies due to the presence of measurement noise. By assuming $\mathbf{Q} = \mathbf{I}$, however, $\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega) \rightarrow \mathbf{I}$ for $\omega \rightarrow \infty$, which would always violate the stability condition in Theorem 3.7. We will therefore consider a finite frequency range up to some $\bar{\omega}$ and assume $\mathbf{Q} = \mathbf{0}$ above, which is always possible.

PD-type ILC laws

As the name implies, PD-type learning laws combine proportional and derivative action, i.e,

$$\mathbf{Le}_j = \mathbf{K}_p \mathbf{e}_j(t) + \mathbf{K}_d \frac{d\mathbf{e}_j(t)}{dt} \quad (3.74)$$

with the proportional and derivative gain matrices \mathbf{K}_p and \mathbf{K}_d , respectively, that are tunable parameters that have to be chosen such that a desired performance is reached. These heuristic designs are arguably the most widely used ILC laws in the literature since they do not require an accurate model but rely on intensive tuning of the gain matrices. Since this can be a quite tedious task for MIMO systems, PD-type ILC laws are typically used for SISO systems only, which yields the learning filter

$$\mathbf{L}(s) = k_p + k_d s \quad (3.75)$$

in the frequency domain. While there are no tuning guidelines, it is usually suggested to start with small gain values and increase them until a growth of (usually) high-frequency components indicate that the learning law left the stable parameter regime [3.0].

The fact that such simple ILC laws perform quite well for a large class of systems is somewhat surprising. However, this can be explained by observing that a PD-type laws (3.74) is the exact inverse of a first-order lag (i.e., P-T₁) element. Since most technical systems exhibit some kind of low-pass behavior, PD-type laws can be seen as an approximate inverse of their behavior. Note that the learning filter (3.75) yields infinite gains for $\omega \rightarrow \infty$, which is problematic in the presence of measurement noise.

Example 3.1. Consider a system described by the scalar transfer function

$$G(s) = \frac{s + 1/2}{s^2 + s + 3} \quad (3.76)$$

together with the PD-type law (3.75). To evaluate stability for different parameter choices of k_p and k_d , the absolute value of $|1 - L(j\omega)G(j\omega)|$ is plotted in Fig. 3.3 together with the resulting convergence behavior of the ILC law. Note that the ILC law ultimately reaches the precision of the numerical simulation. Observing that the inverse of the system's transfer function is

$$G^{-1}(s) = 1/2 + s + \frac{11}{2(2s + 1)}, \quad (3.77)$$

one could try to use $k_p = 1/2$ and $k_d = 1$, which is in fact unstable according to Theorem 3.7.

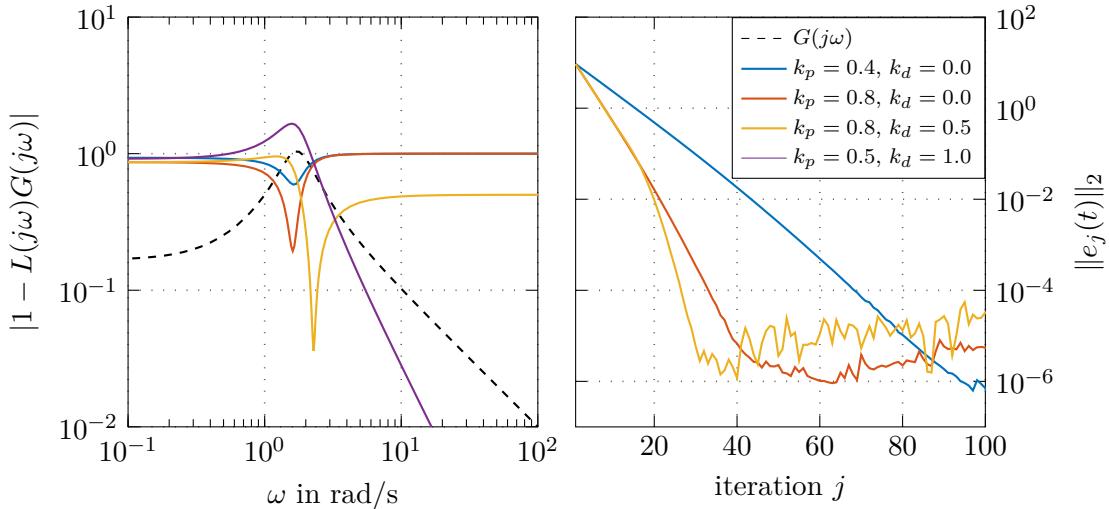


Figure 3.3: Stability and convergence behavior of a PD-type learning law for different parameters k_p and k_d .

Inversion-based ILC law

A particularly intuitive choice is to use $\mathbf{L} = \mathbf{G}^{-1}$, which yields

$$\mathbf{u}_{j+1} = \mathbf{Q}(\mathbf{u}_j + \mathbf{G}^{-1}\mathbf{e}_j) = \mathbf{Q}(\mathbf{u}_j + \mathbf{G}^{-1}(\mathbf{y}_d - \mathbf{G}\mathbf{u}_j - \mathbf{y}_0)) = \mathbf{Q}(\mathbf{G}^{-1}(\mathbf{y}_d - \mathbf{y}_0)). \quad (3.78)$$

The input \mathbf{u}_j thus remains unchanged after the first iteration (dead-beat behavior) with a resulting output error

$$\mathbf{e}_{j+1} = \mathbf{e}_\infty = \mathbf{y}_d - \mathbf{G}\mathbf{Q}(\mathbf{G}^{-1}(\mathbf{y}_d - \mathbf{y}_0)) - \mathbf{y}_0. \quad (3.79)$$

Using $\mathbf{Q} = \mathbf{I}$ further implies $\mathbf{e}_\infty = \mathbf{0}$. There are several caveats to this result: First of all, \mathbf{G} is never known exactly which limits the usability of inversion-based ILC laws and mandates an additional \mathbf{Q} -filter. Further, $\mathbf{G}(s)$ is strictly proper for any physical process, wherefore its inverse $\mathbf{G}^{-1}(s)$ is unbound for $s = j\omega \rightarrow \infty$ and thus \mathbf{e}_∞ is undefined for arbitrary $\mathbf{y}_d - \mathbf{y}_0 \in \mathcal{L}_2(\mathbb{R})$. Finally, \mathbf{G}^{-1} is unstable for non-minimum phase systems although this can be alleviated by (non-causal) stable-inversion methods [3.0].

Pseudo-Inversion-based ILC law

Using an inversion-based ILC law is essentially determining an exact feedforward input signal from output data. Unlike for feedforward purposes, however, we cannot demand certain levels of regularity of the output signals due to the presence of stochastic disturbances and measurement noise. The obvious solution is to regularize the system inversion by using a pseudo-inverse learning filter [3.0], i.e.,

$$\mathbf{L}(s) = (\alpha\mathbf{I} + \mathbf{G}^\dagger(s)\mathbf{G}(s))^{-1}\mathbf{G}^\dagger(s) \quad (3.80)$$

with the regularization parameter $\alpha > 0$ and the transfer matrix of the adjoint system $\mathbf{G}^\dagger(s)$. The regularization parameter separates regions where the learning law is approaching an inversion-based law (i.e., $\mathbf{L}(s) \approx \mathbf{G}^{-1}(s)$) from regions where learning is almost prohibited (i.e., $\mathbf{L}(s) \approx \mathbf{0}$) depending on whether $\|\mathbf{G}^\dagger(s)\mathbf{G}(s)\|$ is much larger or smaller than α , respectively. For systems with low-pass behavior where $\mathbf{G}(s)$ is strictly proper, the learning filter (3.80) thus naturally avoids learning in the high-frequency region, i.e., $\lim_{s \rightarrow \infty} \mathbf{L}(s) = \mathbf{0}$. It still remains to be shown that (3.80) results in a stable learning iteration. Using

$$\mathbf{I} - \mathbf{L}\mathbf{G} = \mathbf{I} - (\alpha\mathbf{I} + \mathbf{G}^\dagger\mathbf{G})^{-1}\mathbf{G}^\dagger\mathbf{G} = \left(\mathbf{I} + \frac{1}{\alpha}\mathbf{G}^\dagger\mathbf{G}\right)^{-1} \quad (3.81)$$

it follows with the definition of the adjoint system that

$$\rho(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega)) = \rho\left(\left(\mathbf{I} + \frac{1}{\alpha}(\mathbf{G}(-j\omega))^T\mathbf{G}(j\omega)\right)^{-1}\right) \quad (3.82)$$

$$= \rho\left(\left(\mathbf{I} + \frac{1}{\alpha}(\mathbf{G}(j\omega))^H\mathbf{G}(j\omega)\right)^{-1}\right). \quad (3.83)$$

Assuming that $\mathbf{G}(j\omega)\mathbf{u} \neq \mathbf{0}$ for any $\mathbf{u} \neq \mathbf{0}$, the matrix $\frac{1}{\alpha}(\mathbf{G}(j\omega))^H \mathbf{G}(j\omega)$ is positive definite and thus all of its eigenvalues are strictly positive. This implies the asymptotic stability of the learning law since $\rho(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega)) < 1$ for all ω .

Example 3.2. Consider again the plant (3.76). According to Figure 3.3, the magnitude of $G(j\omega)$ is above $1 \cdot 10^{-1}$ for the essential part of its dynamic behavior. The resulting Pseudo-Inversion-based learning filter (3.80) for different choices of the regularization parameter $\alpha = \{1 \cdot 10^{-1}, 5 \cdot 10^{-2}, 1 \cdot 10^{-2}, 1 \cdot 10^{-3}\}$ is illustrated in Figure 3.4. As one can see, the learning filter (3.80) is an approximation of the exact inverse $G^{-1}(j\omega)$. While $L(j\omega)$ deviates (significantly) over the whole frequency range for high values of α , small values of α only introduce a roll-off for high-frequency components. Since asymptotic convergence is equivalent to monotonic convergence for SISO plants according to Theorem 3.7 and Theorem 3.9, the learning law (3.80) always converges exponentially. Note that this holds only true for infinite time horizons and that the convergence rate in Theorem (3.9) is in general not equivalent (but determined by) the regularization parameter α .

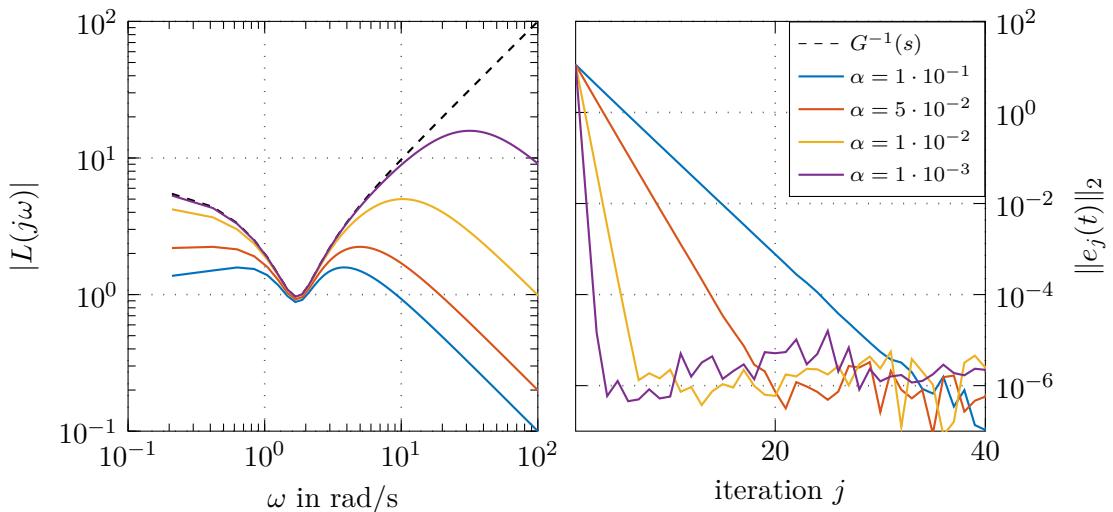


Figure 3.4: Learning filter and convergence behavior of the Pseudo-Inversion-based learning filter (3.80) for different regularization parameters α .

3.3.3 Stochastically optimal learning laws

We already mentioned in the previous section that the approximate inversion due to the learning filter \mathbf{L} is inherently constrained by the presence of stochastic quantities, namely (non-repeating) process disturbances and measurement noise $\mathbf{w}_j(t) \in \mathbb{R}^l$ in

$$\mathbf{y}_j(t) = \mathbf{G}\mathbf{u}_j(t) + \mathbf{w}_j(t) \quad (3.84)$$

that is given by a zero-mean wide-sense stationary (WSS) process with identical stochastic properties in each iteration. To analyze this problem in a stochastic framework, we

introduce the following notation: For two vector-valued stochastic signals $\mathbf{a}_k(t)$ and $\mathbf{b}_l(t)$ with the iteration indices k and l , their cross-correlation function is given by $\mathbf{R}^{a_k b_l}(\tau) = \text{E}\{\mathbf{a}_k(t + \tau)(\mathbf{b}_l(t))^T\}$ where $\text{E}\{\cdot\}$ denotes the expectation operator. The corresponding power spectral density (PSD) reads as $\mathbf{S}^{a_k b_l}(s) = \mathcal{B}\{\mathbf{R}^{a_k b_l}(\tau)\}$. In case $k = l$, the common index is written as subscript, i.e., $\mathbf{R}_k^{ab}(\tau) = \mathbf{R}^{a_k b_k}(\tau)$ and $\mathbf{S}_k^{ab}(s) = \mathbf{S}^{a_k b_k}(s)$.

For generality, we use a *iteration-varying* learning law (again without \mathbf{Q} -filter)

$$\mathbf{u}_{j+1}(t) = \mathbf{u}_j(t) + \mathbf{L}_j \mathbf{e}_j(t). \quad (3.85)$$

Our goal now is to design a learning filter \mathbf{L}_j in a stochastically optimal way, i.e. $\mathbf{u}_{j+1}(t)$ should be determined such that it minimizes the expected value of the mean-square output error $\text{E}\{(\mathbf{e}_{j+1}(t))^T \mathbf{e}_{j+1}(t)\}$. Following section 3.3.1, the output error $\mathbf{e}_j(t) = \mathbf{y}_d(t) - \mathbf{y}_j(t)$ of the iteration j is given by

$$\mathbf{e}_j(t) = \mathbf{G} \boldsymbol{\nu}_j(t) - \mathbf{w}_j(t) \quad (3.86)$$

using the input error $\boldsymbol{\nu}_j(t) = \mathbf{u}_d(t) - \mathbf{u}_j(t)$. Together with the learning law, the evolution of the input error is described by

$$\boldsymbol{\nu}_{j+1}(t) = \boldsymbol{\nu}_j(t) - \mathbf{L}_j \mathbf{e}_j(t) \quad (3.87)$$

and the output error of the iteration $j + 1$ yields

$$\mathbf{e}_{j+1}(t) = (\mathbf{I} - \mathbf{G} \mathbf{L}_j) \mathbf{e}_j(t) + \mathbf{w}_j(t) - \mathbf{w}_{j+1}(t). \quad (3.88)$$

The problem of learning in a stochastically optimal sense can be written as the optimization problem

$$\min_{\mathbf{L}_j(t)} \text{E}\{(\mathbf{e}_{j+1}(t))^T \mathbf{e}_{j+1}(t)\}. \quad (3.89)$$

Since we only consider LTI systems and the stochastic disturbance \mathbf{w}_j is WSS, the optimization problem (3.89) can be treated in the Laplace domain by applying the bilateral Laplace transform, which yields

$$\begin{aligned} \min_{\mathbf{L}_j(t)} \text{E}\{(\mathbf{e}_{j+1}(t))^T \mathbf{e}_{j+1}(t)\} &= \min_{\mathbf{L}_j(t)} \text{Tr}\{\mathbf{R}_{j+1}^{ee}(0)\} \\ &= \min_{\mathbf{L}_j(s)} \frac{1}{j2\pi} \text{Tr}\left\{\int_{-\infty}^{\infty} \mathbf{S}_{j+1}^{ee}(j\omega) d\omega\right\}, \end{aligned} \quad (3.90)$$

where $\text{Tr}\{\cdot\}$ denotes the trace operator. To obtain the output error PSD $\mathbf{S}_{j+1}^{ee}(s)$, one can use (3.88) together with (3.86), which yields

$$\mathbf{e}_{j+1}(t) = (\mathbf{G} - \mathbf{G} \mathbf{L}_j \mathbf{G}) \boldsymbol{\nu}_j(t) + \mathbf{G} \mathbf{L}_j \mathbf{w}_j(t) - \mathbf{w}_{j+1}(t). \quad (3.91)$$

In general, the stochastic quantities $\boldsymbol{\nu}_j$, \mathbf{w}_j and \mathbf{w}_{j+1} will be correlated. We thus make the following assumptions:

- A1 different instances of the disturbance are uncorrelated, i.e., $E\{\mathbf{w}_i(t+\tau)(\mathbf{w}_j(t))^T\} = \mathbf{0}$ for $i \neq j$
- A2 the input error $\boldsymbol{\nu}_j$ is uncorrelated with the exogenous disturbance of the current and the following iteration, i.e., $E\{\boldsymbol{\nu}_j(t+\tau)(\mathbf{w}_i(t))^T\} = \mathbf{0}$ for $i \in \{j, j+1\}$. Together with A1, this is equivalent to $E\{\boldsymbol{\nu}_0(t+\tau)(\mathbf{w}_j(t))^T\} = \mathbf{0}$
- A3 the stochastic properties of the disturbance do not change over iterations, i.e., $E\{\mathbf{w}_j(t+\tau)(\mathbf{w}_j(t))^T\} = \mathbf{S}^{ww}$.

Using these assumption, we obtain

$$\mathbf{S}_{j+1}^{ee} = (\mathbf{G} - \mathbf{GL}_j \mathbf{G}) \mathbf{S}_j^{\nu\nu} (\mathbf{G} - \mathbf{GL}_j \mathbf{G})^\dagger + \mathbf{GL}_j \mathbf{S}^{ww} (\mathbf{L}_j)^\dagger \mathbf{G}^\dagger + \mathbf{S}^{ww} \quad (3.92)$$

The optimization problem (3.90) can be solved by variational calculus. The resulting learning filter

$$\mathbf{L}_k(s) = \mathbf{S}_j^{\nu\nu}(s) \mathbf{G}^\dagger(s) [\mathbf{G}(s) \mathbf{S}_j^{\nu\nu}(s) \mathbf{G}^\dagger(s) + \mathbf{S}^{ww}(s)]^{-1} \quad (3.93)$$

is an iterative version of the well-known (non-causal) Wiener filter, which is used to *estimate* the input deviation that optimally *explains* the measured output error. A common problem of Wiener-filter-based approaches is that the optimal solution (3.93) requires knowledge of the input error PSD $\mathbf{S}_j^{\nu\nu}(s)$, which is typically handled using a-priori knowledge of the problem.

Due to the learning process, the measured output error $\mathbf{e}_j(t)$ will be increasingly dominated by stochastic disturbances. A stochastically optimal learning law will therefore reduce its learning action with increasing iterations. Such a behavior is intrinsic to (3.93) due to the decreasing input error PSD $\mathbf{S}_j^{\nu\nu}(s)$. From (3.86), (3.87) and (3.93) we obtain the iterative relation

$$\mathbf{S}_{j+1}^{\nu\nu}(s) = (\mathbf{I} - \mathbf{L}_j(s) \mathbf{G}(s)) \mathbf{S}_j^{\nu\nu}(s). \quad (3.94)$$

By starting from an initial PSD $\mathbf{S}_0^{\nu\nu}(s)$ and the corresponding learning filter (3.93), one can iterate forward in time to obtain a stochastically optimal ILC method.

Theorem 3.11 (Stochastically optimal learning). *If $\mathbf{S}^{ww}(s)$ and the initial input error PSD $\mathbf{S}_0^{\nu\nu}(s)$ are positive definite and the system's transfer matrix $\mathbf{G}(s)$ does not exhibit transmission zeros on the imaginary axis, the learning filter (3.93) together with (3.94) yields a stable learning law that ensures convergence to the optimal error PSDs*

$$\mathbf{S}_\infty^{\nu\nu}(s) = \mathbf{0}, \quad \mathbf{S}_\infty^{\eta\eta}(s) = \mathbf{S}^{ww}(s). \quad (3.95)$$

The fact that one can iteratively construct a noise-less representation of the unknown desired input \mathbf{u}_∞ from noisy output measurements and an output error containing only noise sounds like a very attractive solution. However, there is a severe limitation to this seemingly nice result: By design, the learning filter is asymptotically vanishing and thus the learning process comes to a halt. Since the forward iteration (3.94) is independent

of actual measurements, this is the case even if some unforeseen disturbance (i.e., not represented in the stochastic model) is causing large output errors.

Since iteration-varying learning laws furthermore increase the effort of implementation, one may thus prefer to accept sub-optimal performance by using a fixed input error PSD $\mathbf{S}_j^{\nu\nu}(s) = \mathbf{S}^{\nu\nu}(s)$.

Theorem 3.12 (Stochastically sub-optimal learning). *If $\mathbf{S}^{ww}(s)$ and the chosen input error PSD $\mathbf{S}^{\nu\nu}(s)$ are positive definite and the system's transfer matrix $\mathbf{G}(s)$ does not exhibit transmission zeros on the imaginary axis, the iteration-invariant learning filter*

$$\mathbf{L}(s) = \mathbf{S}^{\nu\nu}(s)\mathbf{G}^\dagger(s)\left[\mathbf{G}(s)\mathbf{S}^{\nu\nu}(s)\mathbf{G}^\dagger(s) + \mathbf{S}^{ww}(s)\right]^{-1} \quad (3.96)$$

yields a stable learning law that converges to a positive definite asymptotic output error PSD $\mathbf{S}_\infty^{ee}(s)$ given by the solution of

$$(\mathbf{I} - \mathbf{G}(s)\mathbf{L}(s))\mathbf{S}_\infty^{ee}(s)(\mathbf{I} - \mathbf{G}(s)\mathbf{L}(s))^\dagger - \mathbf{S}_\infty^{ee}(s) + \mathbf{G}\mathbf{L}(s)\mathbf{S}^{ww}(s) + \mathbf{S}^{ww}(s)\mathbf{L}^\dagger(s)\mathbf{G}^\dagger(s) = \mathbf{0}. \quad (3.97)$$

This sub-optimal learning law is structurally similar to pseudo-inversion-based learning laws. For the special case of $\mathbf{S}^{\nu\nu}(s) = \sigma_\nu \mathbf{I}$ and $\mathbf{S}^{ww}(s) = \sigma_w \mathbf{I}$, it is easy to show that (3.96) is identical to (3.80) with $\alpha = \frac{\sigma_w}{\sigma_\nu}$. Stochastically optimal learning filters thus regularize the system inversion according to the expected signal-to-noise ratio.

3.3.4 Q-filtering and robustness

One of the main advantages of ILC over other control approaches is its ability to achieve (almost) perfect tracking in the presence of external disturbances and model-plant mismatch. While we have investigated external disturbances in the previous section, we may now shift our attention to an inevitable model-plant mismatch. Assuming that the plant can be described by $\mathbf{G}(s) = \mathbf{G}_0(s)\Delta\mathbf{G}(s)$ with the nominal (design) model $\mathbf{G}_0(s)$ and the unknown deviation $\Delta\mathbf{G}(s)$, it is expected that $\rho(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}_0(s)\Delta\mathbf{G}(s))$ will in general be larger than one for some ω since we designed the learning filter $\mathbf{L}(s)$ such that $\rho(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}_0(s)) < 1$. However, it is clear that the asymptotic stability criterion

$$\sup_\omega \rho(\mathbf{Q}(j\omega)(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}(j\omega))) < 1 \quad (3.98)$$

can always be met by choosing $\mathbf{Q}(s)$ sufficiently small. For example, for a known $\Delta\mathbf{G}(s)$ one could always use $\mathbf{Q}(s) = \kappa\mathbf{I}$ with $\kappa > 1/\sup_\omega \rho(\mathbf{I} - \mathbf{L}(j\omega)\mathbf{G}_0(j\omega)\Delta\mathbf{G}(j\omega))$.

Using a spectrally uniform \mathbf{Q} -filter is usually not recommended, since $\Delta\mathbf{G}(s)$ is typically small in those parts of the frequency range within which we want to learn. To avoid unnecessary large asymptotic tracking errors, a spectrally selective \mathbf{Q} -filter is thus usually advisable. The most common case in practice is that the (identified) nominal model fits quite accurately up to a certain frequency ω_c and starts to deteriorate beyond that frequency, which motivates the use of low-pass filters.

Example 3.3. Consider the plant (3.76) again with an additional model-plant mismatch given by

$$G_0(s) = \frac{s + 1/2}{s^2 + s + 3} \quad \text{and} \quad \Delta G(s) = \frac{1}{1 + s/30 + (s/15)^2} \quad (3.99)$$

and a pseudo-inverse learning law $L(s)$ with $\alpha = 1 \cdot 10^{-2}$ for the nominal model $G_0(s)$. Simulation scenarios for different choices of $Q(s)$ are illustrated in Figure 3.5. For $Q(s) = 1$, we can see that $|1 - L(j\omega)G(j\omega)| > 1$ for $\omega > 10$, which leads to high-frequency oscillations that build up over time. Note that the error seems to converge initially but diverges after 20 iterations. Using a simple first-order low-pass filter

$$Q(s) = \frac{10}{s + 10} \quad (3.100)$$

stabilizes the learning iteration, but at the cost of a significantly higher asymptotic output error. Interestingly, for the more aggressive choice (see left-hand side of Figure 3.5)

$$Q(s) = \frac{10^2}{(s + 10)(-s + 10)}, \quad (3.101)$$

the asymptotic output error is vastly improved. The main reason for this is that simple first-order low-pass filter introduces a phase shift to the learned signal, which in turn results in a slight temporal mismatch between $y_d(t)$ and $y_\infty(t)$ that dominates all other error sources. The latter choice of $Q(s)$ is a so-called *zero-phase* filter.

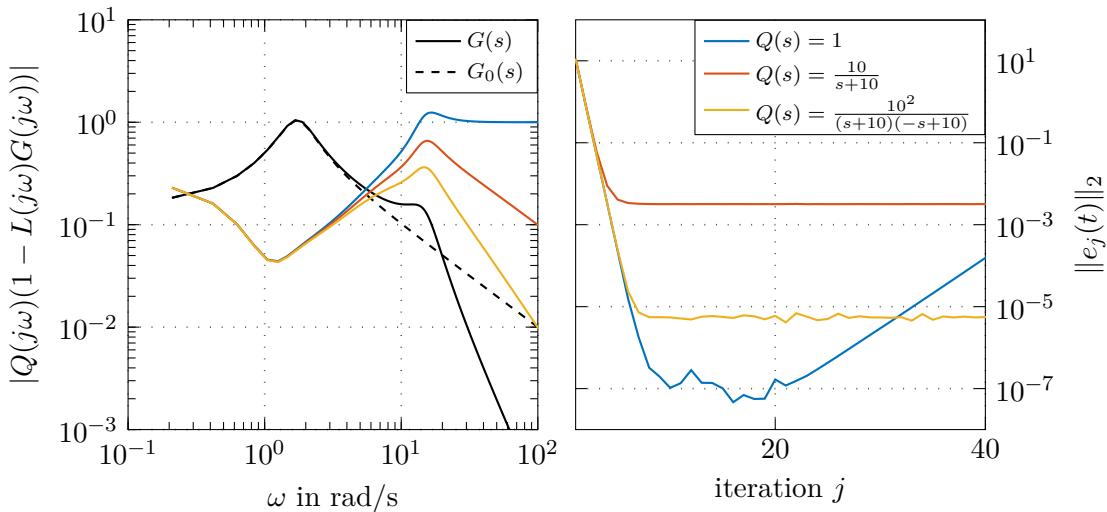


Figure 3.5: Robustness to model-plant mismatch by using a **Q**-filter.

Since there is no systematic design for **Q** in the MIMO case, we will restrict ourselves to the SISO case in the following. In line with Example 3.3, it is generally recommended

to use *zero-phase* filters for $Q(s)$, i.e.,

$$\text{Im}\{Q(j\omega)\} = 0. \quad (3.102)$$

Exercise 3.5. Show that a zero-phase filter $\mathbf{Q}(s)$ that is not a constant gain matrix is necessarily noncausal.

Gaussian filter

A particularly simple type of zero-phase filter is the Gaussian filter given by the impulse response

$$Q(t) = \frac{1}{\sigma_q \sqrt{2\pi}} \exp\left(-\frac{t^2}{2\sigma_q^2}\right) \quad (3.103)$$

with the standard deviation σ_q that is related to the 3-dB bandwidth f_c

$$\sigma_q = \frac{\sqrt{\ln(2)}}{2\pi f_c} \quad (3.104)$$

of the corresponding frequency-domain representation

$$Q(j\omega) = \exp\left(-\frac{\sigma_q^2 \omega^2}{2}\right). \quad (3.105)$$

Forward-Backward-Filtering

A popular alternative with arbitrary spectral behavior is to use

$$Q(s) = Q_f^\dagger(s)Q_f(s) \quad (3.106)$$

where $Q_f(s)$ is stable and causal. As a result, $Q(s)$ can be easily implemented by integrating the realization of $Q_f(s)$ forward and $Q_f^\dagger(s)$ backward in time (see MATLAB command `filtfilt`). This procedure was used in Example 3.3 to implement the zero-phase Q -filter.

Time delay

One particular type of model-plant mismatch in ILC applications is a temporal delay by time T , i.e.,

$$\Delta G(s) = \exp(-sT). \quad (3.107)$$

Using a pseudo-inversion-based learning law, it follows that

$$\begin{aligned} |1 - L(j\omega)G(j\omega)| &= |1 - L(j\omega)G_0(j\omega) \exp(-j\omega T)| \\ &= \left| 1 - \frac{|G_0(j\omega)|^2}{\alpha + |G_0(j\omega)|^2} \exp(-j\omega T) \right|. \end{aligned} \quad (3.108)$$

As one can see, the resulting ILC law is always unstable in the absence of a Q -filter for sufficiently high frequencies ω . Any predictable time-delay should thus be compensated by shifting the time axis respectively.

3.3.5 Implementation aspects

After having designed a suitable learning filter $\mathbf{L}(s)$, there are still some open questions on how to implement the resulting learning law for a practical problem that is typically defined on a finite time horizon $t \in [0, t_f]$.

Spectral factorization

A direct method to achieve this is to use spectral factorization methods to decompose the learning filter into the product of two parts: one that is causal (i.e., its impulse response matrix is in \mathcal{H}_2), and one that is anti-causal (i.e., its impulse response matrix is in \mathcal{H}_2^\perp). Since the anti-causal part can also be written as the adjoint of a causal system, one obtains the factorization

$$\mathbf{L}(s) = \mathbf{L}_b^\dagger(s)\mathbf{L}_f(s) \quad (3.109)$$

As a result, the learning filter $\mathbf{L}(s)$ can be implemented by finding a state-space representation of $\mathbf{L}_f(s)$ that is solved forward in time and a state-space representation of $\mathbf{L}_b^\dagger(s)$ that is solved backwards in time. Note that this solution procedure requires corresponding initial and terminal conditions for the state-space representations that are not determined from the infinite time-horizon design.

FIR-filter approximation

Finding a factorization (3.109) can be quite tedious except for SISO systems. An alternative solution is to determine the impulse response matrix $\mathbf{L}(t)$ directly. Since $\mathbf{L}(t)$ is in general not of finite support, it has to be truncated in time to a finite-impulse response (FIR) approximation to implement the convolution (3.59). Along this line, $\mathbf{L}(t)$ can either be transformed back analytically using the inverse bilateral Laplace transform and truncated afterwards or one chooses to use numerical methods. By assumption, the system's transfer matrix $\mathbf{G}(s)$ has no poles on the imaginary axis $s = j\omega$, wherefore $\mathbf{L}(s)$ has a ROC around the imaginary axis where $\mathbf{L}(t)$ is given by

$$\mathbf{L}(t) = \mathcal{F}^{-1}\{\mathbf{L}(j\omega)\}. \quad (3.110)$$

A convolution with $\mathbf{L}(t)$ is therefore equivalent to a forward and backward integration of the spectral factorization (3.109). The continuous Fourier transform can be approximately computed using the discrete Fourier transform (DFT) on a sampled temporal and spectral grid. By assuming that $\mathbf{L}(j\omega) = \mathbf{0}$ outside the finite interval $[-\frac{\omega_s}{2}, \frac{\omega_s}{2}]$, we can approximate the inverse Fourier transform for $t \in [-\frac{t_s}{2}, \frac{t_s}{2}]$ on the discrete grids $t_n = (n - N/2)\frac{t_s}{N}$ and

$\omega_m = (m - N/2) \frac{\omega_s}{N}$ with $0 \leq n, m < N$ by

$$\begin{aligned}\mathbf{L}(t_n) &= \frac{1}{2\pi} \int_{-\omega_s/2}^{\omega_s/2} \mathbf{L}(j\omega) \exp[j\omega t_n] d\omega \\ &\approx \frac{\omega_s}{2\pi N} \sum_{m=0}^{N-1} \mathbf{L}(j\omega_m) \exp\left[j \frac{t_s \omega_s}{N^2} \left(n - \frac{N}{2}\right) \left(m - \frac{N}{2}\right)\right] \\ &\approx \frac{\omega_s}{2\pi N} \exp\left[-j \frac{t_s \omega_s}{2N} \left(n - \frac{N}{2}\right)\right] \sum_{m=0}^{N-1} \mathbf{L}(j\omega_m) \exp\left[-j \frac{t_s \omega_s}{2N} m\right] \exp\left[j \frac{t_s \omega_s}{N^2} mn\right]\end{aligned}\quad (3.111)$$

Calculating the right-hand side of (3.111) can be computationally quite expansive. For the case $\omega_s t_s = 2\pi N$, however, this can be reformulated as¹

$$\begin{aligned}\mathbf{L}(t_n) &\approx \frac{\omega_s}{2\pi N} \exp\left[-j\pi \left(n - \frac{N}{2}\right)\right] \sum_{m=0}^{N-1} \mathbf{L}(j\omega_m) \exp[-j\pi m] \exp\left[j \frac{2\pi}{N} mn\right] \\ &= \frac{\omega_s}{2\pi N} \exp\left[-j\pi \left(n - \frac{N}{2}\right)\right] \text{DFT}\{\mathbf{L}(j\omega_m) \exp[-j\pi m]\}\end{aligned}\quad (3.112)$$

which is computationally very efficient by using FFT algorithms to compute the DFT for N chosen as a power of 2. Note that this method does not require $\mathbf{L}(s)$ to be a rational transfer matrix, see Example 3.4.

Boundary effects

Finally, to implement such learning laws requires information that is not provided by the infinite time horizon design. This lack was already quite explicit when factorizing the learning law and finding state-space representations for forward and backward integration, where suitable boundary conditions at $t = 0$ and $t = t_f$ are required. Using the convolutional representation (3.59) with a FIR filter of length t_s , i.e., $\mathbf{L}(t) \neq \mathbf{0}$ only if $t \in [-t_s/2, t_s/2]$, yields

$$(\mathbf{L}\mathbf{e}_j)(t) = \int_{-t_s}^{t_f+t_s} \mathbf{L}(t-\tau) \mathbf{e}_j(\tau) d\tau \quad \text{for } t \in [0, t_f]. \quad (3.113)$$

Rather than additional boundary values, the convolutional representation requires values of $\mathbf{e}_j(t)$ for $t \in [-t_s, t_f + t_s]$, which extends beyond the available measurements. Since we assumed that the system is correctly (and identically) initialized, $\mathbf{e}_j(0) = \mathbf{0}$ and thus $\mathbf{e}_j(t) = \mathbf{0}$ for $t \in [-t_s, 0]$ is a obvious choice. For $t = t_f$, the two most widely used options are:

1. $\mathbf{e}(t) = 0$ for $t \in [t_f, t_f + t_s]$ (truncation)
2. $\mathbf{e}(t) = \mathbf{e}(t_f)$ for $t \in [t_f, t_f + t_s]$ (extension)

¹This links the temporal resolution t_s/N with the spectral resolution ω_s/N , which can be problematic for some applications with large bandwidths where a sufficient temporal resolution requires a very high number of discretization points.

Note that this choice will impact the stability of the learning law in general [3.0]. However, an infinite time horizon convolution operator Ψ according to

$$\Psi \mathbf{e} = \int_{-\infty}^{\infty} \Psi(t - \tau) \mathbf{e}(\tau) d\tau, \quad (3.114)$$

with $\mathbf{e} \in \mathcal{L}_2(-\infty, \infty)$ is a contraction mapping, i.e., $\|\Psi \mathbf{e}\|_2 < \|\mathbf{e}\|_2$, then its truncated version

$$\Psi_T \mathbf{e}' = \int_0^{t_f} \Psi(t - \tau) \mathbf{e}'(\tau) d\tau \quad (3.115)$$

is a contraction mapping on the truncated space $\mathbf{e}' \in \mathcal{L}_2([0, t_f])$, too. This can be seen by using the standard truncation operator

$$(\mathbf{T}\mathbf{e})(t) = \begin{cases} \mathbf{e}(t) & \text{for } t \in [0, t_f], \\ \mathbf{0} & \text{otherwise,} \end{cases} \quad (3.116)$$

since

$$\|\Psi_T \mathbf{e}'\|_{2,[0,T]} = \|\mathbf{T}\Psi \mathbf{T}\mathbf{e}'\|_2 \leq \|\Psi \mathbf{T}\mathbf{e}'\|_2 \leq \|\Psi\|_\infty \|\mathbf{T}\{\mathbf{e}'\}\|_2 = \|\Psi\|_\infty \|\mathbf{e}'\|_{2,[0,T]} \quad (3.117)$$

using the identity $\Psi_T = \mathbf{T}\Psi\mathbf{T}$ and the fact that the induced norm $\|\mathcal{T}\|_\infty = 1$.

Example 3.4. Consider the coupled system of parabolic PDEs

$$\frac{\partial \mathbf{x}}{\partial t}(z, t) = \underbrace{\begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}}_{=\Lambda} \frac{\partial^2}{\partial z^2} \mathbf{x}(z, t) + \underbrace{\begin{bmatrix} 0 & \sigma_{12} \\ \sigma_{21} & 0 \end{bmatrix}}_{=\Sigma} \mathbf{x}(z, t) \quad (3.118a)$$

defined on the spatial domain $z \in [0, 1]$, with the initial condition $\mathbf{x}(z, 0) = \mathbf{0}$, the boundary conditions

$$\mathbf{x}(0, t) = \mathbf{u}(t) \quad \frac{\partial \mathbf{x}}{\partial z}(1, t) = \underbrace{\begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & \Gamma_{22} \end{bmatrix}}_{=\Gamma} \mathbf{x}(1, t), \quad (3.118b)$$

and the output equation

$$\mathbf{y}(t) = \mathbf{x}(1, t). \quad (3.118c)$$

To obtain a transfer function description of the system (3.118), we apply the bilateral Laplace transform to (3.118a). This gives the spatial ODE

$$\frac{\partial^2}{\partial z^2} \mathbf{x}(z, s) = \underbrace{\Lambda^{-1}(s\mathbf{I} - \Sigma)}_{=\mathbf{H}(s)} \mathbf{x}(z, s). \quad (3.119)$$

The dependence of $\mathbf{H}(s)$ on the Laplace variable s will be omitted for the following calculations. The matrix \mathbf{H} can be diagonalized by using the state transform $\mathbf{x}(z, s) =$

$\mathbf{P}\chi(z, s)$, which yields

$$\frac{\partial^2}{\partial z^2} \chi(z, s) = \tilde{\mathbf{H}}^2 \chi(z, s), \quad (3.120)$$

with the diagonal matrix $\tilde{\mathbf{H}} = \sqrt{\mathbf{P}^{-1} \mathbf{H} \mathbf{P}}$ and the corresponding boundary conditions

$$\chi(0, s) = \mathbf{P}^{-1} \mathbf{u}(s), \quad \frac{\partial \chi}{\partial z}(1, s) = \tilde{\boldsymbol{\Gamma}} \chi(1, s) \quad (3.121)$$

with $\tilde{\boldsymbol{\Gamma}} = \mathbf{P}^{-1} \boldsymbol{\Gamma} \mathbf{P}$. Applying the ansatz

$$\chi(z, s) = \cosh(\tilde{\mathbf{H}}z) \begin{bmatrix} c_{11} \\ c_{21} \end{bmatrix} + \sinh(\tilde{\mathbf{H}}z) \begin{bmatrix} c_{12} \\ c_{22} \end{bmatrix} \quad (3.122)$$

to (3.120), (3.121) yields the solution

$$\mathbf{x}(z, s) = \mathbf{P} [\cosh(\tilde{\mathbf{H}}z) - \sinh(\tilde{\mathbf{H}}z) \boldsymbol{\Xi}] \mathbf{P}^{-1} \mathbf{u}(s), \quad (3.123)$$

with

$$\boldsymbol{\Xi} = [\tilde{\mathbf{H}} \cosh(\tilde{\mathbf{H}}) - \tilde{\boldsymbol{\Gamma}} \sinh(\tilde{\mathbf{H}})]^{-1} [\tilde{\mathbf{H}} \sinh(\tilde{\mathbf{H}}) - \tilde{\boldsymbol{\Gamma}} \cosh(\tilde{\mathbf{H}})]. \quad (3.124)$$

Finally, with (3.118c) the transfer matrix is given by

$$\mathbf{G}_u = \mathbf{P} [\cosh(\tilde{\mathbf{H}}) - \sinh(\tilde{\mathbf{H}}) \boldsymbol{\Xi}] \mathbf{P}^{-1}. \quad (3.125)$$

In the absence of stochastic disturbances, a simple pseudo-inversion-based law (3.80) with $\alpha = 1 \cdot 10^{-3}$ is chosen. Using the parameter values $\lambda_1 = \lambda_2 = 1$, $\sigma_{12} = 1/2$, $\sigma_{21} = -1$, $\Gamma_{11} = 0$, $\Gamma_{12} = 1/2$, $\Gamma_{21} = 1$, and $\Gamma_{22} = 1/10$, one obtains a numerical solution of the impulse response matrix $\mathbf{L}(t)$ using (3.112) as shown in Figure 3.6.

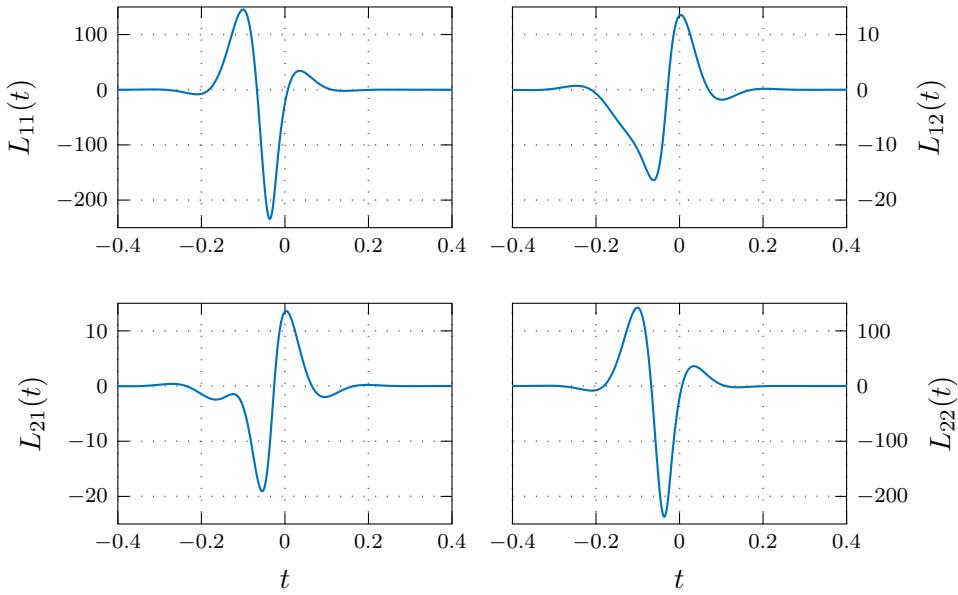


Figure 3.6: Entries of the (non-causal) learning kernel $\mathbf{L}(t)$ calculated as a FIR-approximation using the FFT.

Using the calculated learning kernel to track a desired output trajectory $\mathbf{y}^d(t) = [y_1^d(t), y_2^d(t)]^T$, $t \in [0, 10]$ given by

$$y_1^d(t) = \Theta_5(t - 2.5), \quad y_2^d(t) = 3 \frac{d}{dt}(\Theta_5(t) + \Theta_5(t - 5)) \quad (3.126)$$

with the smoothed step function

$$\Theta_T(t) = \begin{cases} 0 & t \leq 0 \\ 1 & t \geq T \\ \frac{\int_0^t \theta_T(\tau) d\tau}{\int_0^T \theta_T(\tau) d\tau} & t \in (0, T), \end{cases} \quad (3.127)$$

whereby

$$\theta_T(t) = \begin{cases} 0 & t \notin (0, T) \\ \exp\left[-((1 - t/T)t/T)^{-1.5}\right] & t \in (0, T), \end{cases} \quad (3.128)$$

the system (3.118) converges up to numeric precision of the solver after the first iteration with the resulting state profile $\mathbf{x}_1(z, t)$ shown in (3.7). Using higher values for α reduces the learning rate as expected.

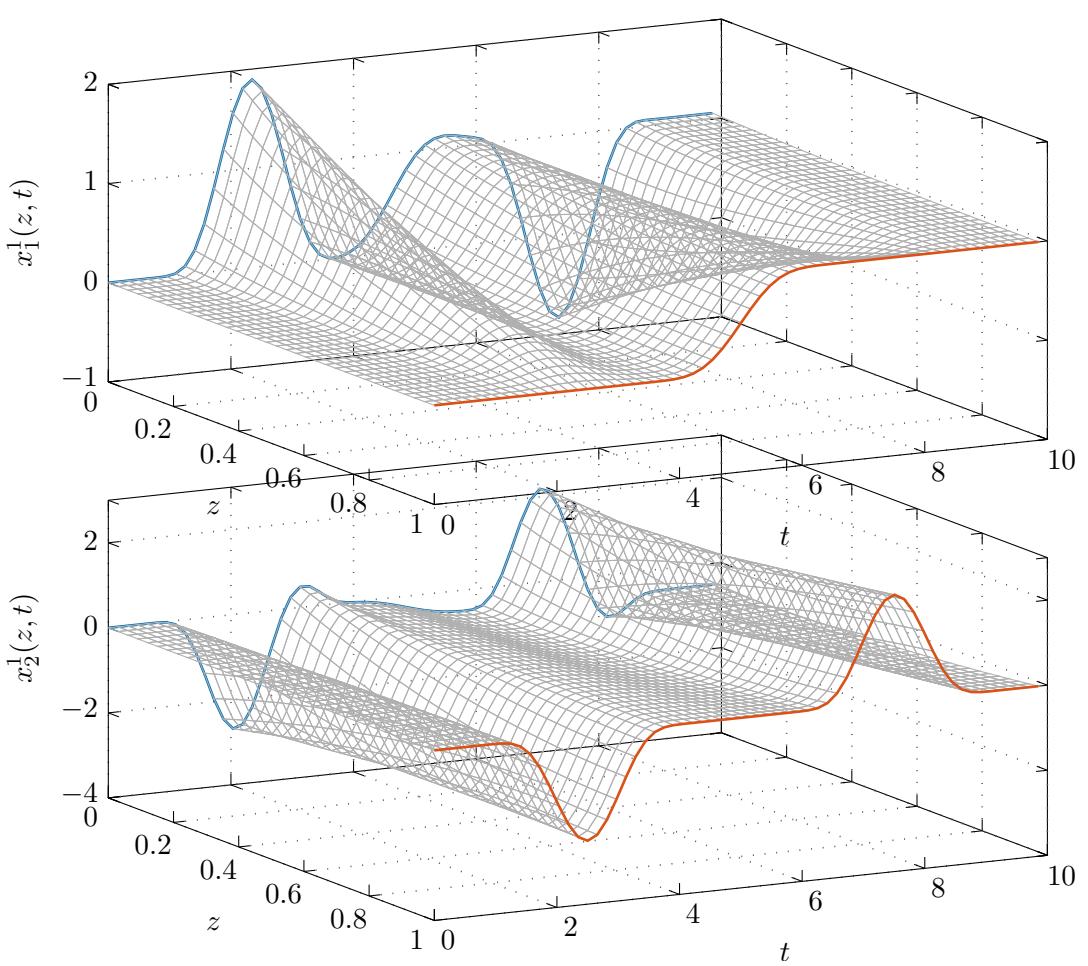


Figure 3.7: State $\mathbf{x}_1(z, t)$ of the PDE system (3.118) with the input $\mathbf{u}_1(t) = \mathbf{x}_1(0, t)$ (blue) that yields $\mathbf{y}_1(t) = \mathbf{x}_1(1, t) \approx \mathbf{y}^d(t)$ (red) after the first iteration.

3.4 Discrete-time systems on finite time horizons

Frequency-domain methods are a very valuable tool to gain insight and intuition, but they are limited to infinite time horizons by design. For the finite time horizon, one has to either analyze the input-output behavior using the (linear) integral operator of the convolution in time-domain or to revert to state-space methods. For discrete-time systems, however, there exist an alternative approach: since the signal spaces of the system input and output are reduced to finite dimensions, every linear mapping between these spaces can be represented as a matrix. This is the basis of the so-called *lifted system representation*, which is commonly used to analyze ILC schemes.

3.4.1 Lifted system representation

In this section, we restrict the general state-space description (3.54) to the SISO case without measurement noise, i.e.,

$$\dot{\mathbf{x}}_j(t) = \mathbf{A}\mathbf{x}_j(t) + \mathbf{b}u_j(t), \quad \mathbf{x}_j(0) = \mathbf{x}_0 \quad (3.129a)$$

$$y_j(t) = \mathbf{c}^T\mathbf{x}_j(t) + du_j(t), \quad (3.129b)$$

Note that the lifted-system representation is in principle open to time-varying and MIMO systems. Using a zero-order hold (ZOH) model with a sampling time T_a , we obtain the discrete time system

$$\mathbf{x}_j[k+1] = \Phi\mathbf{x}_j[k] + \gamma u_j[k], \quad \mathbf{x}_j[0] = \mathbf{x}_0 \quad (3.130a)$$

$$y_j[k] = \mathbf{c}^T\mathbf{x}_j[k] + du_j[k] + v_j[k] \quad (3.130b)$$

with time index $k = 0, 1, \dots$ and the sampled state of the j -th iteration $\mathbf{x}_j[k] = \mathbf{x}_j(kT_a)$, the input $u_j[k] = u_j(kT_a)$, the output $y_j[k] = y_j(kT_a)$, and

$$\Phi = \exp(\mathbf{A}T_a) \quad \text{und} \quad \Gamma = \int_0^{T_a} \exp(\mathbf{A}\tau) d\tau \mathbf{b} = (\exp(\mathbf{A}T_a) - \mathbf{I})\mathbf{A}^{-1}\mathbf{b}. \quad (3.131)$$

The corresponding input-output representation in discrete time reads

$$y_j[0] = \mathbf{c}^T\mathbf{x}_0 + du_j[0], \quad (3.132a)$$

$$y_j[k] = \mathbf{c}^T\Phi^k\mathbf{x}_0 + \mathbf{c}^T \sum_{m=0}^{k-1} (\Phi^{k-m-1}\Gamma u_j[m]) + du_j[k], \quad k = 1, 2, \dots. \quad (3.132b)$$

that can be rewritten using

$$g[k] = \begin{cases} d & \text{for } k = 0 \\ \mathbf{c}^T\Phi^{k-1}\Gamma & \text{for } k = 1, 2, \dots \end{cases}, \quad (3.133)$$

to obtain the discrete-time input-output representation

$$y_j[k] = \mathbf{c}^T\Phi^k\mathbf{x}_0 + \sum_{m=0}^k g[m]u_j[k-m]. \quad (3.134)$$

This equation is the analogous result to (3.55) using the discrete convolution with the impulse response $g[k]$, except that on a finite horizon we cannot assume that the system is relaxed initially and thus there is a remaining contribution of the initial state \mathbf{x}_0 . By introducing the shift operator $\delta : \delta(z_j[k]) = z_j[k+1]$ and assuming $\mathbf{x}_0 = \mathbf{0}$ one can again define the transfer operator

$$\begin{aligned} G(\delta) &= \frac{y_j[k]}{u_j[k]} = d + \mathbf{c}^T(\delta\mathbf{I} - \Phi)^{-1}\gamma = d + \delta^{-1}\mathbf{c}^T(\mathbf{I} - \delta^{-1}\Phi)^{-1}\gamma \\ &= d + \delta^{-1}\mathbf{c}^T(\mathbf{I} + \delta^{-1}\Phi^1 + \delta^{-2}\Phi^2 + \dots)\gamma = d + \delta^{-1}\mathbf{c}^T \sum_{k=0}^{\infty} (\delta^{-k}\Phi^k)\gamma \\ &= d + \sum_{k=1}^{\infty} \mathbf{c}^T\Phi^{k-1}\gamma\delta^{-k} = \sum_{k=0}^{\infty} g[k]\delta^{-k}. \end{aligned} \quad (3.135)$$

Note 3.1. The infinite-horizon methods introduced in the previous section for continuous-time systems can be transferred to the discrete-time case directly using (3.135).

Note 3.2. The relative degree of a discrete-time systems is defined as follows:

Definition 3.10. A system (3.130) with $d = 0$ is of relative degree r if

$$(A) \quad \mathbf{c}^T \Phi^k \Gamma = 0, \quad k = 0, 1, \dots, r - 2$$

$$(B) \quad \mathbf{c}^T \Phi^{r-1} \Gamma \neq 0.$$

The relative degree of a discrete-time systems thus merely corresponds to the time index k of $y_j[k]$ at which the input $u_j[0]$ appears for the first time, i.e.,

$$y_j[0] = \mathbf{c}^T \mathbf{x}_0 \tag{3.136a}$$

$$y_j[1] = \mathbf{c}^T \Phi \mathbf{x}_0 + \underbrace{\mathbf{c}^T \gamma}_{=0} u_j[0] \tag{3.136b}$$

$$y_j[2] = \mathbf{c}^T \Phi^2 \mathbf{x}_0 + \underbrace{\mathbf{c}^T \Phi \gamma}_{=0} u_j[0] + \underbrace{\mathbf{c}^T \gamma}_{=0} u_j[1] \tag{3.136c}$$

⋮

$$y_j[r-1] = \mathbf{c}^T \Phi^{r-1} \mathbf{x}_0 + \underbrace{\mathbf{c}^T \sum_{m=0}^{r-2} \Phi^{r-m-1} \gamma}_{=0} u_j[m] \tag{3.136d}$$

$$y_j[r] = \mathbf{c}^T \Phi^r \mathbf{x}_0 + \underbrace{\mathbf{c}^T \sum_{m=0}^{r-1} \Phi^{r-m-1} \gamma}_{\neq 0} u_j[m] \tag{3.136e}$$

The case of a discrete-time system (3.130) obtained by (ZOH-) sampling of a continuous-time system (3.129) begs the questions how the relative degrees of these two systems compare. Using a series expansion of the exponential matrix $\exp(\mathbf{A}\tau)$ yields

$$\mathbf{c}^T \Gamma = \mathbf{c}^T \int_0^{T_a} \sum_{m=0}^{\infty} \frac{(\mathbf{A}\tau)^m}{m!} \mathbf{b} d\tau \tag{3.137a}$$

$$= \mathbf{c}^T \int_0^{T_a} \mathbf{b} + \mathbf{A}\mathbf{b}\tau + \dots + \frac{1}{(n-1)!} \mathbf{A}^{n-1} \mathbf{b} \tau^{n-1} + \mathcal{O}(T_a^n) d\tau \tag{3.137b}$$

$$= \mathbf{c}^T \mathbf{b} T_a + \frac{1}{2!} \mathbf{c}^T \mathbf{A} \mathbf{b} T_a^2 + \dots + \frac{1}{n!} \mathbf{c}^T \mathbf{A}^{n-1} \mathbf{b} T_a^n + \mathcal{O}(T_a^{n+1}) \tag{3.137c}$$

which is always unequal to zero. A discrete-time system derived from a continuous-time system via ZOH is thus always of relative degree $r = 1$ for $d = 0$ (and $r = 0$ for $d \neq 0$).

Before trying to obtain a matrix representation of the system's input-output behavior (3.134), it makes sense to carefully define suitable input and output sequences that shall

be related by this mapping. Specifically, the initial input $u_j[0]$ will only start to act on the output $y_j[m]$ with some temporal delay $m = r + w$, where r is the relative degree of the system (3.130) and w is an additional delay due to sampling, conversion and processing of data that is unavoidable in real-world applications. For a finite time horizon $t \in [0, t_f]$ with $t_f = NT_a$ we will thus consider the input and output sequences

$$u_j[k], \quad k = 0, 1, \dots, N - 1 \quad (3.138a)$$

$$y_j[k], \quad k = m, m + 1, \dots, N + m - 1 \quad (3.138b)$$

as well as a desired output sequence

$$y_d[k], \quad k = m, m + 1, \dots, N + m - 1. \quad (3.139a)$$

Note that this is essentially shifting the time axis of the output sequence relative to the input sequence to compensate for the total system delay (cf. the time delay analysis in Section 3.3.4). Rewriting the corresponding input and output sequences in vector notation, i.e.,

$$\mathbf{u}_j^T = [u_j[0] \ u_j[1] \ \dots \ u_j[N - 1]] \in \mathbb{R}^N \quad (3.140a)$$

$$\mathbf{y}_j^T = [y_j[m] \ y_j[m + 1] \ \dots \ y_j[m + N - 1]] \in \mathbb{R}^N \quad (3.140b)$$

$$\mathbf{y}_d^T = [y_d[m] \ y_d[m + 1] \ \dots \ y_d[m + N - 1]] \in \mathbb{R}^N \quad (3.140c)$$

$$\mathbf{y}_0^T = [y_0[m] \ y_0[m + 1] \ \dots \ y_0[m + N - 1]] \in \mathbb{R}^N \quad (3.140d)$$

$$\mathbf{e}_j^T = \mathbf{y}_d^T - \mathbf{y}_j^T = [e_j[m] \ e_j[m + 1] \ \dots \ e_j[m + N - 1]] \in \mathbb{R}^N, \quad (3.140e)$$

and using (3.134) yields the so-called *lifted system representation*

$$\mathbf{y}_j = \mathbf{y}_0 + \mathbf{G}\mathbf{u}_j \quad (3.141)$$

with the matrix

$$\mathbf{G} = \begin{bmatrix} g[m] & 0 & \dots & 0 \\ g[m + 1] & g[m] & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ g[m + N - 1] & g[m + N - 2] & \dots & g[m] \end{bmatrix} \in \mathbb{R}^{N \times N}. \quad (3.142)$$

For LTI systems this is a so-called *Toeplitz* matrix where the entries \mathbf{G}_{ij} only depend on the difference $i - j$. An equivalent representation for time-varying systems can be found in [3.0]. Note that since \mathbf{G} is a lower triangular matrix due to causality of the system (3.129) and we assured by definition that $g[m] \neq 0$, \mathbf{G} is of full rank and thus always invertible. Analogous to (3.4), a linear ILC law for the lifted system representation is given by

$$\mathbf{u}_{j+1} = \mathbf{Q}(\mathbf{u}_j + \mathbf{L}\mathbf{e}_j) \quad (3.143)$$

with the \mathbf{Q} -filtering matrix

$$\mathbf{Q} = \begin{bmatrix} q[0] & q[-1] & \cdots & q[-(N-1)] \\ q[1] & q[0] & \cdots & q[-(N-2)] \\ \vdots & \vdots & \ddots & \vdots \\ q[N-1] & q[N-2] & \cdots & q[0] \end{bmatrix} \in \mathbb{R}^{N \times N} \quad (3.144)$$

and the learning gain matrix

$$\mathbf{L} = \begin{bmatrix} l[0] & l[-1] & \cdots & l[-(N-1)] \\ l[1] & l[0] & \cdots & l[-(N-2)] \\ \vdots & \vdots & \ddots & \vdots \\ l[N-1] & l[N-2] & \cdots & l[0] \end{bmatrix} \in \mathbb{R}^{N \times N}. \quad (3.145)$$

Exercise 3.6. Reformulate a PD-type learning law analogous to (3.74) with a Gaussian \mathbf{Q} -filter (3.103) in the lifted framework. What options do you have? What about the boundaries of the time horizon?

The system description (3.141) and the learning law (3.143) are very similar to the infinite-horizon case (3.57) and (3.58) except for the constant term \mathbf{y}_0 . Defining $\tilde{\mathbf{y}}_d = \mathbf{y}_d - \mathbf{y}_0$, one obtains at the input iteration

$$\mathbf{u}_{j+1} = \mathbf{\Psi}\mathbf{u}_j + \mathbf{\Lambda}\tilde{\mathbf{y}}_d \quad (3.146)$$

with $\mathbf{\Psi} = \mathbf{Q}(\mathbf{I} - \mathbf{LG})$ and $\mathbf{\Lambda} = \mathbf{QL}$ and the corresponding iteration of the output error $\mathbf{e}_j = \mathbf{y}_d - \mathbf{y}_j$ as

$$\mathbf{e}_{j+1} = \mathbf{G}\mathbf{\Psi}\mathbf{G}^{-1}\mathbf{e}_j + (\mathbf{I} - \mathbf{GQG}^{-1})\mathbf{y}_d. \quad (3.147)$$

Both iterations are algebraically identical to the infinite-horizon case. We can thus directly transfer stability and convergence results to the lifted system representation, which are restated in the following for completeness.

Theorem 3.13 (Asymptotic stability of the ILC law). *The input iteration (3.146) of the ILC law (3.143) is asymptotically stable if*

$$\rho(\mathbf{Q}(\mathbf{I} - \mathbf{LG})) < 1 \quad (3.148)$$

and \mathbf{u}_j converges to \mathbf{u}_∞ .

Theorem 3.14 (Asymptotic stability of the output iteration). *The output iteration (3.147) of the ILC law (3.143) is asymptotically stable iff the input iteration is stable, i.e.,*

$$\rho(\mathbf{Q}(\mathbf{I} - \mathbf{LG})) < 1 \quad (3.149)$$

and \mathbf{e}_j then converges to the asymptotic tracking error

$$\mathbf{e}_\infty = \left(\mathbf{I} - \mathbf{G}(\mathbf{I} - \boldsymbol{\Psi})^{-1} \boldsymbol{\Lambda} \right) \mathbf{y}_d = \left(\mathbf{I} - \mathbf{G}\boldsymbol{\Phi}\mathbf{G}^{-1} \right)^{-1} \left(\mathbf{I} - \mathbf{G}\mathbf{Q}\mathbf{G}^{-1} \right) \mathbf{y}_d . \quad (3.150)$$

Theorem 3.15 (Monotonic convergence of the input iteration). *The input iteration (3.146) of the ILC law (3.143) converges monotonically to \mathbf{u}_∞ , i.e., it holds that*

$$\|\mathbf{u}_{j+1} - \mathbf{u}_\infty\| \leq \alpha \|\mathbf{u}_j - \mathbf{u}_\infty\| \quad (3.151)$$

for $0 \leq \alpha < 1$ if

$$\|\boldsymbol{\Psi}\| = \bar{\sigma}(\mathbf{Q}(\mathbf{I} - \mathbf{L}\mathbf{G})) = \alpha < 1 . \quad (3.152)$$

Theorem 3.16 (Monotonic convergence of the output iteration). *The output iteration (3.64) of the ILC law (3.58) converges monotonically to \mathbf{e}_∞ , i.e., it holds that*

$$\|\mathbf{e}_{j+1} - \mathbf{e}_\infty\| \leq \beta \|\mathbf{e}_j - \mathbf{e}_\infty\| \quad (3.153)$$

for $0 \leq \alpha < 1$ if

$$\|\mathbf{G}\boldsymbol{\Psi}\mathbf{G}^{-1}\| = \bar{\sigma}(\mathbf{G}\mathbf{Q}(\mathbf{I} - \mathbf{L}\mathbf{G})\mathbf{G}^{-1}) = \beta < 1 . \quad (3.154)$$

Note 3.3. These results are structurally very similar to the infinite time horizon case with a number of significant differences: By applying the Laplace transform on infinite time horizons, one is considering stability for every $\omega \in \mathbb{R}$ independently. While stability can be transferred to the finite time horizon as shown in the previous section, frequency-domain criteria are rather conservative. Conversely, stability criteria using the lifted system representation are sharp by accurately accounting for boundary effects. The dimension of \mathbf{Q} and \mathbf{L} is determined by the length of the sampled time horizon N , which can be problematic for long time horizons.

3.4.2 ILC as an online optimization strategy

Using measurements of a system's behavior to iteratively improve its performance with respect to some cost function can also be seen as an optimization problem that is solved online. Consider the problem

$$\min_{\mathbf{u}} \frac{1}{2} \mathbf{e}^T \mathbf{P} \mathbf{e} + \frac{1}{2} \mathbf{u}^T \mathbf{W} \mathbf{u} + \mathbf{u}^T \mathbf{F} \mathbf{e} \quad (3.155a)$$

$$\text{subject to } \mathbf{e} = \mathbf{y}_d - \mathbf{G} \mathbf{u} , \quad (3.155b)$$

with the symmetric, positive (semi-) definite weighting matrices \mathbf{P} und \mathbf{W} . One can assume that $\mathbf{F}\mathbf{G}$ is a skew-symmetric matrix without loss of generality since any symmetric component could always be absorbed into \mathbf{W} . This becomes apparent when plugging

(3.155b) into (3.155a), which yields the equivalent problem

$$\min_{\mathbf{u}} J(\mathbf{u}) = \frac{1}{2} \mathbf{u}^T \bar{\mathbf{A}} \mathbf{u} + \mathbf{u}^T \bar{\mathbf{b}} + \bar{\mathbf{c}} \quad (3.156)$$

with

$$\bar{\mathbf{A}} = \mathbf{G}^T \mathbf{P} \mathbf{G} + \mathbf{W} \quad (3.157a)$$

$$\bar{\mathbf{b}} = (\mathbf{F} - \mathbf{G}^T \mathbf{P}) \mathbf{y}_d \quad (3.157b)$$

$$\bar{\mathbf{c}} = \frac{1}{2} \mathbf{y}_d^T \mathbf{P} \mathbf{y}_d. \quad (3.157c)$$

The quadratic expression $\mathbf{u}^T \mathbf{F} \mathbf{G} \mathbf{u}$ vanishes due to $\mathbf{F} \mathbf{G}$ being skew-symmetric. Since we further assumed that \mathbf{P} and \mathbf{W} are symmetric and positive (semi-) definite matrices, the same holds true for $\bar{\mathbf{A}}$. In case $\bar{\mathbf{A}}$ is indeed positive definite, the cost function $J(\mathbf{u})$ is strictly convex and thus has a unique global minimum at $\mathbf{u}_\infty = -\bar{\mathbf{A}}^{-1} \bar{\mathbf{b}}$ which is determined by the necessary and sufficient first-order condition $\nabla J(\mathbf{u}_\infty) = \bar{\mathbf{A}} \mathbf{u}_\infty + \bar{\mathbf{b}} = \mathbf{0}$.

Alternatively, such an optimization problem can be solved iteratively using a gradient-descent method with constant step width α , e.g.,

$$\mathbf{u}_{j+1} = \mathbf{u}_j - \alpha \nabla J(\mathbf{u}_j) = (\mathbf{I} - \alpha \bar{\mathbf{A}}) \mathbf{u}_j - \alpha \bar{\mathbf{b}}, \quad (3.158)$$

where $0 < \alpha < 2/\|\bar{\mathbf{A}}\|$ ensures (monotone) convergence, cf. Theorem 2.1.13 in [3.0]. Plugging (3.157b) into (3.158) and suppressing \mathbf{y}_d in favor of \mathbf{u}_j and \mathbf{e}_j using (3.155b) yields an ILC-like update law

$$\mathbf{u}_{j+1} = \mathbf{Q}(\mathbf{u}_j + \mathbf{L} \mathbf{e}_j) \quad (3.159)$$

with

$$\mathbf{Q} = \mathbf{I} - \alpha(\mathbf{W} + \mathbf{F} \mathbf{G}) \quad (3.160a)$$

$$\mathbf{Q} \mathbf{L} = \alpha(\mathbf{G}^T \mathbf{P} - \mathbf{F}). \quad (3.160b)$$

3.4.3 Norm-optimal ILC strategies

Along the same line of thought, norm-optimal ILC methods [3.0] avoid the explicit design of a Q -filter and learning filter by solving the optimization problem

$$\min_{\mathbf{u}_{j+1}} J(\mathbf{u}_{j+1}) = \underbrace{\frac{1}{2} \mathbf{e}_{j+1}^T \mathbf{V} \mathbf{e}_{j+1}}_{J_1(\mathbf{u}_{j+1})} + \underbrace{\frac{1}{2} \mathbf{u}_{j+1}^T \mathbf{S} \mathbf{u}_{j+1}}_{J_2(\mathbf{u}_{j+1})} + \frac{1}{2} (\mathbf{u}_{j+1} - \mathbf{u}_j)^T \mathbf{R} (\mathbf{u}_{j+1} - \mathbf{u}_j) \quad (3.161)$$

$$\text{u.B.v. } \mathbf{e}_{j+1} = \mathbf{y}_d - \mathbf{y}_{j+1} = \mathbf{e}_j + \mathbf{G} \mathbf{u}_j - \mathbf{G} \mathbf{u}_{j+1}$$

for *every* iteration. Note that this is in contrast to the previous section, where ILC was rewritten as an *iterative solution* of a *single* optimization problem. We assume that \mathbf{V} and \mathbf{S} are symmetric, positive semidefinite matrices and \mathbf{R} , $\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R}$, $\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S}$ are symmetric, positive definite matrices.

Plugging the constraint into the cost function of (3.161) yields

$$\begin{aligned} J_1(\mathbf{u}_{j+1}) &= \frac{1}{2}(\mathbf{e}_j + \mathbf{G}\mathbf{u}_j - \mathbf{G}\mathbf{u}_{j+1})^T \mathbf{V} (\mathbf{e}_j + \mathbf{G}\mathbf{u}_j - \mathbf{G}\mathbf{u}_{j+1}) \\ &= \frac{1}{2} \left(\mathbf{e}_j^T \mathbf{V} \mathbf{e}_j + \mathbf{u}_j^T \mathbf{G}^T \mathbf{V} \mathbf{e}_j - 2\mathbf{u}_{j+1}^T \mathbf{G}^T \mathbf{V} \mathbf{e}_j \right. \\ &\quad \left. + \mathbf{e}_j^T \mathbf{V} \mathbf{G} \mathbf{u}_j + \mathbf{u}_j^T \mathbf{G}^T \mathbf{V} \mathbf{G} \mathbf{u}_j - 2\mathbf{u}_{j+1}^T \mathbf{G}^T \mathbf{V} \mathbf{G} \mathbf{u}_j + \mathbf{u}_{j+1}^T \mathbf{G}^T \mathbf{V} \mathbf{G} \mathbf{u}_{j+1} \right) \end{aligned} \quad (3.162)$$

and

$$\begin{aligned} J_2(\mathbf{u}_{j+1}) &= \frac{1}{2} \mathbf{u}_{j+1}^T \mathbf{S} \mathbf{u}_{j+1} + \frac{1}{2} (\mathbf{u}_{j+1}^T - \mathbf{u}_j^T) \mathbf{R} (\mathbf{u}_{j+1} - \mathbf{u}_j) \\ &= \frac{1}{2} \left(\mathbf{u}_{j+1}^T \mathbf{S} \mathbf{u}_{j+1} + \mathbf{u}_{j+1}^T \mathbf{R} \mathbf{u}_{j+1} - 2\mathbf{u}_{j+1}^T \mathbf{R} \mathbf{u}_j + \mathbf{u}_j^T \mathbf{R} \mathbf{u}_j \right). \end{aligned} \quad (3.163)$$

Using the first-order optimality condition

$$\left(\frac{\partial}{\partial \mathbf{u}_{j+1}} J \right) (\mathbf{u}_{j+1}) = \mathbf{0} \quad (3.164)$$

results in input update

$$(\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S} + \mathbf{R}) \mathbf{u}_{j+1} = (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R}) \mathbf{u}_j + \mathbf{G}^T \mathbf{V} \mathbf{e}_j \quad (3.165)$$

that is equivalent to an ILC law (3.143) with \mathbf{Q} -filtering and learning matrices

$$\mathbf{Q} = (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S} + \mathbf{R})^{-1} (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R}) \quad (3.166a)$$

$$\mathbf{L} = (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R})^{-1} \mathbf{G}^T \mathbf{V}. \quad (3.166b)$$

By considering

$$\begin{aligned} \mathbf{Q}(\mathbf{I} - \mathbf{L}\mathbf{G}) &= (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S} + \mathbf{R})^{-1} (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R}) \left(\mathbf{I} - (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{R})^{-1} \mathbf{G}^T \mathbf{V} \mathbf{G} \right) \\ &= (\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S} + \mathbf{R})^{-1} \mathbf{R}. \end{aligned} \quad (3.167)$$

one can show that the derived norm-optimal ILC scheme is asymptotically stable, i.e.,

$$\rho((\mathbf{G}^T \mathbf{V} \mathbf{G} + \mathbf{S} + \mathbf{R})^{-1} \mathbf{R}) < 1. \quad (3.168)$$

The values of \mathbf{V} , \mathbf{S} and \mathbf{R} are tuning parameters that are often simplified by using diagonal matrices, i.e., $\mathbf{V} = v\mathbf{I} > 0$, $\mathbf{S} = s\mathbf{I} > 0$ und $\mathbf{R} = r\mathbf{I} > 0$. It follows that:

- Large values of v increases the weighting of the output error \mathbf{e}_j , which increases the convergence rate and reduces \mathbf{e}_∞ . The learning gain is becoming more aggressive and the action of the \mathbf{Q} -filter is reduced. The limit $v \rightarrow \infty$ yields $\mathbf{L} \rightarrow \mathbf{G}^{-1}$ and $\mathbf{Q} \rightarrow \mathbf{I}$.
- Large values of s penalize the control input \mathbf{u}_j , which increases the asymptotic output error \mathbf{e}_∞ . Note that s only affects \mathbf{Q} where $s \rightarrow 0$ yields $\mathbf{Q} \rightarrow \mathbf{I}$.
- Large values of r penalize changes of the control input $\mathbf{u}_{j+1} - \mathbf{u}_j$, which decreases the convergence rate. For $r \rightarrow 0$ it follows that $\mathbf{L} \rightarrow \mathbf{G}^{-1}$.

3.5 Literatur

- [3.0] Z. Bien and J.-X. Xu, *Iterative learning control - analysis, design, integration and applications*. London: Springer, 1998.
- [3.0] D. Owens and J. Hätönen, “Iterative learning control - an optimization paradigm,” *Annual Reviews in Control*, vol. 29, no. 1, pp. 57–70, 2005.
- [3.0] D. A. Bristow, M. Tharayil, and A. Alleyne, “A survey of iterative learning control,” *Control Systems, IEEE*, vol. 26, no. 3, pp. 96–114, 2006.
- [3.0] H.-S. Ahn, Y.-Q. Chen, and K. Moore, “Iterative learning control: Brief survey and categorization,” *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, vol. 37, no. 6, pp. 1099–1121, 2007.
- [3.0] H.-S. Ahn, K. L. Moore, and Y. Chen, *Iterative learning control - robustness and monotonic convergence for interval system*. London: Springer, 2007.
- [3.0] Y. Wang, F. Gao, and F. J. Doyle III, “Survey on iterative learning control, repetitive control, and run-to-run control,” *Journal of Process Control*, vol. 19, no. 10, pp. 1589–1600, 2009.
- [3.0] J. M. Ortega, “Stability of difference equations and convergence of iterative processes,” *SIAM Journal on Numerical Analysis*, vol. 10, no. 2, pp. 268–282, 1973.
- [3.0] D. A. Bristow, “Iterative learning control for precision motion of microscale and nanoscale tracking systems,” Ph.D. dissertation, University of Illinois at Urbana-Champaign, 2007.
- [3.0] J. Ghosh and B. Paden, “Iterative learning control for nonlinear nonminimum phase plants with input disturbances,” in *Proceedings of the American Control Conference*, Jun. 1999, pp. 2584–2589. (visited on 05/24/2016).
- [3.0] J. Ghosh and B. Paden, “A pseudoinverse-based iterative learning control,” *IEEE Transactions on Automatic Control*, vol. 47, no. 5, pp. 831–837, 2002.
- [3.0] J. Wallén, S. Gunnarsson, and M. Norrlöf, “Analysis of boundary effects in iterative learning control,” *International Journal of Control*, vol. 86, no. 3, pp. 410–415, 2013.
- [3.0] J. H. Lee, K. S. Lee, and W. C. Kim, “Model-based iterative learning control with a quadratic criterion for time-varying linear systems,” *Automatica*, vol. 36, no. 5, pp. 641–657, 2000.
- [3.0] Y. Nesterov, *Introductory lectures on convex optimization: A basic course*. Dordrecht: Kluwer Academic Publishers, 2004.
- [3.0] K. Barton and A. Alleyne, “A norm optimal approach to time-varying ilc with application to a multi-axis robotic testbed,” *Control Systems Technology, IEEE Transactions on*, vol. 19, no. 1, pp. 166–180, Jan. 2011.

4 Stochastic optimal control and reinforcement learning

There are different types of machine learning including supervised learning, self-supervised learning and unsupervised learning. *Reinforcement learning refers to a learner or agent that interacts with its environment and modifies its actions, or control policies, based on stimuli received in response to its actions.* This is based on evaluative information from the environment and could be called action-based learning. Reinforcement learning implies a cause and effect relationship between actions and reward. It implies goal directed behavior at least insofar as the agent has an understanding of reward versus lack of reward or punishment. *Optimal Control (OC) and Reinforcement Learning (RL)* both deal with *sequential decision-making in deterministic and stochastic environments*, but they approach the problem from different perspectives and have distinct characteristics:

Foundation and historical roots:

- OC: Stems from the field of control theory. It traditionally focuses on the mathematical modeling and understanding of systems, often with a well-defined model of the environment in mind.
- RL: Emerged from the realm of artificial intelligence and, to some extent, psychology, inspired by behavioral learning theories. RL often operates in scenarios where the model of the environment is unknown.

Knowledge about the environment:

- OC: Typically assumes a known model of the environment. This model describes state transitions and rewards (or costs) as a function of states and actions.
- RL: Can operate with or without an explicit model of the environment. Model-free RL methods, like Q-learning or policy gradient methods, directly learn a policy or value function without needing to know the dynamics or reward structure.

Solution techniques:

- OC: Methods in OC, such as dynamic programming (Value Iteration, Policy Iteration), rely heavily on the known model to find the optimal policy.
- RL: Employs a wider range of techniques, including but not limited to dynamic programming. Many algorithms, especially model-free ones, rely on interaction with the environment to learn. RL also incorporates deep learning (Deep RL) to handle large-scale problems with high-dimensional input spaces.

Application:

- OC: Traditionally applied to problems with smaller state and action spaces due to the computational intensity of exact methods. However, approximations are possible for larger problems.
- RL: Given its roots in AI, RL has been widely applied to large-scale and high-dimensional problems, especially with the advent of deep learning. Examples include playing Atari games, Go, and various robotics tasks.

Exploration vs. exploitation:

- OC: When formulated classically, the OC problem typically assumes full knowledge of the system, and hence the concepts of exploration and exploitation are not central.
- RL: The trade-off between exploration (trying new actions to discover their effects) and exploitation (choosing known good actions) is a central theme in RL, especially in model-free settings.

Goal:

- OC: The primary goal is to find the optimal policy (or control law) given the system dynamics and reward function.
- RL: While the end goal is also to find an optimal policy, RL places a significant emphasis on learning from interaction. This can involve learning about the environment's dynamics, the reward structure, or directly about the optimal policy or value function.

4.1 Theoretical foundation

Since we are in a stochastic setting, it is important to understand the *basics of probability theory*. We use sans-serif letters such as x , \mathbf{x} , and \mathbf{X} to represent *random variables* (scalars, vectors and matrices) and serif letters such as x , \mathbf{x} , and \mathbf{X} to represent the corresponding *deterministic variables* or events. For an introduction to probability theory, please refer to the brief introduction given in the Appendix A and for a consistent summary to [4.0]. A comprehensive overview can be found in [4.0].

4.1.1 Stochastic dynamical systems

We consider stochastic, nonlinear dynamical systems with discrete time evolution of the form

$$\mathbf{x}_{k+1} = \mathbf{f}_d(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w}_k), \quad k = 0, \dots, N-1 , \quad (4.1)$$

and initial condition \mathbf{x}_0 drawn from the initial condition distribution $p_{\mathbf{x}_0}(\mathbf{x}_0)$. Here,

- k denotes the discrete time or epoch index,

- $\mathbf{x}_k \in \mathcal{X} \subset \mathbb{R}^n$ is the stochastic state vector and
- $\mathbf{u}_k \in \mathcal{U}(\mathbf{x}_k) \subset \mathcal{U} \subset \mathbb{R}^m$ is the stochastic input vector. The input vector \mathbf{u}_k is restricted to a nonempty subset $\mathcal{U}(\mathbf{x}_k) \subset \mathcal{U}$ that depends on \mathbf{x}_k .
- $\mathbf{w}_k \in \mathcal{W}$ is the disturbance vector.
- The horizon length is given by N .
- The mapping $\mathbf{f}_d : \mathcal{X} \times \mathcal{U} \times \mathcal{W} \rightarrow \mathcal{X}$ describes how the system state \mathbf{x} evolves over time.
- The sequence $\mathbf{W}_{[0:N-1]} = (\mathbf{w}_0, \mathbf{w}_1, \dots, \mathbf{w}_{N-1})$ is assumed to be uncorrelated.

Let us use $p(\mathbf{i}_k)$ to denote the probability distribution, with observation \mathbf{i}_k , of the stochastic vector $\mathbf{i}_k \in \mathcal{I}$ at time index k , i.e., $\mathbf{i}_k \sim p(\mathbf{i}_k)$. A probability distribution is a description of how likely a stochastic variable is to take on each of its possible states. The method of describing this distribution varies based on the nature of the variable, being either discrete or continuous. Using this concept, it is actually possible to write the dynamics of the distribution with the *law of total probability*, cf. (A.12), as

$$\begin{aligned} p(\mathbf{x}_{k+1}) &= \mathbb{E}_{\mathbf{x}_k, \mathbf{u}_k \sim p(\mathbf{x}_k, \mathbf{u}_k)} \{ p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) \} \\ &= \int_{\mathcal{X}} \int_{\mathcal{U}} p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) p(\mathbf{x}_k, \mathbf{u}_k) d\mathbf{u}_k d\mathbf{x}_k , \end{aligned} \quad (4.2)$$

where $p_x : \mathcal{X} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ encodes the stochastic dynamics as a conditional distribution of the next state \mathbf{x}_{k+1} as a function of the current state \mathbf{x}_k when applying the control input \mathbf{u}_k [4.0, p. 13]. Thus, (4.1) can be alternatively written as

$$\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k), \quad k = 0, 1, \dots, N-1 . \quad (4.3)$$

The fact that at this point two expressions (4.1) and (4.3) are given for the system dynamics is due to the fact that, depending on the problem, one of the two representations is to be preferred in order to obtain mathematically consistent and clear expressions. In the control engineering context, (4.1) is preferred over (4.3), whereas (4.3) is widely used in the reinforcement learning context. We will utilize both notations. Note that by construction, the system (4.1) satisfies the Markov property, see [4.0], which states that \mathbf{x}_{k+1} depends only on quantities of the previous time step k , i.e.

$$p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) = p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_k, \mathbf{u}_k), \quad k = 0, \dots, N-1 . \quad (4.4)$$

Note 4.1. In the reinforcement learning literature, the control input \mathbf{u}_k is referred to as action \mathbf{a}_k and the state \mathbf{x}_k is denoted by \mathbf{s}_k .

4.1.2 Rollouts, rewards and return

Let us introduce the *state sequence*

$$\mathbf{X}_{[0:N]} = (\mathbf{x}_0, \dots, \mathbf{x}_N) \quad (4.5)$$

and *input sequence*

$$\mathbf{U}_{[0:N-1]} = (\mathbf{u}_0, \dots, \mathbf{u}_{N-1}), \quad (4.6)$$

which describes the time evolution of the state and input variables. We denote

$$\boldsymbol{\tau}_{[0,N]} = (\mathbf{x}_0, \mathbf{u}_0, \mathbf{x}_1, \mathbf{u}_1, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}, \mathbf{x}_N) \quad (4.7)$$

as the N-step *rollout*¹. In optimal control, we are typically interested in minimizing a total cost. In the reinforcement learning literature, the reward is used instead of the cost. Consequently, the accumulative costs (total cost) are minimized and accumulative rewards (return) are maximized. Throughout this work we will adopt the logic and notation commonly found in the reinforcement learning literature. Note that the rewards can be positive and negative.

Definition 4.1. (Discounted return, see [4.0]) The *future discounted return* of a rollout $\boldsymbol{\tau}_{[k:N]}$ is

$$R_k(\boldsymbol{\tau}_{[k:N]}) = \sum_{i=k}^N \gamma^{i-k} r_i = \gamma^{N-k} r_N + \sum_{i=k}^{N-1} \gamma^{i-k} r_i \quad (4.8)$$

with discount factor $\gamma \in [0, 1]$. It discounts and accumulates the *random immediate reward* $r_k : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ as well as the *random terminal reward* $r_N : \mathcal{X} \rightarrow \mathbb{R}$.

The smaller γ is chosen, the lower future rewards in r_k will be weighted, thereby giving more significance to immediate rewards. Conversely, rewards received for high occupancies in the future are weighted more heavily. Thus, this parameter essentially represents the distance of foresight into the future. Note that if R_k is the value received at time index k , then

$$\begin{aligned} R_k &= r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots + \gamma^{N-k-1} r_{N-1} + \gamma^{N-k} r_N \\ &= r_k + \gamma(r_{k+1} + \gamma r_{k+2} + \gamma^2 r_{k+3} + \dots + \gamma^{N-k-2} r_{N-1} + \gamma^{N-k-1} r_N) \\ &= r_k + \gamma R_{k+1}. \end{aligned} \quad (4.9)$$

So, the further away a reward is from the initial state \mathbf{x}_k , the less actual reward we will receive from it. For the sake of completeness, we introduce the *reward sequence*

$$\mathbf{r}_{[0:N]} = (r_0, \dots, r_N). \quad (4.10)$$

There is a finite and infinite horizon setting in RL. They mainly differ in the time frame over which the agent plans and makes decisions:

¹Rollouts are also frequently called trajectories.

- In the *finite horizon setting*, the agent operates over a fixed, known number of time steps, i.e., N is finite.
- In the *infinite horizon setting*, the agent considers an infinite number of future steps, i.e., $N = \infty$.

Note 4.2. There is no consensus within the community about whether the reward corresponding to the state \mathbf{x}_k is gained at time index k as in [4.0], or time index $k + 1$, as in [4.0]. Here, it is assumed that the reward is gained at time index k and is denoted by r_k .

There are two settings for learning and optimization:

- The *episodic/sequential setting*, where the experience is broken up into a series of episodes/sequences. It is our goal to maximize the cumulative reward within a single episode. This setting is typical in scenarios like games or simulations where tasks are naturally episodic and reset after reaching a conclusion.
- The *continuing setting*, where the task doesn't have clear episode boundaries. It is our goal to maximize the cumulative reward over an infinite or very long time horizon, i.e. $N = \infty$. This is more reflective of real-world scenarios like stock market trading or ecosystem management, where the process is ongoing and doesn't reset periodically. In this setting, the notion of discounting future rewards often becomes essential to ensure that the cumulative reward doesn't diverge to infinity.

4.1.3 Different terminologies and interaction models

In control engineering, when addressing a (deterministic) *Optimal Control Problem* (OCP), we refer to a simulation model, a cost function, and a numerical solution method to solve the dynamical optimization problem. Hence, we adopt the Optimizer-Model-Cost interaction framework, as illustrated in Figure 4.1a. In reinforcement learning, we consider the agent-environment-reward interaction model², rooted in the principles of the (stochastic) *Markov Decision Process* (MDP). This interaction model is depicted in Figure 4.1b. In RL, the accumulative rewards (return) are maximized and, in OC, the accumulative costs (total cost) are minimized. In control engineering, the environment corresponds to the controlled system, plant or model and the agent to the decision maker or controller. In essence, while there exists a linguistic divergence between control engineering and computational intelligence terminologies, their underlying semantics largely converge. This dichotomy is primarily a byproduct of the distinct historical and academic roots of the two fields. For more information about the different terminologies in Control Systems Engineering and Computational Intelligence, see [4.0, p. 43].

4.1.4 Markov Decision Process

Markov Decision Processes (MDP) formally describe an environment for reinforcement learning. A Markov Process is a memoryless random process, i.e. a sequence of random

²Frequently referred to simply as the agent-environment interaction model.

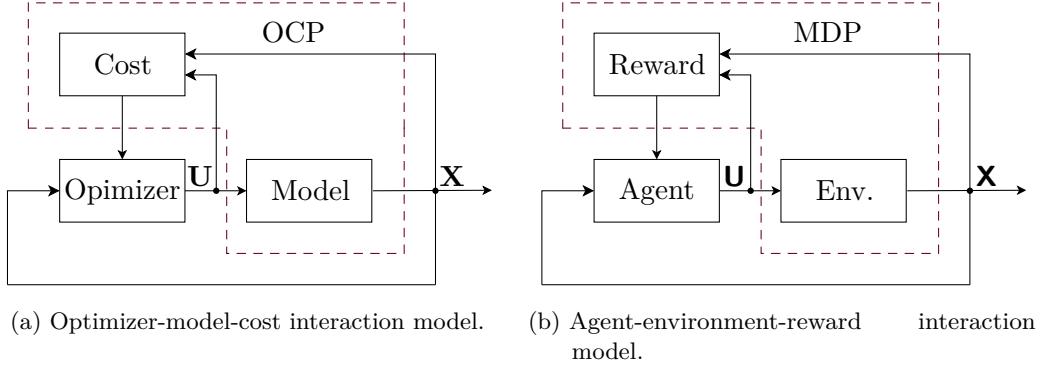


Figure 4.1: Interaction models.

states with the Markov property (4.4).

Definition 4.2. (Markov Decision Process) A (stationary) Markov Decision Process (MDP) is a fully observable, probabilistic state model. A finite-horizon, discount-reward MDP \mathcal{M} is a tuple $\langle \mathcal{X}, \mathcal{U}, \mathcal{R}, p, p_{x_0}, \gamma, N \rangle$ containing:

- \mathcal{X} is the state space
- \mathcal{U} is the input space
- \mathcal{R} is the reward space
- p is the dynamic function
- p_{x_0} is the initial condition function
- $\gamma \in [0, 1]$ is a discount factor and
- N is the horizon length.

Note 4.3. Markov models are categorized as either fully or partially observable as well as either autonomous or non-autonomous, cf. Table 4.1. In the partially observable setting, the agent only has access to the output^a \mathbf{y}_k at each time step k . This model is called *Partially Observable Markov Decision Process* (POMDP) and in addition to the MDP, there is an output-transition probability $o(\mathbf{y}_k | \mathbf{x}_k, \mathbf{u}_k)$. In the autonomous case, there is no input \mathbf{u}_k .

^aCalled observation in the reinforcement learning literature.

The MDP and agent together thereby give rise to an episode

$$\mathbf{x}_0, \mathbf{u}_0, \mathbf{r}_0, \mathbf{x}_1, \mathbf{u}_1, \mathbf{r}_1, \dots, \mathbf{x}_{N-1}, \mathbf{u}_{N-1}, \mathbf{r}_{N-1}, \mathbf{x}_N, \mathbf{r}_N , \quad (4.11)$$

where $\mathbf{x}_k \in \mathcal{X}$, $\mathbf{u}_k \in \mathcal{U}$ and $\mathbf{r}_k \in \mathcal{R}$. The states, inputs, and rewards are either *discrete* or

Table 4.1: Types of Markov Models.

	Fully observable	Partially observable
Auto.	Markov Chain	Hidden Markov Model
Non-auto.	Markov Decision Process	Partially Observable Markov Decision Process

continuous random variables.

- Discrete variables: These are variables that take values on a finite set. For example, inputs in a Tic-Tac-Toe game are discrete because there are a limited number of squares where a player can place a symbol.

Example 4.1 (Tic-Tac-Toe game). In the Tic-Tac-Toe game, the board configuration at any given time serves as the outcomes. Each outcome can be represented as a 3×3 matrix or a vector with 9 elements with entries from $\{X, O, \text{empty}\}$. That's a total of $3^9 = 19683$ different ways the 3×3 grid can be filled in. The sample space is the set of all possible outcomes and can be written as $\xi \in \Xi = \{X, O, \text{empty}\}^9$. For each element x of the random vector \mathbf{x} , a possible mapping from the sample space the real numbers is $x(X) = 1$, $x(O) = 2$, $x(\text{empty}) = 3$. Hence, $\mathbf{x}(\xi) \in \mathbb{R}^9$ and $\mathbf{x} \in \mathcal{X} = \{1, 2, 3\}^9 \subseteq \mathbb{R}^9$. The least number of moves required to win a Tic-Tac-Toe game is 5, while the maximum is 9. Assuming player one, using X, goes first and player two, using O, follows, the game proceeds with alternating moves until the board is filled with five X's and four O's. The inputs are discrete and correspond to the placement of $\{X, O\}$ in an empty cell. At the start of the game, there are 9 possible inputs, one for each cell on the board. Thus, the first player has $9 + 7 + 5 + 3 + 1 = 25$ and the second has $8 + 6 + 4 + 2 = 20$ possible moves, not accounting for games ending before the board is full. In a Tic-Tac-Toe game, the reward structure is often based on the game's outcome, which is inherently discrete. The goal of the game is to get three in a row - horizontally, vertically, or diagonally. Play continues until a player achieves this goal or all the spaces are filled with X's and O's. Win-Lose-Draw: A simple way to define the reward is to give a positive value for a win (1), a negative value for a loss (-1), and a smaller value or zero for a draw (0), i.e. $\mathcal{R} = \{1, -1, 0\}$, with $|\mathcal{R}| = 3$. In more sophisticated models, future rewards can be discounted to prioritize short-term gains and a negative reward could be given for each move to encourage the agent to win in fewer steps. This discussion reveals that even elementary games such as Tic-Tac-Toe are founded on complex mathematical concepts.

Example 4.2 (Frozen Lake). In the Frozen Lake game, the environment is a grid composed of safe tiles (frozen surface) and dangerous tiles (holes). For example, it can be represented as a 4×4 matrix or a vector with 16 elements. In the environment's code, each tile is represented by a letter as follows: (S: starting point, safe), (F: frozen surface, safe), (H: hole, stuck forever) and (G: goal, safe).

<table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>X</td><td>O</td><td>X</td></tr> <tr><td>X</td><td>X</td><td>O</td></tr> <tr><td>O</td><td>O</td><td>X</td></tr> </table>	X	O	X	X	X	O	O	O	X	<table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>X</td><td>O</td><td>X</td></tr> <tr><td>X</td><td>O</td><td></td></tr> <tr><td>O</td><td>O</td><td>X</td></tr> </table>	X	O	X	X	O		O	O	X	<table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>X</td><td>O</td><td>X</td></tr> <tr><td>X</td><td>O</td><td>O</td></tr> <tr><td>O</td><td>X</td><td>X</td></tr> </table>	X	O	X	X	O	O	O	X	X
X	O	X																											
X	X	O																											
O	O	X																											
X	O	X																											
X	O																												
O	O	X																											
X	O	X																											
X	O	O																											
O	X	X																											
(a) Win $r = 1$.	(b) Loss $r = -1$.	(c) Draw $r = 0$.																											

Figure 4.2: Reward structure of the Tic-Toc-Toe game.

Each tile can take one of four values $\{S, F, H, G\}$. By default, the environment is always in the same configuration. The agent A must navigate from the starting point (S) to the goal (G) without falling into holes. The agent's location in the grid is the state and hence there are $16 = |\mathcal{X}| = \{0, 1, \dots, 15\}$ possible states. Possible inputs are $\{\leftarrow, \downarrow, \rightarrow, \uparrow\}$. We can map the inputs to the real line via $\{u(\leftarrow) = 0, u(\downarrow) = 1, u(\rightarrow) = 2, u(\uparrow) = 3\}$. A sequence of inputs leads the agent to the goal. The agent must learn the optimal path, avoiding holes, and reaching the goal with a minimum number of moves. The game is deterministic in the non-slippery version. In the slippery version, the action the agent takes only has 33% chance of succeeding. In case of failure, one of the three other inputs is randomly taken instead.

<table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>F</td><td>F</td><td>F</td><td>G</td></tr> <tr><td>F</td><td>H</td><td>H</td><td>F</td></tr> <tr><td>F</td><td>F</td><td>H</td><td>F</td></tr> <tr><td>S</td><td>F</td><td>H</td><td>H</td></tr> </table>	F	F	F	G	F	H	H	F	F	F	H	F	S	F	H	H	<table border="1" style="margin-left: auto; margin-right: auto;"> <tr><td>0</td><td>1</td><td>2</td><td>3</td></tr> <tr><td>4</td><td>5</td><td>6</td><td>7</td></tr> <tr><td>8</td><td>9</td><td>10</td><td>11</td></tr> <tr><td>12</td><td>13</td><td>14</td><td>15</td></tr> </table>	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
F	F	F	G																														
F	H	H	F																														
F	F	H	F																														
S	F	H	H																														
0	1	2	3																														
4	5	6	7																														
8	9	10	11																														
12	13	14	15																														
(a) Configuration of the environment.	(b) Current location of the agent as state encoding.																																

Figure 4.3: 4×4 Frozen Lake game.

- Continuous variables: These are variables that can take an infinite number of values within a given range. For example, the angle of a rotatory pendulum is a continuous state because it can vary within a range.

Example 4.3 (Rotatory pendulum). The state of a pendulum is typically characterized by its angle θ and angular velocity ω . These variables are continuous and can vary within a specific range, i.e. $\mathbf{x}^\top = [\theta, \omega] \in \mathbb{R}^2$ and

$\mathcal{X} = [0 \text{ rad}, 2\pi \text{ rad}] \times [-2 \text{ rad/s}, 2 \text{ rad/s}]$ for example. The input for a pendulum could be the torque τ applied to it, which is also a continuous variable, i.e. $\mathbf{u} = \tau \in \mathbb{R}^1$ and $\mathcal{U} = [-1 \text{ Nm}, 1 \text{ Nm}]$. A negative reward can be given based on the angular deviation from the upright position, usually $\theta = 0 \text{ rad}$. The closer the pendulum is to being upright, the smaller the absolute value of the negative reward. To encourage efficient control, the reward can also include a term that penalizes large torque inputs, hence a possible deterministic reward is $r = -(|\theta| + d\tau^2)$, with coefficient $d > 0$.

Value-discrete case

Suppose \mathcal{X} , \mathcal{U} , and \mathcal{R} are *finite sets* of cardinality $|\mathcal{X}|$, $|\mathcal{U}|$, and $|\mathcal{R}|$. Let us assume that the random variables \mathbf{x}' and \mathbf{r} have well defined discrete probability distributions dependent only on the preceding state \mathbf{x} and input \mathbf{u} . That is, for particular values of these random variables, there is a *dynamic function* $\mathbf{p} : \mathcal{X} \times \mathcal{R} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$, which is equivalent to the conditional probability $\mathbb{P}\{\cdot | \cdot\}$. It allows us to predict the future state \mathbf{x}' and current reward r :

$$\mathbf{p}(\mathbf{x}', r | \mathbf{x}, \mathbf{u}) = \mathbb{P}\{\mathbf{x}_{k+1} = \mathbf{x}', \mathbf{r}_k = r | \mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}\}. \quad (4.12)$$

The dynamic function gives rise to the *state-transition function* $\mathbf{p}_x : \mathcal{X} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$, which is given by

$$\mathbf{p}_x(\mathbf{x}' | \mathbf{x}, \mathbf{u}) = \sum_{r \in \mathcal{R}} \mathbf{p}(\mathbf{x}', r | \mathbf{x}, \mathbf{u}) \quad (4.13)$$

and the *reward function* $\mathbf{p}_r : \mathcal{R} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$, which is

$$\mathbf{p}_r(r | \mathbf{x}, \mathbf{u}) = \sum_{\mathbf{x}' \in \mathcal{X}} \mathbf{p}(\mathbf{x}', r | \mathbf{x}, \mathbf{u}). \quad (4.14)$$

The function \mathbf{p}_x defines the dynamics of the MDP

$$\mathbf{x}' \sim \mathbf{p}_x(\mathbf{x}' | \mathbf{x}, \mathbf{u}) \quad (4.15)$$

and the function \mathbf{p}_r gives the rewards of the MDP

$$\mathbf{r} \sim \mathbf{p}_r(r | \mathbf{x}, \mathbf{u}). \quad (4.16)$$

In an unknown environment, we do not have perfect knowledge about \mathbf{p}_x and \mathbf{p}_r . For MDPs with a *discrete state space*, with a little abuse of notation³, we determine the *expected rewards* for state-input pairs i -and- l as a two-argument function $r : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$:

$$r_i^{\mathbf{u}_l} = r(\mathbf{x}_i, \mathbf{u}_l) = \mathbb{E}_{\mathbf{r} \sim \mathbf{p}_r(r | \mathbf{x}_i, \mathbf{u}_l)}\{\mathbf{r}\} = \sum_{r \in \mathcal{R}} r \mathbf{p}_r(r | \mathbf{x}_i, \mathbf{u}_l). \quad (4.17)$$

³To avoid introducing a second index, we denote a specific value-discrete state by $\mathbf{x}_i \in \mathcal{X}$ and a specific value-discrete input by $\mathbf{u}_l \in \mathcal{U}$. In contrast, \mathbf{x}_k and \mathbf{u}_k represent the state and input at a specific time index k .

We define the *i-to-j-for-l state-transition probabilities* as

$$p_{ij}^{\mathbf{u}_l} = \mathbf{p}_{\mathbf{x}}(\mathbf{x}_j \mid \mathbf{x}_i, \mathbf{u}_l) = \mathbb{P}[\mathbf{x}_{k+1} = \mathbf{x}_j, \mid \mathbf{x}_k = \mathbf{x}_i, \mathbf{u}_k = \mathbf{u}_l] , \quad (4.18)$$

which can be summarized in the *state-transition probability tensor* \mathbf{P} , with tensor elements

$$\mathbf{P}^{\mathbf{u}_l}[i, j] = p_{ij}^{\mathbf{u}_l} . \quad (4.19)$$

The rows of the state-transition probability sum up to one for each input \mathbf{u}_l , i.e.

$$\sum_{j=0}^{|\mathcal{X}|-1} p_{ij}^{\mathbf{u}_l} = 1 . \quad (4.20)$$

for all $\mathbf{x}_i \in \mathcal{X}$ and $\mathbf{u}_l \in \mathcal{U}$.

Example 4.4. For $\mathcal{X} = \{x_0, x_1, x_2\}$ and $\mathcal{U} = \{u_0\}$, with state-transition probability matrix

$$\mathbf{P}^{u_0} = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} = \begin{bmatrix} p_{00}^{u_0} & 0 & 0 \\ p_{10}^{u_0} & 0 & 0 \\ p_{20}^{u_0} & 0 & 0 \end{bmatrix} , \quad (4.21)$$

we can draw a state-transition diagram shown in Figure 4.4.

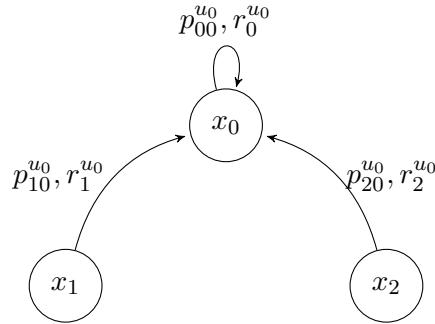


Figure 4.4: A state-transition diagram.

Value-continuous case

Now suppose \mathcal{X} , \mathcal{U} , and \mathcal{R} are *continuous sets*. Then, the *dynamic function* induces a *state-transition function* $\mathbf{p}_{\mathbf{x}} : \mathcal{X} \times \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ of the form

$$\mathbf{p}_{\mathbf{x}}(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) = \int_{\mathcal{R}} \mathbf{p}(\mathbf{x}_{k+1}, r_k \mid \mathbf{x}_k, \mathbf{u}_k) dr_k \quad (4.22)$$

and a *reward function* $\mathbf{p}_r : \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$, which is

$$\mathbf{p}_r(r_k \mid \mathbf{x}_k, \mathbf{u}_k) = \int_{\mathcal{X}} \mathbf{p}(\mathbf{x}_{k+1}, r_k \mid \mathbf{x}_k, \mathbf{u}_k) d\mathbf{x}_{k+1} . \quad (4.23)$$

The function p_x defines the dynamics of the MDP

$$\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) \quad (4.24)$$

and the function p_r gives the rewards of the MDP

$$r_k \sim p_r(r_k \mid \mathbf{x}_k, \mathbf{u}_k) . \quad (4.25)$$

We can once again determine the *expected rewards* for state-input pairs as a two-argument function $r_k : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$,

$$r_k(\mathbf{x}_k, \mathbf{u}_k) = \mathbb{E}_{r_k \sim p_r(r_k \mid \mathbf{x}_k, \mathbf{u}_k)}\{r_k\} = \int_{\mathcal{R}} r_k p_r(r_k \mid \mathbf{x}_k, \mathbf{u}_k) dr_k . \quad (4.26)$$

4.1.5 Control policy

A control law or a control policy⁴ is a function that tells us which is the best input to choose in each state. A policy can be deterministic or stochastic.

Definition 4.3. (Policy, see [4.0, p. 15]). A *stochastic policy* $\pi : \mathcal{X} \times \mathcal{U} \rightarrow [0, 1]$ is a distribution over inputs given states

$$\mathbf{u}_k \sim \pi(\mathbf{u}_k \mid \mathbf{x}_k) . \quad (4.27)$$

A *deterministic policy* $\pi : \mathcal{X} \rightarrow \mathcal{U}$ is a sequence of functions

$$(\boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{N-1}) , \quad (4.28)$$

where $\boldsymbol{\mu}_k$ maps the state vector \mathbf{x}_k to control inputs

$$\mathbf{u}_k = \boldsymbol{\mu}_k(\mathbf{x}_k) . \quad (4.29)$$

The policy is admissible if $\boldsymbol{\mu}_k(\mathbf{x}_k) \in \mathcal{U}(\mathbf{x}_k)$ holds for all $\mathbf{x}_k \in \mathcal{X}_k$ and k . The deterministic policy (4.29) can be written as

$$\pi(\mathbf{u}_k \mid \mathbf{x}_k) = \begin{cases} 1 & \text{if } \mathbf{u}_k = \boldsymbol{\mu}_k(\mathbf{x}_k) \\ 0 & \text{else} \end{cases} . \quad (4.30)$$

Example 4.5. In the discrete case, a stochastic policy denoted as $\pi(\mathbf{u} \mid \mathbf{x})$ is a conditional distribution over the inputs $\mathbf{u} \in \mathcal{U}$ given the state $\mathbf{x} \in \mathcal{X}$, $\pi(\mathbf{u} \mid \mathbf{x}) = P(\mathbf{u} \mid \mathbf{x})$. As an example, if a robot has four inputs $\mathcal{U} = \{\mathbf{u}(\leftarrow), \mathbf{u}(\downarrow), \mathbf{u}(\rightarrow), \mathbf{u}(\uparrow)\}$. The policy at a state $x \in \mathcal{X}$ for such a set of inputs \mathcal{U} is a distribution where the probabilities of the four inputs could be $[0.4, 0.2, 0.1, 0.3]$; at some other state $x' \in \mathcal{X}$ the probabilities $\pi(u \mid x')$ of the same four inputs could be $[0.1, 0.1, 0.2, 0.6]$. Note that we should have $\sum_{u \in \mathcal{U}} \pi(u \mid x) = 1$ for any state x . A deterministic policy is a

⁴It is also often referred to simply as policy.

special case of a stochastic policy in that the distribution $\pi(u | x)$ only gives non-zero probability to one particular input, e.g., $[1, 0, 0, 0]$ for our example with four inputs.

Theorem 4.1. (*Probability of observing a rollout, see [4.0, p. 11]*) The probability distribution of observing a N -step rollout τ under a policy π is

$$p(\tau | \pi) = p^\pi(\tau) = p_{x_0}(x_0) \prod_{k=0}^{N-1} p_x(x_{k+1} | x_k, u_k) \pi(u_k | x_k) . \quad (4.31)$$

Proof. We can proof (4.31) by induction using the product rule (A.11) and the Markov property (4.4). With $p^\pi(x_0) \doteq p_{x_0}(x_0)$ and $p^\pi(u_k | x_k) \doteq \pi(u_k | x_k)$, we get

$$\begin{aligned} p^\pi(x_1, x_0, u_0) &= p_{x_0}(x_0) \underbrace{p(x_1, u_0 | x_0)}_{=p_x(x_1 | x_0, u_0) \pi(u_0 | x_0)} \\ p^\pi(x_2, x_1, x_0, u_1, u_0) &= p^\pi(x_1, x_0, u_0) \underbrace{p(x_2, u_1 | x_1, x_0, u_0)}_{=p_x(x_2 | x_1, u_1) \pi(u_1 | x_1)} \\ &= p_{x_0}(x_0) \prod_{k=0}^1 p_x(x_{k+1} | x_k, u_k) \pi(u_k | x_k) \quad (4.32) \\ &\vdots \\ p^\pi(x_N, \dots, x_0, u_{N-1}, \dots, u_0) &= p_{x_0}(x_0) \prod_{k=0}^{N-1} p_x(x_{k+1} | x_k, u_k) \pi(u_k | x_k) . \end{aligned}$$

□

Whenever any policy π , whether deterministic or probabilistic, is implemented, the resulting process is a Markov process. In value-discrete case, the associated *state-transition probability tensor* is denoted by \mathbf{P}^π . In the deterministic case, $\pi \in \Pi_d$, then at time index k , the corresponding input u_k equals $\mu_k(x_k)$ and we have

$$\mathbf{P}^\pi[i, j] = p_{ij}^\pi = \mathbb{P}[x_{k+1} = x_j | x_k = x_i, u_k = \mu_l] . \quad (4.33)$$

In the stochastic case, $\pi \in \Pi_p$ and

$$p_{ij}^\pi = \mathbb{E}_{\mathbf{u} \sim \pi(\mathbf{u} | \mathbf{x})} \left\{ p_{ij}^{\mathbf{u}} \right\} . \quad (4.34)$$

4.1.6 Value functions

In reinforcement learning, value functions play a central role in representing and estimating the expected future rewards or returns an agent can obtain from states and inputs. They serve as a foundational concept for many algorithms and provide insight into the quality of different decisions.

Definition 4.4. (Action-Value Function Q^π for a MDP \mathcal{M} and policy π , see [4.0]). The *action-value function*^a $Q^\pi : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ for a MDP \mathcal{M} and policy π gives the expected return if you start in state \mathbf{x} , take an arbitrary control input \mathbf{u} (which may not have to come from the control law), and then forever after act according to policy π , i.e.

$$Q_k^\pi(\mathbf{x}, \mathbf{u}) = \mathbb{E}_\pi \left\{ R_k(\boldsymbol{\tau}_{[k:N]}) \mid \mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u} \right\}. \quad (4.35)$$

^aAlso Q-value function.

The action-value function depends on $\mathbf{x}, \mathbf{u}, \pi$ and $p_x(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})$ but is independent of $\mathbf{X}_{[k+1:N]}$ and $\mathbf{U}_{[k+1:N-1]}$ since they are eliminated by taking the expectation. The action-value function $Q_k^\pi(\mathbf{x}, \mathbf{u})$ evaluates how good it is to pick an input \mathbf{u} being in state \mathbf{x} .

Definition 4.5. (State-Value Function V^π , see [4.0]). The *state-value function* $V^\pi : \mathcal{X} \rightarrow \mathbb{R}$ gives the expected return if you start in state \mathbf{x} and always act according to policy π , i.e.

$$V_k^\pi(\mathbf{x}) = \mathbb{E}_\pi \left\{ R_k(\boldsymbol{\tau}_{[k:N]}) \mid \mathbf{x}_k = \mathbf{x} \right\} \quad (4.36a)$$

$$= \mathbb{E}_{\mathbf{u}_k \sim \pi(\mathbf{u}_k | \mathbf{x}_k)} \{ Q_k^\pi(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_k = \mathbf{x} \}. \quad (4.36b)$$

Intuitively, for a fixed policy π , the state-value function $V_k^\pi(\mathbf{x})$ evaluates how good the situation is in state \mathbf{x} . Clearly, in the deterministic case, the action-value function and state-value function are related via

$$V_k^\pi(\mathbf{x}_k) = Q_k^\pi(\mathbf{x}_k, \boldsymbol{\mu}_k(\mathbf{x}_k)). \quad (4.37)$$

Moreover, by definition, the optimal value function produces the maximum return

$$V_k^*(\mathbf{x}) = \max_{\pi \in \Pi} V_k^\pi(\mathbf{x}) \quad \text{and} \quad Q_k^*(\mathbf{x}, \mathbf{u}) = \max_{\pi \in \Pi} Q_k^\pi(\mathbf{x}, \mathbf{u}). \quad (4.38)$$

Definition 4.6. (Advantage Function A^π for a MDP \mathcal{M} , see [4.0]). The *advantage function* $A^\pi : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ of a MDP \mathcal{M} for a given policy π , indicates how much better the control input \mathbf{u} is in state \mathbf{x} compared to the average, i.e.

$$A_k^\pi(\mathbf{x}, \mathbf{u}) = Q_k^\pi(\mathbf{x}, \mathbf{u}) - V_k^\pi(\mathbf{x}). \quad (4.39)$$

Advantage functions measure the relative benefit of choosing a specific input over others in a given state.

Note 4.4. The action-value, state-value, and advantage functions are functions of the time index k in the finite horizon setting. To simplify the notation, we omit the time index k whenever the functions are conditioned on $\mathbf{x}_k = \mathbf{x}$ and/or $\mathbf{u}_k = \mathbf{u}$.

Example 4.6 (Frozen Lake game continued). In the infinite horizon and value discrete case, the value function can be organized in a table. For instance, for a 4×4 Frozen Lake grid, there are 16 states (each cell is a state) and 4 inputs $\{\leftarrow, \downarrow, \rightarrow, \uparrow\}$. The Q -table is a 16×4 matrix with rows representing states and columns representing inputs. So want to calculate $16 \times 4 = 64$ action-state values. Table 4.2 shows the Q -Table for the example at hand.

State	$u = u(\leftarrow) = 0$	$u = u(\downarrow) = 1$	$u = u(\rightarrow) = 2$	$u = u(\uparrow) = 3$
$x = 0$	$Q^\pi(0, 0)$	$Q^\pi(0, 1)$	$Q^\pi(0, 2)$	$Q^\pi(0, 3)$
$x = 1$	$Q^\pi(1, 0)$	$Q^\pi(1, 1)$	$Q^\pi(1, 2)$	$Q^\pi(1, 3)$
\dots	\dots	\dots	\dots	\dots
$x = 14$	$Q^\pi(14, 0)$	$Q^\pi(14, 1)$	$Q^\pi(14, 2)$	$Q^\pi(14, 3)$
$x = 15$	$Q^\pi(15, 0)$	$Q^\pi(15, 1)$	$Q^\pi(15, 2)$	$Q^\pi(15, 3)$

Table 4.2: Q-Table for the 4×4 Frozen Lake game.

4.1.7 Bellman Expectation Equation

The *Bellman Expectation Equation* serves as the foundational framework for reinforcement learning. It is derived next for the continuous case in terms of the state-value function and action-value function.

State-value function

The state-value function can be decomposed into an expected immediate reward plus the discounted values of the successor state. To show this, we start from the definition of the state-value function (4.36a) and substitute $R_k = r_k + \gamma R_{k+1}$ to get

$$\begin{aligned} V_k^\pi(\mathbf{x}_k) &= \mathbb{E}_\pi\{R_k \mid \mathbf{x}_k = \mathbf{x}_k\} \\ &= \mathbb{E}_\pi\{r_k \mid \mathbf{x}_k = \mathbf{x}_k\} + \gamma \mathbb{E}_\pi\{R_{k+1} \mid \mathbf{x}_k = \mathbf{x}_k\}. \end{aligned} \quad (4.40)$$

We already introduced the *expected reward* in (4.17) as

$$r_k(\mathbf{x}_k, \mathbf{u}_k) = \mathbb{E}_\pi\{r_k \mid \mathbf{x}_k = \mathbf{x}_k, \mathbf{u}_k = \mathbf{u}_k\} = \mathbb{E}_{r_k \sim p_r(r_k | \mathbf{x}_k, \mathbf{u}_k)}\{r_k\} = \int_{\mathcal{R}} r_k p_r(r_k \mid \mathbf{x}_k, \mathbf{u}_k) dr_k. \quad (4.41)$$

Consequently, the *expected reward* following π is

$$r_k^\pi(\mathbf{x}_k) = \mathbb{E}_\pi\{r_k \mid \mathbf{x}_k = \mathbf{x}_k\} = \mathbb{E}_{\mathbf{u}_k \sim \pi(\mathbf{u}_k | \mathbf{x}_k)}\{r_k(\mathbf{x}_k, \mathbf{u}_k)\} = \int_{\mathcal{R}} r_k p_r(r_k \mid \mathbf{x}_k) dr_k. \quad (4.42)$$

Note that the distribution $p_r(r_k | \mathbf{x}_k)$ is a marginal distribution of a distribution that also contains the variables \mathbf{u}_k and \mathbf{x}_{k+1} , respectively. Thus, we have

$$\begin{aligned} p_r(r_k | \mathbf{x}_k) &= \int_{\mathcal{X}} \int_{\mathcal{U}} p(\mathbf{x}_{k+1}, \mathbf{u}_k, r_k | \mathbf{x}_k) d\mathbf{u}_k d\mathbf{x}_{k+1} \\ &= \int_{\mathcal{X}} \int_{\mathcal{U}} p(\mathbf{x}_{k+1}, r_k | \mathbf{u}_k, \mathbf{x}_k) \pi(\mathbf{u}_k | \mathbf{x}_k) d\mathbf{u}_k d\mathbf{x}_{k+1}, \end{aligned} \quad (4.43)$$

where we used $\pi(\mathbf{u}_k | \mathbf{x}_k) \doteq p(\mathbf{u}_k | \mathbf{x}_k)$. For the second part follows after some reformulations

$$\begin{aligned} \mathbb{E}_\pi\{\mathbf{R}_{k+1} | \mathbf{x}_k = \mathbf{x}_k\} &= \int_{\mathcal{R}} R_{k+1} p(R_{k+1} | \mathbf{x}_k) dR_{k+1} \\ &\stackrel{(A.12)}{=} \int_{\mathcal{X}} \int_{\mathcal{R}} R_{k+1} p(R_{k+1} | \mathbf{x}_{k+1}, \mathbf{x}_k) p(\mathbf{x}_{k+1} | \mathbf{x}_k) dR_{k+1} d\mathbf{x}_{k+1} \\ &= \int_{\mathcal{X}} \underbrace{p(\mathbf{x}_{k+1} | \mathbf{x}_k) \int_{\mathcal{R}} R_{k+1} p(R_{k+1} | \mathbf{x}_{k+1}, \mathbf{x}_k) dR_{k+1}}_{= \mathbb{E}_\pi\{\mathbf{R}_{k+1} | \mathbf{x}_{k+1} = \mathbf{x}_{k+1}\} = V_{k+1}^\pi(\mathbf{x}_{k+1})} d\mathbf{x}_{k+1} \\ &\stackrel{(A.12)}{=} \int_{\mathcal{U}} \int_{\mathcal{X}} p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k) \pi(\mathbf{u}_k | \mathbf{x}_k) V_{k+1}^\pi(\mathbf{x}_{k+1}) d\mathbf{x}_{k+1} d\mathbf{u}_k \\ &= \mathbb{E}_{\substack{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k) \\ \mathbf{u}_k \sim \pi(\mathbf{u}_k | \mathbf{x}_k)}} \{V_{k+1}^\pi(\mathbf{x}_{k+1})\} \end{aligned} \quad (4.44)$$

Finally, we have found the *Bellman Expectation Equation*

$$V_k^\pi(\mathbf{x}_k) = r_k^\pi(\mathbf{x}_k) + \gamma \mathbb{E}_{\substack{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k) \\ \mathbf{u}_k \sim \pi(\mathbf{u}_k | \mathbf{x}_k)}} \{V_{k+1}^\pi(\mathbf{x}_{k+1})\}. \quad (4.45)$$

Furthermore, in the *deterministic case*, the Bellman Expectation Equation (4.45) simplifies to the Bellman Equation

$$V_k^\pi(\mathbf{x}_k) = r_k^\pi(\mathbf{x}_k) + \gamma V_{k+1}^\pi(\mathbf{x}_{k+1}). \quad (4.46)$$

Within this framework, Figure 4.5 illustrates the significance of the Bellman Equation. Here,

- $V_k^\pi(\mathbf{x}_k)$ can be viewed as a predicted performance,
- $r_k^\pi(\mathbf{x}_k)$ represents the observed immediate reward, and
- $V_{k+1}^\pi(\mathbf{x}_{k+1})$ offers a current prediction of future control inputs.

This concept is called *bootstrapping* in reinforcement learning. It refers to the idea of using the estimates of future value functions to update the current estimate of the value function. In other words, we're "bootstrapping" our value updates based on other estimated values – "Learn a guess from a guess". These concepts are further explored and utilized in the ensuing discussion on temporal difference learning.

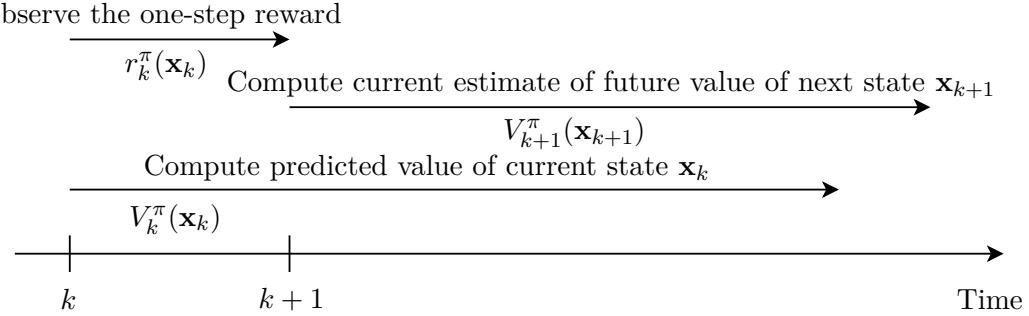


Figure 4.5: The temporal difference perspective of the Bellman Equation illustrates how the equation embodies the processes of action-taking, observation, assessment, and enhancement inherent in reinforcement learning [4.0, p. 469].

Action-value function

We can draw similar conclusions for the action-value function. From (4.35), we get

$$\begin{aligned} Q_k^\pi(\mathbf{x}_k, \mathbf{u}_k) &= \mathbb{E}_\pi\{\mathbf{R}_k \mid \mathbf{x}_k = \mathbf{x}_k, \mathbf{u}_k = \mathbf{u}_k\} \\ &= \mathbb{E}_\pi[r_k \mid \mathbf{x}_k = \mathbf{x}_k, \mathbf{u}_k = \mathbf{u}_k] + \gamma \mathbb{E}_\pi\{\mathbf{R}_{k+1} \mid \mathbf{x}_k = \mathbf{x}_k, \mathbf{u}_k = \mathbf{u}_k\} \\ &= r_k(\mathbf{x}_k, \mathbf{u}_k) \\ &\quad + \gamma \int_{\mathcal{X}} p(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) \int_{\mathcal{U}} \pi(\mathbf{u}_{k+1} \mid \mathbf{x}_{k+1}) Q_{k+1}^\pi(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) d\mathbf{u}_{k+1} d\mathbf{x}_{k+1}. \end{aligned} \quad (4.47)$$

Hence, in equivalence to (4.45), we find the recursion

$$Q_k^\pi(\mathbf{x}_k, \mathbf{u}_k) = r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\substack{\mathbf{x}_{k+1} \sim p_{\mathbf{x}}(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) \\ \mathbf{u}_{k+1} \sim \pi(\mathbf{u}_{k+1} \mid \mathbf{x}_{k+1})}} \{Q_{k+1}^\pi(\mathbf{x}_{k+1}, \mathbf{u}_{k+1})\} \quad (4.48a)$$

$$= r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_{\mathbf{x}}(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k)} \{V_{k+1}^\pi(\mathbf{x}_{k+1})\}. \quad (4.48b)$$

The last reformulation is due to the relation (4.37).

4.1.8 Bellman Optimality Equation

The previously obtained results motivate the next theorem:

Theorem 4.2. (Bellman Optimality Equation, see [4.0, p. 12]) Let us define the optimal action-state value function $Q_k^* : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ by

$$Q_k^*(\mathbf{x}_k, \mathbf{u}_k) = r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_{\mathbf{x}}(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k)} \{V_{k+1}^*(\mathbf{x}_{k+1})\}. \quad (4.49)$$

Then, Q_k^* satisfies the following Bellman Optimality Equation

$$Q_k^*(\mathbf{x}_k, \mathbf{u}_k) = r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_{\mathbf{x}}(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)} \left\{ \max_{\mathbf{u}_{k+1} \in \mathcal{U}} Q_{k+1}^*(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \right\} \quad (4.50)$$

with

$$V_k^*(\mathbf{x}_k) = \max_{\mathbf{u}_k \in \mathcal{U}} Q_k^*(\mathbf{x}_k, \mathbf{u}_k) . \quad (4.51)$$

Moreover, every policy $\pi \in \Pi_d$ such that

$$\mu_k^*(\mathbf{x}_k) = \arg \max_{\mathbf{u}_k \in \mathcal{U}} Q_k^*(\mathbf{x}_k, \mathbf{u}_k) \quad (4.52)$$

is optimal.

Proof. Since Q^* is defined by (4.49), it follows that

$$\max_{\mathbf{u}_k \in \mathcal{U}} Q_k^*(\mathbf{x}_k, \mathbf{u}_k) = \max_{\mathbf{u}_k \in \mathcal{U}} \left[r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_{\mathbf{x}}(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)} \{ V_{k+1}^*(\mathbf{x}_{k+1}) \} \right] = V_k^*(\mathbf{x}_k) . \quad (4.53)$$

This establishes (4.51) and (4.52). Substituting from (4.51) into (4.49) gives (4.50). \square

Obviously, the Bellman Optimality Equation can be used to find the optimal value function and policy.

4.2 Model-based reinforcement learning

In model-based reinforcement learning, we assume to know the dynamic and reward function p_x and p_r , respectively.

4.2.1 Infinite horizon Dynamic Programming

Next we show the *Value and Policy Iteration* used in reinforcement learning. In this section, we will restrict our self to the infinite horizon ($N = \infty$) and value-discrete case. In this case, the Bellman Expectation equation (4.48) with (4.18) can be written as

$$\begin{aligned} Q^\pi(\mathbf{x}_i, \mathbf{u}_l) &= r(\mathbf{x}_i, \mathbf{u}_l) + \gamma \sum_{\mathbf{x}_j \in \mathcal{X}} p_x(\mathbf{x}_j \mid \mathbf{x}_i, \mathbf{u}_l) V^\pi(\mathbf{x}_j) \\ &= r(\mathbf{x}_i, \mathbf{u}_l) + \gamma \sum_{j=0}^{|\mathcal{X}|-1} p_{ij}^{\mathbf{u}_l} V^\pi(\mathbf{x}_j). \end{aligned} \quad (4.54)$$

Value Iteration

Value iteration is the computation of the state-value function V^π by the Dynamic Programming described in Theorem 4.2. The iteration is performed for the entire state space starting from arbitrarily initialized residual reward [4.0, p. 83]. The iteration rule in Theorem 4.2 represents a fixed point iteration for V^π and converges against V^* , see [4.0]. Therefore, we introduce the *state-value vector*

$$\mathbf{v}^\pi = \left[\begin{array}{ccc} V^\pi(\mathbf{x}_0) & \dots & V^\pi(\mathbf{x}_{|\mathcal{X}|-1}) \end{array} \right]^\top \in \mathbb{R}^{|\mathcal{X}|} \quad (4.55)$$

and define the *Bellman Iteration Map*⁵ $\mathbf{B} : \mathbb{R}^{|\mathcal{X}|} \rightarrow \mathbb{R}^{|\mathcal{X}|}$ via

$$\mathbf{y} \mapsto (\mathbf{B}(\mathbf{y}))[i] = \max_{\mathbf{u}_l \in \mathcal{U}} \left[r(\mathbf{x}_i, \mathbf{u}_l) + \gamma \sum_{j=0}^{|\mathcal{X}|-1} p_{ij}^{\mathbf{u}_l} \mathbf{y}[j] \right]. \quad (4.56)$$

Theorem 4.3. (Theorem 3 in [4.0]) *The map \mathbf{B} is monotone and a contraction with respect to the ℓ_∞ -norm. Therefore, the fixed point $\mathbf{v}^{(\infty)}$ of the map \mathbf{B} satisfies the relation*

$$\mathbf{v}^{(\infty)}[i] = \max_{\mathbf{u}_l \in \mathcal{U}} \left[r(\mathbf{x}_i, \mathbf{u}_l) + \gamma \sum_{j=0}^{|\mathcal{X}|-1} p_{ij}^{\mathbf{u}_l} \mathbf{v}^{(\infty)}[j] \right]. \quad (4.57)$$

See [4.0] for a proof. Given an initial guess $\mathbf{v}^{(0)} \in \mathbb{R}^{|\mathcal{X}|}$, one can iteratively employ the Bellman iteration. These iterations will converge to the unique fixed point, denoted as $\mathbf{v}^{(\infty)}$, of the operator \mathbf{B} . The importance of this iterative process is elaborated in the subsequent theorem.

⁵Called Backup operator in the reinforcement learning literature.

Theorem 4.4. (*Theorem 4 in [4.0]*) Define $\mathbf{v}^{(\infty)} \in \mathbb{R}^n$ to be the unique fixed point of B , and define $\mathbf{v}^* \in \mathbb{R}^n$ to equal $V^*(\mathbf{x})$, $\mathbf{x} \in \mathcal{X}$, where $V^*(\mathbf{x})$ is defined in (4.38). Then $\mathbf{v}^{(\infty)} = \mathbf{v}^*$.

See [4.0] for a proof. Therefore, the optimal value vector can be computed using the Bellman iteration. However, knowing the optimal value vector does not, by itself, give us an optimal policy. Therefore, after the convergence Bellman iteration, a policy determination according to (4.52) is performed.

The *Infinite Horizon Value Iteration Algorithm 1* show how to use Value Iteration to find an optimal policy π^* .

Algorithm 1: Infinite Horizon Value Iteration Algorithm

```

Data:  $\theta$  is a small number
Result: Find  $V^*(\mathbf{x})$  and  $\pi^*$ ,  $\forall \mathbf{x} \in \mathcal{X}$ 
/* Initialization */ 
1 Initialize  $V(\mathbf{x})$  arbitrarily;
2  $V^\pi(\mathbf{x}) \leftarrow 0$ ,  $\forall \mathbf{x} \in \mathcal{X}$ ;
   /* Loop until convergence */
3  $\Delta \leftarrow 0$ ;
   /* Optimal Value Determination */
4 while  $\Delta < \theta$  do
5   for each  $\mathbf{x} \in \mathcal{X}$  do
6      $v \leftarrow V^\pi(\mathbf{x})$ ;
7      $V^\pi(x) \leftarrow \max_{\mathbf{u} \in \mathcal{U}} [r(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_x(\mathbf{x}' | \mathbf{x}, \mathbf{u})} \{V^\pi(\mathbf{x}')\}]$ ;
8      $\Delta \leftarrow \max(\Delta, |v - V^\pi(\mathbf{x})|)$ ;
9   end
10 end
   /* Optimal Policy Determination */
11  $\mu^*(\mathbf{x}) \leftarrow \arg \max_{\mathbf{u} \in \mathcal{U}} [r(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}_{\mathbf{x}' \sim p_x(\mathbf{x}' | \mathbf{x}, \mathbf{u})} \{V^*(\mathbf{x}')\}]$ ;
12 return  $\pi^*$ ,  $\forall \mathbf{x} \in \mathcal{X}$ 
```

Policy Iteration

When the model is fully known, following Bellman Expectation equation, we can use Dynamic Programming (DP) to iteratively evaluate value functions and improve the policy.

- *Policy Evaluation* (value update) is to compute the state-value functions $V^\pi(\mathbf{x})$ for a given policy π according to the Bellman Expectation equation (4.45).
- Based on the state-value functions, *Policy Improvement* (policy update) generates a better policy $V^{\pi'}(\mathbf{x}) \geq V^\pi(\mathbf{x})$ by acting greedily.

Once again, we show Policy Iteration for the infinite horizon ($N = \infty$) and value-discrete case (\mathcal{X} and \mathcal{U} finite) only.

Policy Evaluation

Define the *state-value vector* $\mathbf{v}^\pi \in \mathbb{R}^{|\mathcal{X}|}$ via

$$\mathbf{v}^\pi = \begin{bmatrix} V^\pi(\mathbf{x}_0) & \dots & V^\pi(\mathbf{x}_{|\mathcal{X}|-1}) \end{bmatrix}^\top, \quad (4.58)$$

and the *reward vector* $\mathbf{r}^\pi \in \mathbb{R}^{|\mathcal{X}|}$ via

$$\mathbf{r}^\pi = \begin{bmatrix} r_0^\pi & \dots & r_{|\mathcal{X}|-1}^\pi \end{bmatrix}^\top. \quad (4.59)$$

Then, the value-discrete infinite horizon Bellman Expectation equation (4.45) reads in vector notation as

$$\mathbf{v}^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}^\pi. \quad (4.60)$$

Policy Evaluation involves computing the state-value vector \mathbf{v}^π for a given policy π solving the Bellman Expectation equation. We also refer to it as the prediction problem. In this context, it is important to note that the induced matrix norm $\|\mathbf{P}^\pi\|_\infty$ is equal to 1, and since $\gamma < 1$, the matrix $\mathbf{I} - \gamma \mathbf{P}^\pi$ is nonsingular. Consequently, for any reward vector \mathbf{r}^π , there exists a unique \mathbf{v}^π that satisfies (4.45). In principle it is possible to deduce from (4.45)

$$\mathbf{v}^\pi = (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}^\pi. \quad (4.61)$$

However, applying this formula can be challenging in practical applications of Markov Decision Problems, mainly due to the large size of the state space \mathcal{X} , represented by the integer $|\mathcal{X}|$. Additionally, matrix inversion has a cubic complexity, making it impractical to invert the matrix $\mathbf{I} - \gamma \mathbf{P}^\pi$. Consequently, alternative approaches need to be explored. The Contraction Mapping Theorem provides a feasible solution in such cases.

Theorem 4.5. (Iterative policy evaluation, Theorem 2 in [4.0]) The map $\mathbf{y} \mapsto \mathbf{T}(\mathbf{y}) = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{y}$ is monotone and is a contraction with respect to the ℓ_∞ -norm, with contraction constant γ . Therefore, we can choose some vector $\mathbf{y}^{(0)}$ arbitrarily, and then define

$$\mathbf{y}^{(i+1)} = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{y}^{(i)}. \quad (4.62)$$

Then $\mathbf{y}^{(i)}$ converges to the state value vector \mathbf{v}^π .

See [4.0] for a proof. Similar results can be obtained for the value-continuous case, see [4.0]. Therefore, the transition tensor \mathbf{P}^π is replaced by a transition operator.

Policy Improvement

Given a state-value function $V^\pi(\mathbf{x})$ of a policy π , it is possible to perform a *policy improvement step*, cf. [4.0, pp. 84–85]. The result is a potentially better strategy π' , where $V^{\pi'}(\mathbf{x}) \geq V^\pi(\mathbf{x})$ holds.

Theorem 4.6. (*Policy Improvement Theorem*) We consider two policies $\pi(\mathbf{u} \mid \mathbf{x})$ and $\pi'(\mathbf{u} \mid \mathbf{x})$. Let us define

$$Q^\pi(\mathbf{x}, \pi') = \mathbb{E}_{\mathbf{u} \sim \pi'(\mathbf{u} \mid \mathbf{x})} \{Q^\pi(\mathbf{x}, \mathbf{u})\} . \quad (4.63)$$

If $\forall \mathbf{x} \in \mathcal{X}$, we have that $Q^\pi(\mathbf{x}, \pi') \geq V^\pi(\mathbf{x})$, then it holds that

$$V^\pi(\mathbf{x}) \leq Q^\pi(\mathbf{x}, \pi') \leq V^{\pi'}(\mathbf{x}), \quad \forall \mathbf{x} . \quad (4.64)$$

This means that π' is at least as good a policy as π .

Proof. By expanding Q^π , we can get that $\forall \mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} V^\pi(\mathbf{x}) &\leq Q^\pi(\mathbf{x}, \pi') = \mathbb{E}_{\mathbf{u} \sim \pi'(\mathbf{u} \mid \mathbf{x})} \{Q^\pi(\mathbf{x}, \mathbf{u})\} \\ &= \mathbb{E}_{\mathbf{u} \sim \pi'(\mathbf{u} \mid \mathbf{x})} \left\{ r_k + \gamma \underbrace{V^\pi(\mathbf{x}')}_{\leq Q^\pi(\mathbf{x}', \pi')} \right\} \\ &\leq \mathbb{E}_{\mathbf{u} \sim \pi'(\mathbf{u} \mid \mathbf{x})} \{r_k + \gamma Q^\pi(\mathbf{x}', \pi')\} \\ &= \mathbb{E}_{\mathbf{u}, \mathbf{u}' \sim \pi'} \{r_k + \gamma r_{k+1} + \gamma^2 V^\pi(\mathbf{x}'')\} \\ &\leq \dots \\ &\leq \mathbb{E}_{\mathbf{u}, \mathbf{u}', \mathbf{u}'' \dots \sim \pi'} \{r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots\} \\ &= V^{\pi'}(\mathbf{x}) . \end{aligned} \quad (4.65)$$

□

To obtain π' , from now on, the previous strategy π is no longer followed, but instead the *greedy policy* is chosen, which maximizes the expected reward:

$$\begin{aligned} \mu'(\mathbf{x}) &= \mathbf{u}' = \arg \max_{\mathbf{u} \in \mathcal{U}} Q^\pi(\mathbf{x}, \mathbf{u}) \\ &= \arg \max_{\mathbf{u} \in \mathcal{U}} [r^\pi + \gamma \mathbb{E}_{\mathbf{x}' \sim p_x(\mathbf{x}' \mid \mathbf{x}, \mathbf{u})} \{V^\pi(\mathbf{x}')\}] . \end{aligned} \quad (4.66)$$

The *Infinite Horizon Policy Iteration Algorithm 2* shows how to use policy evaluation and policy improvement to find an optimal policy π^* .

Algorithm 2: Infinite Horizon Policy Iteration Algorithm

```

Data:  $\theta$  is a small number
/* Initialization */
```

1 Initialize $V^\pi(\mathbf{x})$ arbitrarily;

2 Randomly initialize policy $\mu(\mathbf{x})$;

/* Policy Evaluation */

3 $\Delta \leftarrow 0$;

4 **while** $\Delta < \theta$ **do**

5 **for** each $\mathbf{x} \in \mathcal{X}$ **do**

6 $v^\pi \leftarrow V^\pi(\mathbf{x})$;

7 $V^\pi(\mathbf{x}) \leftarrow r^\pi + \gamma \mathbb{E}_{\mathbf{x}' \sim p_x(\mathbf{x}'|\mathbf{x}, \mathbf{u})} \{V^\pi(\mathbf{x}')\}$;

8 $\Delta \leftarrow \max(\Delta, |v^\pi - V^\pi(\mathbf{x})|)$;

9 **end**

10 **end**

/* Policy Improvement */

11 policy is stable \leftarrow true;

12 **for** each $\mathbf{x} \in \mathcal{X}$ **do**

13 old-policy $\leftarrow \pi$;

14 $\mu(\mathbf{x}) \leftarrow \arg \max_{\mathbf{u} \in \mathcal{U}} [r^\pi + \gamma \mathbb{E}_{\mathbf{x}' \sim p_x(\mathbf{x}'|\mathbf{x}, \mathbf{u})} \{V^\pi(\mathbf{x}')\}]$;

15 **if** old-policy $\neq \mu(\mathbf{x})$ **then**

16 policy is stable \leftarrow false;

17 **end**

18 **end**

19 **if** policy is stable **then**

20 return $V^\pi \approx V^*$ and $\pi \approx \pi^*$;

21 **else**

22 go to Policy Evaluation;

23 **end**

Finally, we compare Value Iteration and Policy Iteration and draw some conclusions:

Aspect	Value Iteration	Policy Iteration
Convergence Speed	Faster to optimal value function	Fewer iterations required and converges to the optimal policy faster
Implementation	Simpler, no separate policy needed	More complex with policy evaluation and update
Efficiency per Iteration	Slower due to all states sweep	More efficient with early policy stop
Computational Expense	Varied, generally moderate	Computationally expensive per iteration

Table 4.3: Comparison of Value Iteration and Policy Iteration.

4.2.2 Finite horizon Dynamic Programming

Finite horizon Dynamic Programming (DP) is based on Bellman's optimality principle, which states that every optimal solution of (4.38) is composed of optimal partial solutions. The basic idea of DP is to divide the problem into subproblems, which can be solved more easily, and to assemble these partial solutions into the overall solution. We will show the Dynamic Programming Algorithm for MDPs only.

Theorem 4.7. (*Bellman's Optimality Principle, see [4.0, p. 20]*).

Let $\pi^* = (\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*)$ be the optimal policy for

$$V_0^\pi(\mathbf{x}_0) = \mathbb{E}_\pi \left\{ R_0(\tau_{[0:N]}) \mid \mathbf{x}_0 = \mathbf{x}_0 \right\}. \quad (4.67)$$

Considering the subproblem of (4.67) where, starting from the state \mathbf{x}_{k+1} , the expected returns

$$V_{k+1}^\pi(\mathbf{x}_{k+1}) = \mathbb{E}_\pi \left\{ R_{k+1}(\tau_{[k+1:N]}) \mid \mathbf{x}_{k+1} = \mathbf{x}_{k+1} \right\}. \quad (4.68)$$

are to be maximized, then the truncated policy $(\mu_{k+1}^*, \mu_{k+2}^*, \dots, \mu_{N-1}^*)$ is optimal for the subproblem.

The underlying idea behind the Principle of Optimality is straightforward. If the truncated policy $(\mu_{k+1}^*, \mu_{k+2}^*, \dots, \mu_{N-1}^*)$ were not optimal as stated, we would be able to increase the return further by switching to an optimal policy for the subproblem once we reach \mathbf{x}_{k+1} . Based on Theorem 4.2, the *Dynamic Programming Algorithm* can be stated. The algorithm constructs optimal value functions

$$V_N^*(\mathbf{x}_N), V_{N-1}^*(\mathbf{x}_{N-1}), \dots, V_0^*(\mathbf{x}_0) \quad (4.69)$$

starting from $V_N^*(\mathbf{x}_N)$ and proceeding backwards in time.

Proposition 4.1. (Dynamic Programming Algorithm, see [4.0, p. 23] - backward pass). For every initial state \mathbf{x}_0 , the maximum expected return $V_0^*(\mathbf{x}_0)$ of the basic problem is equal to $V_0^\pi(\mathbf{x}_0)$, given by the last step of the following algorithm, which proceeds from $\gamma^N r_N$ backward in time from $k = N - 1$ to $k = 0$, for all $\mathbf{x}_k \in \mathcal{X}$:

$$\begin{aligned} V_N^*(\mathbf{x}_N) &= \gamma^N r_N(\mathbf{x}_N) \\ V_k^*(\mathbf{x}_k) &= \max_{\mathbf{u}_k \in \mathcal{U}} Q_k^*(\mathbf{x}_k, \mathbf{u}_k), \end{aligned} \quad (4.70)$$

with

$$Q_k^*(\mathbf{x}_k, \mathbf{u}_k) = r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)} \{ V_{k+1}^*(\mathbf{x}_{k+1}) \}. \quad (4.71)$$

Furthermore, if $\mathbf{u}_k^* = \mu_k^*(\mathbf{x}_k)$ satisfies the Bellman iteration (4.70) for each \mathbf{x}_k and k , the policy $\pi^* = (\mu_0^*, \mu_1^*, \dots, \mu_{N-1}^*)$ is optimal.

Once the value functions $V_0^*(\mathbf{x}_0), \dots, V_N^*(\mathbf{x}_N)$ have been obtained, we can use the following forward algorithm to construct an optimal control input sequence $(\mathbf{u}_0^*, \mathbf{u}_1^*, \dots, \mathbf{u}_{N-1}^*)$ and the corresponding state trajectory $(\mathbf{x}_0^*, \mathbf{x}_1^*, \dots, \mathbf{x}_{N-1}^*)$ for a given initial state \mathbf{x}_0 .

Proposition 4.2. (Construction of an optimal rollout, see [4.0, p. 12] - forward pass). Starting at the initial condition $\mathbf{x}_0^* \sim p_{x_0}(\mathbf{x}_0^*)$, compute forward in time from $k = 0$ to $k = N - 1$ the control input via

$$\mathbf{u}_k^* = \mu_k^*(\mathbf{x}_k^*) = \arg \max_{\mathbf{u}_k \in \mathcal{U}} Q_k^*(\mathbf{x}_k^*, \mathbf{u}_k) \quad (4.72)$$

with

$$Q_k^*(\mathbf{x}_k^*, \mathbf{u}_k) = r_k(\mathbf{x}_k^*, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1}^* | \mathbf{x}_k^*, \mathbf{u}_k)} \{V_{k+1}^*(\mathbf{x}_{k+1})\} \quad (4.73)$$

and optimal state trajectory according to

$$\mathbf{x}_{k+1}^* \sim p_x(\mathbf{x}_{k+1}^* | \mathbf{x}_k^*, \mathbf{u}_k^*) . \quad (4.74)$$

See [4.0, p. 25] and [4.0, p. 11] for more information on the Dynamic Programming Algorithm. The *Finite Horizon Value Iteration Algorithm 3* shows how to use Value Iteration to find an optimal policy π^* .

Algorithm 3: Finite Horizon Value Iteration Algorithm

```

Data:  $\epsilon$  is a small number
Result: Find  $V_k^*(\mathbf{x}_k)$  and  $\pi^*$ ,  $\forall \mathbf{x}_k \in \mathcal{X}$ 
/* Initialization *//
1 Initialize  $V_k^\pi(\mathbf{x}_k)$  arbitrarily, expect the  $V_N^\pi(\mathbf{x}_N)$ ;
2  $V_k^\pi(\mathbf{x}_k) \leftarrow 0$ ,  $\forall \mathbf{x}_k \in \mathcal{X}$ ;
   /* Loop until convergence *//
3  $\Delta \leftarrow 0$ ;
   /* Optimal Value Determination *//
4 for  $k = N - 1, \dots, 0$  do
5   for each  $\mathbf{x}_k \in \mathcal{X}$  do
6      $v_k \leftarrow V_k^\pi(\mathbf{x}_k)$ ;
7      $V_k^\pi(\mathbf{x}_k) \leftarrow \max_{\mathbf{u}_k \in \mathcal{U}} [r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)} \{V_{k+1}^\pi(\mathbf{x}_{k+1})\}]$ ;
8      $\Delta \leftarrow \max(\Delta, |v_k - V_k^\pi(\mathbf{x}_k)|)$ ;
9   end
10 end
   /* Optimal Policy Determination *//
11  $\mu_k^*(\mathbf{x}_k) \leftarrow \arg \max_{\mathbf{u}_k \in \mathcal{U}} [r_k(\mathbf{x}_k, \mathbf{u}_k) + \gamma \mathbb{E}_{\mathbf{x}_{k+1} \sim p_x(\mathbf{x}_{k+1} | \mathbf{x}_k, \mathbf{u}_k)} \{V_{k+1}^*(\mathbf{x}_{k+1})\}]$ ;
12 return  $\pi^*$ ,  $\forall \mathbf{x}_k \in \mathcal{X}$ 
```

We will now have a look at some examples, in particular we investigate the well known *Linear Quadratic Regulator* (LQR) problem from Control Engineering.

Example 4.7 (The Linear Quadratic Regulator problem - infinite time horizon and deterministic case). Consider the problem of the discrete-time linear quadratic regulator (LQR), characterized by deterministic dynamics

$$\mathbf{x}_{k+1} = \Phi \mathbf{x}_k + \Gamma \mathbf{u}_k , \quad (4.75)$$

with k the discrete time index. Let the pair (Φ, Γ) be reachable. In this example, the state space $\mathcal{X} = \mathbb{R}^n$ and input space $\mathcal{U} = \mathbb{R}^m$ are infinite. The return over an infinite horizon, based on deterministic immediate rewards r_i , is expressed as

$$R_k(\tau_{[k:\infty]}) = \frac{1}{2} \sum_{i=k}^{\infty} r_i(\mathbf{x}_i, \mathbf{u}_i) = \frac{1}{2} \sum_{i=k}^{\infty} (\mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i + \mathbf{u}_i^\top \mathbf{R} \mathbf{u}_i) . \quad (4.76)$$

The return R_k depends on all future states $\mathbf{x}_k, \mathbf{x}_{k+1}, \dots$ and all future control inputs $\mathbf{u}_k, \mathbf{u}_{k+1}, \dots$. We select a stationary control law $\mathbf{u}_k = \mu(\mathbf{x}_k)$ and write the associated action-value function as

$$Q^\pi(\mathbf{x}_k, \mathbf{u}_k) = r(\mathbf{x}_k, \mathbf{u}_k) + V^\pi(\mathbf{x}_{k+1}) \quad (4.77a)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \Phi \mathbf{x}_k + \Gamma \mu(\mathbf{x}_k) . \quad (4.77b)$$

and the associated state-value function as

$$V^\pi(\mathbf{x}_k) = \frac{1}{2} \sum_{i=k}^{\infty} r_i^\pi = \frac{1}{2} \sum_{i=k}^{\infty} (\mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i + \mu^\top(\mathbf{x}_i) \mathbf{R} \mu(\mathbf{x}_i)) \quad (4.78a)$$

$$\text{s.t. } \mathbf{x}_{k+1} = \Phi \mathbf{x}_k + \Gamma \mu(\mathbf{x}_k) . \quad (4.78b)$$

Note that the state-value function for a fixed control law depends only on the initial state \mathbf{x}_k . A difference equation equivalent to the infinite sum (4.78a) is given by

$$\begin{aligned} V^\pi(\mathbf{x}_k) &= \frac{1}{2} (\mathbf{x}_k^\top \mathbf{Q} \mathbf{x}_k + \mu^\top(\mathbf{x}_k) \mathbf{R} \mu(\mathbf{x}_k)) + \frac{1}{2} \sum_{i=k+1}^{\infty} (\mathbf{x}_i^\top \mathbf{Q} \mathbf{x}_i + \mu^\top(\mathbf{x}_i) \mathbf{R} \mu(\mathbf{x}_i)) \\ &= \frac{1}{2} (\mathbf{x}_k^\top \mathbf{Q} \mathbf{x}_k + \mu^\top(\mathbf{x}_k) \mathbf{R} \mu(\mathbf{x}_k)) + V^\pi(\mathbf{x}_{k+1}) . \end{aligned} \quad (4.79)$$

Equation (4.79) is exactly the Bellman Equation (4.45) for the LQR problem. Assuming a quadratic state-value function

$$V^\pi(\mathbf{x}_k) = \frac{1}{2} \mathbf{x}_k^\top \mathbf{P} \mathbf{x}_k , \quad (4.80)$$

with positive-definite matrix \mathbf{P} , yields the Bellman Equation

$$2V^\pi(\mathbf{x}_k) = \mathbf{x}_k^\top \mathbf{P} \mathbf{x}_k = \mathbf{x}_k^\top \mathbf{Q} \mathbf{x}_k + \mu^\top(\mathbf{x}_k) \mathbf{R} \mu(\mathbf{x}_k) + \mathbf{x}_{k+1}^\top \mathbf{P} \mathbf{x}_{k+1} , \quad (4.81)$$

which, using the state equation (4.75), can be written as

$$2V^\pi(\mathbf{x}_k) = \mathbf{x}_k^\top \mathbf{Q} \mathbf{x}_k + \boldsymbol{\mu}^\top(\mathbf{x}_k) \mathbf{R} \boldsymbol{\mu}(\mathbf{x}_k) + (\Phi \mathbf{x}_k + \Gamma \boldsymbol{\mu}(\mathbf{x}_k))^\top \mathbf{P} (\Phi \mathbf{x}_k + \Gamma \boldsymbol{\mu}(\mathbf{x}_k)) . \quad (4.82)$$

Assuming a stationary state feedback control law $\boldsymbol{\mu}(\mathbf{x}_k) = -\mathbf{K} \mathbf{x}_k$ with control gain \mathbf{K} , we get

$$2V^\pi(\mathbf{x}_k) = \mathbf{x}_k^\top \mathbf{P} \mathbf{x}_k = \mathbf{x}_k^\top \mathbf{Q} \mathbf{x}_k + \mathbf{x}_k^\top \mathbf{K}^\top \mathbf{R} \mathbf{K} \mathbf{x}_k + \mathbf{x}_k^\top (\Phi - \Gamma \mathbf{K})^\top \mathbf{P} (\Phi - \Gamma \mathbf{K}) \mathbf{x}_k . \quad (4.83)$$

For all state trajectories, we find

$$(\Phi - \Gamma \mathbf{K})^\top \mathbf{P} (\Phi - \Gamma \mathbf{K}) - \mathbf{P} + \mathbf{Q} + \mathbf{K}^\top \mathbf{R} \mathbf{K} = \mathbf{0} . \quad (4.84)$$

This is a Lyapunov equation in Josephs form. That is, the Bellman Equation (4.81) for the discrete-time LQR problem is equivalent to a Lyapunov equation.

Example 4.8 (The Linear Quadratic Regulator problem - finite time horizon and stochastic case). In the discrete-time LQR design, the optimal policy π^* is determined for a linear time-invariant system with discrete-time evolution

$$\mathbf{x}_{k+1} = \Phi_k \mathbf{x}_k + \Gamma_k \mathbf{u}_k + \mathbf{w}_k . \quad (4.85)$$

The disturbance \mathbf{w}_k is characterized by the following properties, see [4.0],

$$\begin{aligned} \mathbb{E}\{\mathbf{w}_k\} &= \mathbf{0} \\ \mathbb{E}\{\mathbf{w}_k \mathbf{w}_k^\top\} &= \mathbf{W} . \end{aligned} \quad (4.86)$$

The deterministic rewards are given by

$$r_k(\mathbf{x}_k, \mathbf{u}_k) = \frac{1}{2} (\mathbf{x}_k^\top \mathbf{Q}_k \mathbf{x}_k + \mathbf{u}_k^\top \mathbf{R}_k \mathbf{u}_k) \quad (4.87)$$

$$r_N(\mathbf{x}_N) = \frac{1}{2} \mathbf{x}_N^\top \mathbf{Q}_N \mathbf{x}_N , \quad (4.88)$$

where $\mathbf{Q}_k \in \mathbb{R}^{n \times n}$ is positive semi-definite and $\mathbf{R}_k \in \mathbb{R}^{m \times m}$ is positive definite for all $k = 0, \dots, N$. To determine the optimal policy π^* in the form

$$\mathbf{u}_k = \boldsymbol{\mu}_k(\mathbf{x}_k) = -\mathbf{K}_k \mathbf{x}_k , \quad (4.89)$$

we apply the Dynamic Programming Algorithm from Theorem 4.1.

Starting from the terminal reward $r_N(\mathbf{x}_N)$, a backward recursion is performed based on the Bellman iteration (4.70). We initialize by

$$V_N(\mathbf{x}_N) = \frac{1}{2} \mathbf{x}_N^\top \mathbf{Q}_N \mathbf{x}_N \equiv \frac{1}{2} \mathbf{x}_N^\top \mathbf{P}_N \mathbf{x}_N . \quad (4.90)$$

Substituting the system dynamics $\mathbf{x}_N = \Phi_k \mathbf{x}_{N-1} + \Gamma_k \mathbf{u}_{N-1} + \mathbf{w}_{N-1}$ results in

$$\begin{aligned} V_{N-1}(\mathbf{x}_{N-1}) &= \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} [r_{N-1}(\mathbf{x}_{N-1}, \mathbf{u}_{N-1}) + \mathbb{E}_{\mathbf{w}_{N-1}} \{V_N^*(\mathbf{x}_N)\}] \\ &= \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \left[\mathbf{x}_{N-1}^\top \mathbf{Q}_{N-1} \mathbf{x}_{N-1} + \mathbf{u}_{N-1}^\top \mathbf{R}_{N-1} \mathbf{u}_{N-1} + \mathbf{x}_{N-1}^\top \mathbf{P}_N \mathbf{x}_{N-1} \right] . \end{aligned} \quad (4.91)$$

After rearranging, we have

$$\begin{aligned} V_{N-1}^*(\mathbf{x}_{N-1}) &= \frac{1}{2} \min_{\mathbf{u}_{N-1} \in \mathbb{R}^m} \left[\mathbf{x}_{N-1}^\top (\mathbf{Q}_{N-1} + \Phi_{N-1}^\top \mathbf{P}_N \Phi_{N-1}) \mathbf{x}_{N-1} \right. \\ &\quad + \mathbf{u}_{N-1}^\top (\mathbf{R}_{N-1} + \Gamma_{N-1}^\top \mathbf{P}_N \Gamma_{N-1}) \mathbf{u}_{N-1} \\ &\quad + 2\mathbf{u}_{N-1}^\top (\Gamma_{N-1}^\top \mathbf{P}_N \Phi_{N-1}) \mathbf{x}_{N-1} \\ &\quad \left. + \mathbb{E}_{\mathbf{w}_{N-1}} \{ \mathbf{w}_{N-1}^\top \mathbf{P}_N \mathbf{w}_{N-1} \} \right] . \end{aligned} \quad (4.92)$$

Note that this optimization problem is convex in \mathbf{u}_{N-1} as $\mathbf{R}_{N-1} + \Gamma_{N-1}^\top \mathbf{P}_N \Gamma_{N-1} > 0$. Therefore, any local minimum is a global minimum, and therefore we can simply apply the first order optimality conditions. Differentiating gives

$$\left(\frac{\partial}{\partial \mathbf{u}_{N-1}} V_{N-1}^* \right) (\mathbf{x}_{N-1}) = (\mathbf{R}_{N-1} + \Gamma_{N-1}^\top \mathbf{P}_N \Gamma_{N-1}) \mathbf{u}_{N-1} + (\Gamma_{N-1}^\top \mathbf{P}_N \Phi_{N-1}) \mathbf{x}_{N-1} \quad (4.93)$$

and setting this to zero yields

$$\mathbf{u}_{N-1}^* = -(\mathbf{R}_{N-1} + \Gamma_{N-1}^\top \mathbf{P}_N \Gamma_{N-1})^{-1} (\Gamma_{N-1}^\top \mathbf{P}_N \Phi_{N-1}) \mathbf{x}_{N-1} , \quad (4.94)$$

which we write

$$\mathbf{u}_{N-1}^* = -\mathbf{K}_{N-1} \mathbf{x}_{N-1} , \quad (4.95)$$

which is a time-varying linear feedback control law. Plugging this feedback policy into (4.92) results in

$$\begin{aligned} V_{N-1}^*(\mathbf{x}_{N-1}) &= \mathbf{x}_{N-1}^\top (\mathbf{Q}_{N-1} + \mathbf{K}_{N-1}^\top \mathbf{R}_{N-1} \mathbf{K}_{N-1}) \\ &\quad + (\Phi_{N-1} + \Gamma_{N-1} \mathbf{K}_{N-1})^\top \mathbf{P}_N (\Phi_{N-1} + \Gamma_{N-1} \mathbf{K}_{N-1}) \mathbf{x}_{N-1} . \end{aligned} \quad (4.96)$$

Critically, this implies that the cost is always a positive semi-definite quadratic function of the state. Because the optimal policy is always linear, and the optimal cost is always quadratic, the DP recursion may be recursively performed backward in time and the minimization may be performed analytically. Following the same

procedure, we can write the DP recursion for the discrete-time LQR controller. The optimal feedback gain \mathbf{K}_k according to control law $\mathbf{u}_k^* = -\mathbf{K}_k \mathbf{x}_k$ is obtained from a *backward calculation* from $k = N - 1$ to $k = 0$ starting from $\mathbf{P}_N = \mathbf{Q}_N$:

$$\mathbf{K}_k = (\mathbf{R}_k + \boldsymbol{\Gamma}_k^\top \mathbf{P}_{k+1} \boldsymbol{\Gamma}_k)^{-1} (\boldsymbol{\Gamma}_k^\top \mathbf{P}_{k+1} \boldsymbol{\Phi}_k) \quad (4.97)$$

$$\mathbf{P}_k = \mathbf{Q}_k + \mathbf{K}_k^\top \mathbf{R}_k \mathbf{K}_k + (\boldsymbol{\Phi}_k + \boldsymbol{\Gamma}_k \mathbf{K}_k)^\top \mathbf{P}_{k+1} (\boldsymbol{\Phi}_k + \boldsymbol{\Gamma}_k \mathbf{K}_k) . \quad (4.98)$$

To determine the optimal control sequence $\mathbf{U}_{[0:N-1]}^* = (\mathbf{u}_0^*, \mathbf{u}_1^*, \dots, \mathbf{u}_{N-1}^*)$ and optimal state sequence $\mathbf{X}_{[0:N]}^* = (\mathbf{x}_0^*, \mathbf{x}_1^*, \dots, \mathbf{x}_N^*)$, a *forward calculation* from $k = 0$ to $k = N - 1$ is performed starting from $\mathbf{x}_0 = \mathbf{x}_0^*$ using the optimal control law and discrete-time evolution

$$\begin{aligned} \mathbf{u}_0^* &= -\mathbf{K}_0 \mathbf{x}_0^* & \rightarrow \quad \mathbf{x}_1^* &= \boldsymbol{\Phi}_k \mathbf{x}_0^* + \boldsymbol{\Gamma}_k \mathbf{u}_0^* \\ &\vdots &&\vdots \\ \mathbf{u}_{N-1}^* &= -\mathbf{K}_{N-1} \mathbf{x}_{N-1}^* & \rightarrow \quad \mathbf{x}_N^* &= \boldsymbol{\Phi}_k \mathbf{x}_{N-1}^* + \boldsymbol{\Gamma}_k \mathbf{u}_{N-1}^* . \end{aligned} \quad (4.99)$$

The relation offers multiple insights. Even when $\boldsymbol{\Phi}, \boldsymbol{\Gamma}, \mathbf{Q}, \mathbf{R}$ are constant, the policy is still time-varying. This phenomenon arises because early control efforts lead to reduced state costs in all subsequent time steps. As the episode nears its end, the benefits of this early control wane, and the control effort decrease. However, for a linear time-invariant system where $(\boldsymbol{\Phi}, \boldsymbol{\Gamma})$ is reachable, the feedback gain \mathbf{K}_k approach a constant value as the episode length grows infinitely.

4.3 Model-free reinforcement learning

The methods described earlier require knowledge of the system dynamics. However, this is not the case in model-free reinforcement learning, where the dynamics are unknown. As a result, it becomes necessary to develop solution methods that do not rely on this explicit knowledge.

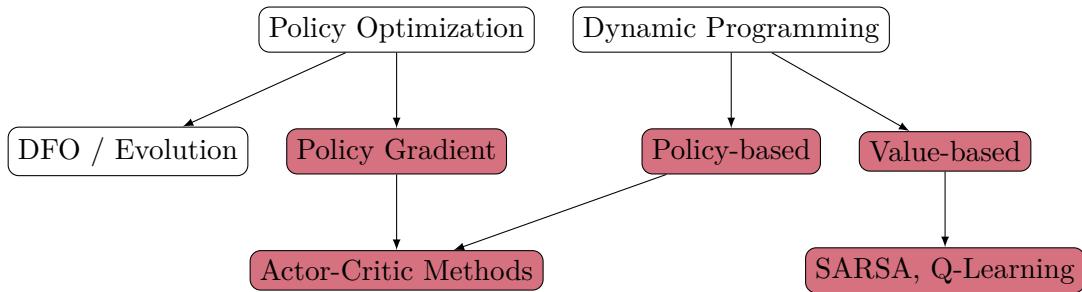


Figure 4.6: Classification of model-free RL algorithms [4.0, p. 4].

In reinforcement learning (RL), there are numerous options for approximating elements like policies, value functions, and dynamic models. Figure 4.6 provides a *classification of model-free RL algorithms*. There are three main approaches: (i.) *Policy Optimization*, (ii.) *Approximate Dynamic Programming*, and (iii.) *Hybrid methods*.

Policy Optimization techniques focus on optimizing the policy, the function mapping the agent's state to its next control input. They consider reinforcement learning as a numerical optimization problem where we optimize the expected reward with respect to the policy's parameters. There are two ways to optimize a policy: *Derivative-free Optimization* (DFO) algorithms and policy gradient methods. DFO algorithms work by perturbing the policy parameters, assessing the performance, and then moving towards better performance. They're simple to implement but scale poorly with increased parameters. On the other hand, *Policy Gradient Methods* estimate the policy improvement direction using various quantities measured by the agent, thus avoiding parameter perturbation. They are more complex to implement but can optimize larger policies than DFO algorithms.

Approximate Dynamic Programming (ADP), is based on learning value functions, predicting the agent's potential reward. ADP algorithms strive to satisfy certain consistency equations that the true value functions obey. Known algorithms for exact solutions in finite state-input RL problems are *policy-based* and *value-based*. They can be combined with function approximation in multiple ways; currently, leading descendants of value iteration focus on *approximating Q-functions*.

Lastly, *Actor-Critic* methods incorporate elements from both policy optimization and dynamic programming. They optimize a policy using value functions for faster optimization and often utilize ideas from approximate dynamic programming for fitting the value functions.

We discuss five key model-free reinforcement learning algorithms. Table 4.4 summarizes the use of these model-free reinforcement learning algorithms. It categorizes these algorithms based on their applicability in discrete and continuous input spaces. On-policy and

off-policy learning refer to how an algorithm learns a policy. On-policy learning improves the policy used to make decisions, while off-policy learning improves a different policy from the one used to generate data.

Algorithm	Policy	Discrete	Continuous	Tabular
Monte Carlo	on	Yes	No	Yes
SARSA	on	Yes	No	Yes
Q -Learning	off	Yes	No	Yes
Deep Q -Network	off	Yes	Yes	No
Policy Gradient	on/off	Yes	Yes	No

Table 4.4: Use of reinforcement learning algorithms in different input spaces.

In the following section, we will specifically focus on what are known as *tabular solution methods*, see [4.0, p. 23], which are applicable in scenarios with an infinite horizon as well as discrete input and state spaces.

4.3.1 Generalized Policy Iteration

Note that in model-free RL, the use of the action-value function $Q^\pi(\mathbf{x}, \mathbf{u})$ is favored over the state-value function $V^\pi(\mathbf{x})$ because $Q^\pi(\mathbf{x}, \mathbf{u})$ eliminates the need to know the state-transition and reward function explicitly, which is consistent with the goal of model-free RL to learn optimal policies without a model of the environment, cf. (4.66). Let us introduce the action-value table/matrix

$$\mathbf{Q}^\pi = \begin{bmatrix} Q^\pi(\mathbf{x}_0, \mathbf{u}_0) & Q^\pi(\mathbf{x}_0, \mathbf{u}_1) & \dots & Q^\pi(\mathbf{x}_0, \mathbf{u}_{|\mathcal{U}|-1}) \\ Q^\pi(\mathbf{x}_1, \mathbf{u}_0) & Q^\pi(\mathbf{x}_1, \mathbf{u}_1) & \dots & Q^\pi(\mathbf{x}_1, \mathbf{u}_{|\mathcal{U}|-1}) \\ \vdots & \ddots & \vdots & \vdots \\ Q^\pi(\mathbf{x}_{|\mathcal{X}|-1}, \mathbf{u}_0) & Q^\pi(\mathbf{x}_{|\mathcal{X}|-1}, \mathbf{u}_1) & \dots & Q^\pi(\mathbf{x}_{|\mathcal{X}|-1}, \mathbf{u}_{|\mathcal{U}|-1}) \end{bmatrix}. \quad (4.100)$$

The *Generalized Policy Iteration* (GPI) algorithm refers to an iterative procedure to improve the policy when combining policy evaluation and improvement:

$$\pi^{(0)} \xrightarrow{\text{evaluate}} (\mathbf{Q}^\pi)^{(0)} \xrightarrow{\text{improve}} \pi^{(1)} \xrightarrow{\text{evaluate}} (\mathbf{Q}^\pi)^{(1)} \dots \xrightarrow{\text{improve}} \pi^* \xrightarrow{\text{evaluate}} (\mathbf{Q}^\pi)^{(*)} \quad (4.101)$$

In GPI, the action-value function is approximated repeatedly to be closer to the true value of the current policy and in the meantime, the policy is improved repeatedly to approach optimality. This policy iteration process works and always converges to the optimality, independent of the granularity and other details of the two processes. Almost all reinforcement learning methods are well described as GPI [4.0, p. 86]:

Policy Improvement

Policy improvement refers to acting greedily and exploiting current knowledge. However, a balance between *exploration and exploitation* is critical in model-free RL. In view of (4.66), given an action-value function $Q^\pi(\mathbf{x}, \mathbf{u})$, the greedy policy

$$\pi(\mathbf{u} | \mathbf{x}) \leftarrow \begin{cases} 1 & \text{if } \mathbf{u} = \arg \max_{\mathbf{w} \in \mathcal{U}} Q^\pi(\mathbf{x}, \mathbf{w}) \\ 0 & \text{else ,} \end{cases} \quad (4.102)$$

is selected to maximize the expected reward based on current knowledge. Several strategies exist to balance exploration and exploitation:

ϵ -Greedy Method

With ϵ -greedy, at each step in state \mathbf{x} , the agent selects a random input with a fixed probability $\epsilon \in [0, 1]$:

- With probability $\frac{\epsilon}{|\mathcal{U}|} + 1 - \epsilon$, choose $\mathbf{u} = \arg \max_{\mathbf{w} \in \mathcal{U}} Q^\pi(\mathbf{x}, \mathbf{w})$ and
- with probability $\frac{\epsilon}{|\mathcal{U}|}$, select any input uniformly at random.

Hence, the ϵ -greedy policy is defined as

$$\pi(\mathbf{u} | \mathbf{x}) \leftarrow \begin{cases} \frac{\epsilon}{|\mathcal{U}|} + 1 - \epsilon & \text{if } \mathbf{u} = \arg \max_{\mathbf{w} \in \mathcal{U}} Q^\pi(\mathbf{x}, \mathbf{w}) \\ \frac{\epsilon}{|\mathcal{U}|} & \text{else .} \end{cases} \quad (4.103)$$

Compared to the greedy method (4.66), the ϵ -greedy method helps in exploring the input space.

Softmax Strategy

A downside of ϵ -greedy exploration is that it randomly picks any possible input for exploration. This means it's just as likely to pick the worst option as it is to pick the almost best one. In contrast, the softmax⁶ strategy utilizes input-selection probabilities which are determined by ranking the action-value function estimates using a Boltzmann distribution. The softmax policy is defined as

$$\pi(\mathbf{u} | \mathbf{x}) \leftarrow \frac{\exp(\beta Q^\pi(\mathbf{x}, \mathbf{u}))}{\sum_{\mathbf{w} \in \mathcal{U}} \exp(\beta Q^\pi(\mathbf{x}, \mathbf{w}))} . \quad (4.104)$$

Here, β is the inverse temperature parameter that controls the stochasticity of the policy. A high β makes the policy more deterministic, choosing the action with the highest Q -value with higher probability. A low β makes the policy more exploratory, distributing the selection probability more uniformly across inputs. In practice, the agent uses the Softmax Strategy to randomly select the next input by sampling from the distribution. More complex strategies exist, but many state-of-the-art algorithms still prefer the simple ϵ -greedy approach for its effectiveness and ease of implementation.

⁶The softmax function is often suggested to be more aptly named softargmax because it provides a softened version of the argmax function.

Policy Evaluation

If the dynamics of the environment are unknown, several methods can be used for *policy evaluation*. In policy evaluation, we apply the update rule

$$\hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \leftarrow \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) + \alpha \delta \quad (4.105)$$

to improve the estimated action-value function $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$. Here, $\alpha \in (0, 1]$ is the learning rate and δ denotes a residual. From a control engineering perspective, (4.105) is an integrator. It updates $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$ until the residual δ approaches zero.

Monte Carlo methods

The principle of Monte Carlo (MC) methods is to sample rollouts τ using the current policy π and approximate the expectation $Q^\pi(\mathbf{x}, \mathbf{u}) = \mathbb{E}_\pi\{\mathbf{R}_k \mid \mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}\}$ to compute the empirical mean return⁷

$$\hat{R}_k = \sum_{j=k}^N \gamma^{j-k} r_j \quad (4.106)$$

for each state-input pair. Then, with residual

$$\delta = \frac{1}{\alpha N(\mathbf{x}, \mathbf{u})} \sum_{k=1}^{N(\mathbf{x}, \mathbf{u})} [\hat{R}_k - \hat{Q}^\pi(\mathbf{x}, \mathbf{u})], \quad (4.107)$$

where $N(\mathbf{x}, \mathbf{u})$ is the total number of times the state-input pair is visited, we get from (4.105) the Monte Carlo update rule

$$\hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \leftarrow \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) + \frac{1}{N(\mathbf{x}, \mathbf{u})} \sum_{k=1}^{N(\mathbf{x}, \mathbf{u})} [\hat{R}_k - \hat{Q}^\pi(\mathbf{x}, \mathbf{u})]. \quad (4.108)$$

An update law for the visitation count $N(\mathbf{x}, \mathbf{u})$ is mathematically formulated as

$$N(\mathbf{x}, \mathbf{u}) \leftarrow N(\mathbf{x}, \mathbf{u}) + \mathbb{1}_{\mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}}.$$

Here, $\mathbb{1}_{\mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}}$ is the binary indicator function⁸. The equation (4.108) represents a simple average of the returns. After each episode, for each state-input pair visited, we calculate the return \hat{R}_k from that point until the end of the episode and then update the Q -value estimate for that state-action pair accordingly. MC methods need to learn from complete episodes to compute and all the episodes must eventually terminate.

SARSA: On-policy Temporal-Difference Learning

On-policy Temporal-Difference (TD) learning combines the principles of Monte Carlo methods with Dynamic Programming (DP). Under the certainty equivalence assumption, we get from the Bellman Expectation Equation, cf. (4.48),

$$Q^\pi(\mathbf{x}, \mathbf{u}) = r(\mathbf{x}, \mathbf{u}) + \gamma Q^\pi(\mathbf{x}', \mathbf{u}') . \quad (4.109)$$

⁷Also referred to as Monte Carlo return.

⁸If $\mathbf{x}_k = \mathbf{x}$ and $\mathbf{u}_k = \mathbf{u}$, the function $\mathbb{1}_{\mathbf{x}_k = \mathbf{x}, \mathbf{u}_k = \mathbf{u}}$ returns 1, otherwise, it returns 0.

The key idea in On-policy TD learning is to update the action-value function $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$ using the TD residual

$$\delta = r(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}^\pi(\mathbf{x}', \mathbf{u}') - \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \quad (4.110)$$

and to follow the SARSA update rule, see, e.g., [4.0, p. 120],

$$\hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \leftarrow \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) + \alpha \left[r(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}^\pi(\mathbf{x}', \mathbf{u}') - \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \right]. \quad (4.111)$$

For $\alpha = 0$, the value $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$ remains unchanged, while for $\alpha = 1$, the old value $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$ is replaced entirely by the so-called TD-target⁹ $r(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}^\pi(\mathbf{x}', \mathbf{u}')$. This becomes evident when we express (4.111) as a first-order low-pass filter in the form

$$\hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \leftarrow (1 - \alpha) \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) + \alpha \left[r(\mathbf{x}, \mathbf{u}) + \gamma \hat{Q}^\pi(\mathbf{x}', \mathbf{u}') \right]. \quad (4.112)$$

On-policy Temporal-Difference learning is called SARSA because of this data sequence used for calculation – State, Action, Reward, State, Action. The estimate \hat{Q}^π is updated based on its own estimation, which is a form of bootstrapping, as described in [4.0, p. 89]. TD learning can learn from incomplete rollouts and hence we don't need to track the rollouts up to termination.

Q-learning: Off-policy Temporal-Difference Learning

Q-learning is an extension of TD-learning that goes beyond predicting the value function Q^π to enable the determination of an optimal policy π^* . The update rule for *Q*-learning involves updating the estimate \hat{Q}^π at each time step, considering both the observed state \mathbf{x} and the corresponding action \mathbf{u} . Note that the Bellman Optimality Equations (4.50) reads as

$$Q^*(\mathbf{x}, \mathbf{u}) = \mathbb{E}_\pi \left\{ r(\mathbf{x}, \mathbf{u}) + \gamma \max_{\mathbf{w} \in \mathcal{U}} Q^*(\mathbf{x}', \mathbf{w}) \right\}. \quad (4.113)$$

Since we do not have access to the optimal values Q^* , the idea in Off-policy TD learning is to update the estimate of the action-value function $\hat{Q}^\pi(\mathbf{x}, \mathbf{u})$ using the residual

$$\delta = r(\mathbf{x}, \mathbf{u}) + \gamma \max_{\mathbf{w} \in \mathcal{U}} \hat{Q}^\pi(\mathbf{x}', \mathbf{w}) - \hat{Q}^\pi(\mathbf{x}, \mathbf{u}). \quad (4.114)$$

The *Q*-learning update law than becomes

$$\hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \leftarrow \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) + \alpha \left[r(\mathbf{x}, \mathbf{u}) + \gamma \max_{\mathbf{w} \in \mathcal{U}} \hat{Q}^\pi(\mathbf{x}', \mathbf{w}) - \hat{Q}^\pi(\mathbf{x}, \mathbf{u}) \right]. \quad (4.115)$$

The estimated action-value function \hat{Q}^π directly approximates optimal action-value function Q^* , independent of the policy being followed. There are a couple of extensions of TD- and *Q*-learning that are designed to address its limitations and improve its applicability and performance. Key extensions include:

⁹In machine learning, a target value refers to the actual value you aim to predict with your model. It's the outcome or label that the algorithm is being trained to forecast in supervised learning.

- Double Q -Learning [4.0, p. 125]: Addresses the overestimation of action-state values by decoupling the improvement and evaluation of the input in the action-state value update.
- Eligibility Traces [4.0, p. 287]: The more theoretical view is that Eligibility Traces are a bridge from TD to Monte Carlo methods (forward view). According to the other view, an Eligibility Trace is a temporary record of the occurrence of an event, such as the visiting of a state or the taking of an action (backward view). The trace marks the memory parameters associated with the event as eligible for undergoing learning changes. $TD(\lambda)$ [4.0, p. 293]: Extends the basic TD method by considering a series of TD errors for all time steps until the end of the episode, weighted by a factor λ . The factor $\lambda \in [0, 1]$ allows you to balance the trade-off between bias (accuracy) and variance (stability) in the learning updates.
- $Q(\lambda)$ [4.0, p. 312]: This extension applies the concept of Eligibility Traces to Q -Learning. In $Q(\lambda)$, every input-value pair (Q -value) has an associated eligibility trace, which decays over time but gets bumped up when visited. The action-state value updates are then applied to all pairs proportionally to their eligibility, allowing for quicker propagation of information through the state-input space.

Finally, we compare Monte Carlo, SARSA and Q -Learning and draw some conclusions:

Aspect	Monte Carlo	SARSA	Q -Learning
Type of Learning	Episodic	On-policy TD	Off-policy TD
Convergence Speed	Slow (high variance)	Moderate	Fast (high bias)
Policy Dependency	Optimal policy	Sensitive to current policy	Learns optimal policy from any policy
Suitability	Episodic tasks	Close to current strategy tasks	Complex environments with suboptimal actions

Table 4.5: Comparison of Monte Carlo, SARSA, and Q -Learning.

4.3.2 Approximate solution methods

For large state and input spaces, methods where the state value function V^π or the action-state value function Q^π are represented by means of a table are no longer manageable. The reason for this is the high memory requirements of the table, since a values must be stored for each state or state-input pair. Moreover, for a continuous state space, an infinite number of values would have to be stored, and the computational cost is considerable has to be applied to each state or state-input pair, respectively. Instead of a table, a function approximator parameterized with the vector ϕ can be used to represent the value functions in the form

$$V^\pi(\mathbf{x}) \approx \hat{V}^\pi(\mathbf{x}; \phi) = \hat{V}_\phi^\pi(\mathbf{x}) \quad (4.116a)$$

$$Q^\pi(\mathbf{x}, \mathbf{u}) \approx \hat{Q}^\pi(\mathbf{x}, \mathbf{u}; \phi) = \hat{Q}_\phi^\pi(\mathbf{x}, \mathbf{u}) . \quad (4.116b)$$

As a result, the memory requirement is significantly reduced to only the parameter vector ϕ . Additionally, by employing a function approximator, the agent can generalize from previously visited states or state-input pairs to unvisited ones. In an ideal scenario where $V^\pi(\mathbf{x})$ and $Q^\pi(\mathbf{x}, \mathbf{u})$ are known, one could utilize supervised learning methods to approximate them with \hat{V}^π and \hat{Q}^π . By generating rollouts

$$\tau^{(e)} = (\mathbf{x}_0^{(e)}, \mathbf{u}_0^{(e)}, \mathbf{x}_1^{(e)}, \mathbf{u}_1^{(e)}, \dots, \mathbf{x}_{N-1}^{(e)}, \mathbf{u}_{N-1}^{(e)}, \mathbf{x}_N^{(e)}) \quad (4.117)$$

for $e = 0, 1, \dots, E$, with data

$$\mathcal{D} = \left\{ \mathbf{x}_k^{(i)}, V^\pi(\mathbf{x}_k^{(e)}) \right\}_{e=0}^E \quad (4.118a)$$

$$\mathcal{D} = \left\{ (\mathbf{x}_k^{(e)}, \mathbf{u}_k^{(e)}), Q^\pi(\mathbf{x}_k^{(e)}, \mathbf{u}_k^{(e)}) \right\}_{e=0}^E , \quad (4.118b)$$

the *regression task* can be formulated as follows

$$\min_{\phi} \frac{1}{|\mathcal{D}|} \sum_{e \in \mathcal{D}} \left(\hat{V}_\phi^\pi(\mathbf{x}_k^{(e)}) - V^\pi(\mathbf{x}_k^{(e)}) \right)^2 \quad (4.119)$$

and

$$\min_{\phi} \frac{1}{|\mathcal{D}|} \sum_{e \in \mathcal{D}} \left(\hat{Q}_\phi^\pi(\mathbf{x}_k^{(e)}, \mathbf{u}_k^{(e)}) - Q^\pi(\mathbf{x}_k^{(e)}, \mathbf{u}_k^{(e)}) \right)^2 . \quad (4.120)$$

Deep Q -Network

Q -learning can experience instability and divergence when combined with nonlinear function approximators and bootstrapping. The Deep Q -Network (DQN) addresses these issues by introducing two key innovations to stabilize and enhance the training process:

- *Experience replay*: Instead of using sequential rollouts for updates, DQN stores a collection of steps $\mathbf{e}_k = (\mathbf{x}_k, \mathbf{u}_k, r_k, \mathbf{x}_{k+1})$ in a replay memory $\mathcal{D}_k = \{\mathbf{e}_0, \dots, \mathbf{e}_k\}$. This memory contains a diverse range of experiences from multiple past episodes. During training, mini-batches of experiences are randomly sampled from this memory

to update the Q -values. This technique not only improves the efficiency of data utilization but also breaks the correlation between consecutive samples and mitigates the non-stationary distribution issues, leading to more stable and reliable learning.

- *Periodically updated the target network:* In standard Q -learning, the same network estimates both the current and the target Q -values, leading to a moving target problem that can cause harmful correlations and oscillations. DQN addresses this by employing two separate networks: a primary network with parameters ϕ for the current Q -value estimation and a target network with parameters ϕ^- that remains fixed for a set number of steps C . The target network's sole purpose is to generate the stable target Q -values for the updates. Every C steps, the primary network's weights are copied to the target network, ensuring that the targets are only periodically updated, which significantly enhances the stability of the learning process.

The objective function for DQN is formulated as follows:

$$L(\phi) = \mathbb{E}_{(\mathbf{x}, \mathbf{u}, r, \mathbf{x}') \sim U(\mathcal{D})} \left[\left(r(\mathbf{x}, \mathbf{u}) + \gamma \max_{\mathbf{w} \in \mathcal{U}} \hat{Q}_{\phi^-}^{\pi}(\mathbf{x}', \mathbf{w}) - \hat{Q}_{\phi}^{\pi}(\mathbf{x}, \mathbf{u}) \right)^2 \right] \quad (4.121)$$

Here, $U(\mathcal{D})$ represents a uniform distribution over the replay memory \mathcal{D} and ϕ^- denotes the parameters of the target network.

4.3.3 Policy Gradients

Policy gradients are a class of algorithms in reinforcement learning that optimize policies directly. They work by calculating the gradient of the expected reward concerning the policy parameters and then adjusting the parameters in the direction that increases the expected reward. This approach is particularly powerful for high-dimensional or continuous input spaces and allows for stochastic policies, offering a more nuanced way of exploring and exploiting the environment. They form the basis for many advanced reinforcement learning methods and are fundamental to understanding and developing new algorithms in the field.

The term *Policy Gradient* (PG) covers methods where a policy is parameterized by the parameter vector θ , i.e.

$$\pi(\mathbf{u}_k \mid \mathbf{x}_k; \theta) = \pi_{\theta}(\mathbf{u}_k \mid \mathbf{x}_k) . \quad (4.122)$$

With this parameterized policy, a multivariate distribution

$$p^{\pi_{\theta}}(\tau) = p_0(\mathbf{x}_0) \prod_{k=0}^{N-1} \pi_{\theta}(\mathbf{u}_k \mid \mathbf{x}_k) p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k) \quad (4.123)$$

can be introduced, which gives the probability density of observing a N -step rollout τ under a the policy π_{θ} . Thus, with the *discounted return* (4.8), we obtain the parameterized action-value function

$$Q_k^{\pi_{\theta}}(\mathbf{x}_k, \mathbf{u}_k) = \mathbb{E}_{\pi_{\theta}}\{\mathsf{R}_k(\tau) \mid \mathbf{x}_k = \mathbf{x}_k, \mathbf{u}_k = \mathbf{u}_k\} , \quad (4.124)$$

and analogously the parameterized state-value function

$$V_k^{\pi_\theta}(\mathbf{x}_k) = \mathbb{E}_{\pi_\theta}\{\mathbf{R}_k(\boldsymbol{\tau}) \mid \mathbf{x}_k = \mathbf{x}_k\}. \quad (4.125)$$

Now, the optimal control problem can be formulated as an optimization over the parameter vector $\boldsymbol{\theta}$, see [4.0], p. 96], in the form

$$\boldsymbol{\theta}^* = \arg \max_{\boldsymbol{\theta}} J(\boldsymbol{\theta}), \quad (4.126)$$

with expected return

$$J(\boldsymbol{\theta}) = \int_{\mathcal{T}} p^{\pi_\theta}(\boldsymbol{\tau}) R_0(\boldsymbol{\tau}) d\boldsymbol{\tau}. \quad (4.127)$$

To solve (4.126) with (4.127), a gradient ascent method can be employed

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} + \alpha \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}^{(i)}}, \quad (4.128)$$

where $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ represents the gradient of the objective function $J(\boldsymbol{\theta})$, and α denotes the step size or learning rate.

4.3.4 Stochastic Policy Gradient

To obtain the Policy Gradient $\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$ for the policy (4.122), the following procedure is performed. From (4.127), follows

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \int_{\mathcal{T}} p^{\pi_\theta}(\boldsymbol{\tau}) R_0(\boldsymbol{\tau}) d\boldsymbol{\tau} = \int_{\mathcal{T}} \nabla_{\boldsymbol{\theta}} p^{\pi_\theta}(\boldsymbol{\tau}) R_0(\boldsymbol{\tau}) d\boldsymbol{\tau}. \quad (4.129)$$

Using the identity¹⁰

$$\nabla_{\boldsymbol{\theta}} p^{\pi_\theta}(\boldsymbol{\tau}) = p^{\pi_\theta}(\boldsymbol{\tau}) \frac{\nabla_{\boldsymbol{\theta}} p^{\pi_\theta}(\boldsymbol{\tau})}{p^{\pi_\theta}(\boldsymbol{\tau})} = p^{\pi_\theta}(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}} \ln(p^{\pi_\theta}(\boldsymbol{\tau})), \quad (4.130)$$

we obtain

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \int_{\mathcal{T}} p^{\pi_\theta}(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}}(p^{\pi_\theta}(\boldsymbol{\tau})) R_0(\boldsymbol{\tau}) d\boldsymbol{\tau} = \mathbb{E}_{\pi_\theta}\{\nabla_{\boldsymbol{\theta}} \ln(p^{\pi_\theta}(\boldsymbol{\tau})) R_0(\boldsymbol{\tau})\}. \quad (4.131)$$

The term $\ln p^{\pi_\theta}(\boldsymbol{\tau})$ can be split into three components using (4.123), i.e.,

$$\ln(p^{\pi_\theta}(\boldsymbol{\tau})) = \ln(p_0(\mathbf{x}_0)) + \sum_{k=0}^{N-1} \ln(\pi_\theta(\mathbf{u}_k \mid \mathbf{x}_k)) + \sum_{k=0}^{N-1} \ln(p_x(\mathbf{x}_{k+1} \mid \mathbf{x}_k, \mathbf{u}_k)). \quad (4.132)$$

Taking the gradient with respect to the parameter vector eliminates the terms that depend on the system dynamics results in

$$\nabla_{\boldsymbol{\theta}} \ln(p^{\pi_\theta}(\boldsymbol{\tau})) = \nabla_{\boldsymbol{\theta}} \sum_{k=0}^{N-1} \ln(\pi_\theta(\mathbf{u}_k \mid \mathbf{x}_k)) = \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln(\pi_\theta(\mathbf{u}_k \mid \mathbf{x}_k)). \quad (4.133)$$

Substituting (4.133) into (4.131) with (4.124) yields the *Stochastic Policy Gradient* (SPG) for finite time horizons¹¹, see [4.0],

¹⁰Known as the log-derivative-trick.

¹¹ $\nabla_{\boldsymbol{\theta}} \ln \pi_\theta$ is the so-called score function.

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) R_0 \right\}. \quad (4.134)$$

There are two important variants of the Stochastic Policy Gradient. Firstly, the Stochastic Policy Gradient can be formulated in term of the return R_k according to (4.8). Causality requires that future inputs and policies at time k do not depend on past rewards at time i . Thus, for $i < k$, we have

$$0 = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \{ \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) \mathbf{r}_i \}. \quad (4.135)$$

With the definition of the return R_k , we get from (4.134)

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) \left(\sum_{i=0}^{k-1} \gamma^{i-k} r_i + \underbrace{\sum_{i=k}^{N-1} \gamma^{i-k} r_i}_{=R_k} \right) \right\} \quad (4.136)$$

and with the causality condition (4.135), we find another formulation of the Stochastic Policy Gradient

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) R_k \right\}. \quad (4.137)$$

Secondly, the Stochastic Policy Gradient can be written in term of the state-value function $Q_k^{\pi_{\boldsymbol{\theta}}}(\mathbf{x}_k, \mathbf{u}_k)$. Define $\boldsymbol{\tau}_{[:,k]} = (\mathbf{x}_0, \mathbf{u}_0, \dots, \mathbf{x}_k, \mathbf{u}_k)$ as the rollout up to time k , and $\boldsymbol{\tau}_{[k,:]}$ as the remainder of the rollout after that. Using the *law of total expectation* (A.9), we can break up (4.134) into:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \sum_{k=0}^{N-1} \mathbb{E}_{\boldsymbol{\tau}_{[:,k]} \sim \pi_{\boldsymbol{\theta}}} \left\{ \mathbb{E}_{\boldsymbol{\tau}_{[k,:]} \sim \pi_{\boldsymbol{\theta}}} \left\{ \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) R_k \mid \boldsymbol{\tau}_{[:,k]} \right\} \right\}. \quad (4.138)$$

The gradient of the log-probability is constant with respect to the inner expectation, due to its dependency on \mathbf{x}_k and \mathbf{u}_k , and it can be extracted:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \sum_{k=0}^{N-1} \mathbb{E}_{\boldsymbol{\tau}_{[:,k]} \sim \pi_{\boldsymbol{\theta}}} \left\{ \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) \mathbb{E}_{\boldsymbol{\tau}_{[k,:]} \sim \pi_{\boldsymbol{\theta}}} \left\{ R_k \mid \boldsymbol{\tau}_{[:,k]} \right\} \right\}. \quad (4.139)$$

Lastly, the Markov property implies

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) Q_k^{\pi_{\boldsymbol{\theta}}}(\mathbf{x}_k, \mathbf{u}_k) \right\}. \quad (4.140)$$

Since $p^{\pi_{\boldsymbol{\theta}}}(\boldsymbol{\tau})$ is unknown, the expectation $\mathbb{E}_{\pi_{\boldsymbol{\theta}}} \{ \cdot \}$ is approximated from $e = 1, \dots, E$ rollouts $\boldsymbol{\tau}^{(e)}$, collected in the set $\mathcal{D} = \{ \boldsymbol{\tau}^{(e)} \}$, by the *empirical mean* as

$$\hat{\mathbf{g}} = \widehat{\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta})} = \frac{1}{|\mathcal{D}|} \sum_{\boldsymbol{\tau} \in \mathcal{D}} \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{x}_k | \mathbf{u}_k)) Q_k^{\pi_{\boldsymbol{\theta}}}(\mathbf{x}_k, \mathbf{u}_k). \quad (4.141)$$

The estimated policy gradient has high variance¹², resulting in poor convergence properties [0]. To reduce variance, a *baseline* $b_k(\mathbf{x}_k) \in \mathbb{R}$ can be subtracted from Q_k^π in (4.134) without affecting the expectation, see [4.0, p. 98] for a proof, resulting in

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} (\ln \pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) (Q_k^{\pi_{\boldsymbol{\theta}}}(\mathbf{x}_k, \mathbf{u}_k) - b_k(\mathbf{x}_k)) \right\}. \quad (4.142)$$

Intuitively, the variance is reduced when subtracting a baseline because the operation effectively re-centers the reward signal around a mean value. A widely used baseline is the value function $b_k(\mathbf{x}_k) = V_k^\pi(\mathbf{x}_k)$. See [4.0] for different baselines and variants of the Stochastic Policy Gradient. From (4.142), we then get

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta}}} \left\{ \sum_{k=0}^{N-1} \nabla_{\boldsymbol{\theta}} \ln (\pi_{\boldsymbol{\theta}}(\mathbf{u}_k | \mathbf{x}_k)) A_k^{\pi_{\boldsymbol{\theta}}}(\mathbf{x}_k, \mathbf{u}_k) \right\}, \quad (4.143)$$

with *Advantage function* according to (4.39). The Advantage function gives the amount of how much the current input is better than what we would usually do in that state.

Note 4.5. A multivariate Gaussian distribution (or multivariate normal distribution) is described by a mean vector $\boldsymbol{\mu}$ and a covariance matrix $\boldsymbol{\Sigma}$. A diagonal Gaussian distribution is a special case where the covariance matrix only has entries on the diagonal. As a result, we can represent it by a vector. Hence, a popular choice for policy model $\pi_{\boldsymbol{\theta}}(\mathbf{u} | \mathbf{x})$ is the *diagonal Gaussian policy model*. A diagonal Gaussian policy always has a neural network that maps from states to mean inputs $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x})$ and standard deviations are typically represented by standalone parameters. Given the mean input $\boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x})$ and standard deviation $\boldsymbol{\sigma}_{\boldsymbol{\theta}}$, and a vector \mathbf{z} of noise from a spherical Gaussian ($\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$), an input sample can be computed with

$$\mathbf{u} = \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}) + \boldsymbol{\sigma}_{\boldsymbol{\theta}} \odot \mathbf{z}, \quad (4.144)$$

where \odot denotes the element wise product of two vectors. In a diagonal Gaussian policy model, exploration is achieved by the additive noise term. The mean determines the expected input, while the diagonal elements of the covariance matrix determine the amount of exploration or noise added to each input. When an input is sampled from this distribution, the noise leads to exploration by encouraging slightly different inputs to be taken, even in the same state. This process allows the model to explore various strategies and potentially discover better ones. Note that the log-likelihood of a m -dimensional input, for a diagonal Gaussian with mean $\boldsymbol{\mu} = \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x})$ and standard deviation $\boldsymbol{\sigma} = \boldsymbol{\sigma}_{\boldsymbol{\theta}}$, is given by

$$\ln (\pi_{\boldsymbol{\theta}}(\mathbf{u} | \mathbf{x})) = -\frac{1}{2} \left(\sum_{i=1}^m \left(\frac{(\mathbf{u}[i] - \boldsymbol{\mu}[i])^2}{\boldsymbol{\sigma}[i]^2} + 2 \ln (\boldsymbol{\sigma}[i]) \right) + m \ln (2\pi) \right). \quad (4.145)$$

¹²We say that a method has high variance if you get different results when you run it multiple times. A method with less variance will give you more similar results when you run it multiple times.

4.3.5 Actor-Critic Methods

Two main components in policy gradient are the policy model and the value function. It makes a lot of sense to learn the value function in addition to the policy, since knowing the value function can assist the policy update, such as by reducing gradient variance in policy gradients, and that is exactly what the Actor-Critic method does. Actor-critic methods consist of two models:

- *Critic* updates the value function parameters ϕ and depending on the algorithm it could be action-value function Q_ϕ^π or state-value function V_ϕ^π .
- *Actor* updates the policy parameters θ for π_θ , in the direction suggested by the critic.

Algorithm 4 shows a vanilla Actor-Critic Policy Gradient Algorithm.

Algorithm 4: Vanilla Actor-Critic Policy Gradient Algorithm

```

/* Initialization */ 
1 Initialize policy parameters  $\theta$  arbitrarily;
2 Initialize value function parameters  $\phi$  arbitrarily;
3 for  $e = 0, 1, 2, \dots, E$  episodes do
    /* Collect rollouts using policy  $\pi^{(e)} = \pi(\theta^{(e)})$  in environment */
    4  $\mathcal{D}^{(e)} = \{\tau^{(1)}, \dots, \tau^{(E)}\}$ ;
    5 for  $k = 0, 1, 2, \dots, N - 1$  time steps do
        /* Compute Monte Carlo return */ 
        6  $\hat{R}_k^{(e)} \leftarrow \sum_{j=k}^{N-1} \gamma^{j-k} r_j$ ;
        /* Compute Advantage estimates */ 
        7  $\hat{A}_k^{(e)} \leftarrow \hat{R}_k^{(e)} - \hat{V}_{\phi^{(e)}}^\pi$ ;
    8 end
    /* Estimate policy gradient */ 
    9  $\hat{\mathbf{g}}^{(e)} \leftarrow \frac{1}{|\mathcal{D}^{(e)}|} \sum_{\tau \in \mathcal{D}^{(e)}} \sum_{k=0}^{N-1} \nabla_\theta \ln (\pi_\theta(\mathbf{u}_k | \mathbf{x}_k)) \Big|_{\theta=\theta^{(e)}} \hat{A}_k^{(e)}$ ;
    /* Actor learning - compute policy update */ 
    10  $\theta^{(e)} \leftarrow \theta^{(i)} + \alpha \hat{\mathbf{g}}^{(e)}$ ;
    /* Critic learning - fit the value function by regression */ 
    11  $\phi^{(e)} \leftarrow \arg \min_\phi \frac{1}{|\mathcal{D}^{(e)}| N} \sum_{\tau \in \mathcal{D}^{(e)}} \sum_{k=0}^{N-1} (\hat{V}_\phi^\pi(\mathbf{x}_k) - \hat{R}_k^{(e)})^2$ ;
12 end

```

4.3.6 Off-Policy Policy Gradient

The vanilla versions of the actor-critic method are exclusively on-policy: training samples are collected according to the *target policy*, which is the same policy we aim to optimize. A simpler strategy involves two policies: the *target policy* π , which is studied and optimized, and the *behavior policy* β , which drives exploratory actions. Here, learning is based on

data "off" the target policy, and the entire mechanism is labeled as off-policy learning. Off-policy methods offer several distinct advantages, including:

- The off-policy approach doesn't require a complete rollout, and it can reuse past episodes (via experience replay), resulting in significantly improved sample efficiency¹³.
- The sample collection employs a behavior policy different from the target policy, facilitating more effective exploration.

Nearly all off-policy methods employ *importance sampling* to estimate expected values from a distribution different than the sample source. In the context of off-policy learning, returns are weighted by the importance sampling ratio, which quantifies the relative likelihood of specific trajectories under the target and behavior policies. Given an initial state \mathbf{x}_k , the probability of a future state-input trajectory under any policy π is given by (4.31). Thus, the relative probability of the trajectory under the target and behavior policies, that is the *importance sampling ratio*, is

$$\rho_{[k:N-1]} = \frac{\prod_{j=k}^{N-1} p_x(\mathbf{x}_{j+1} | \mathbf{x}_j, \mathbf{u}_j) \pi(\mathbf{u}_j | \mathbf{x}_j)}{\prod_{j=k}^{N-1} p_x(\mathbf{x}_{j+1} | \mathbf{x}_j, \mathbf{u}_j) \beta(\mathbf{u}_j | \mathbf{x}_j)} = \prod_{j=k}^{N-1} \frac{\pi(\mathbf{u}_j | \mathbf{x}_j)}{\beta(\mathbf{u}_j | \mathbf{x}_j)}. \quad (4.146)$$

Although the trajectory probabilities depend on the MDP's transition probabilities, which are generally unknown, they appear identically in both the numerator and denominator, and thus cancel. The importance sampling ratio ends up depending only on the two policies and the sequence, not on the MDP.

Remember that we aim to estimate the expected returns under the target policy π , but only possess returns R_k^β from the behavior policy β . These returns exhibit an incorrect expectation denoted by $V^\beta(\mathbf{x}) = \mathbb{E}_\beta\{R_k^\beta | \mathbf{x}_k = \mathbf{x}\}$. To correct this, importance sampling is employed. Using the ratio $\rho_{[k:N-1]}$, we can transform these returns to align with the desired expectation

$$V^\pi(\mathbf{x}) = \mathbb{E}_\pi\left\{\rho_{[k:N-1]} R_k^\beta | \mathbf{x}_k = \mathbf{x}\right\}. \quad (4.147)$$

Now let's see how off-policy policy gradient is computed. Since the training observations are sampled by $\mathbf{u}_k \sim \beta(\mathbf{u}_k | \mathbf{x}_k)$, we can rewrite the gradient following certain calculations, see [4.0] for a proof,

$$\nabla_\theta J(\theta) = \mathbb{E}_\beta\left\{\sum_{k=0}^{N-1} \frac{\pi_\theta(\mathbf{u}_k | \mathbf{x}_k)}{\beta(\mathbf{u}_k | \mathbf{x}_k)} \nabla_\theta \ln(\pi_\theta(\mathbf{u}_k | \mathbf{x}_k)) Q_k^{\pi_\theta}(\mathbf{x}_k, \mathbf{u}_k)\right\}. \quad (4.148)$$

Here, $\frac{\pi_\theta(\mathbf{u}_k | \mathbf{x}_k)}{\beta(\mathbf{u}_k | \mathbf{x}_k)}$ is the *importance weight*. In essence, when implementing policy gradient in the off-policy setting, we can adjust it with a weighted sum, where the weight is the ratio of the target policy to the behavior policy.

¹³Sample efficiency in the context of control technology and machine learning refers to the ability of a system to achieve high performance with a limited number of data samples.

4.3.7 Proximal Policy Optimization

Schulman proposed *Proximal Policy Optimization* (PPO) in 2017, and it has since become a widely used method in continuous model-free RL due to its simplicity and strong performance. PPO-clip is an actor-critic method and makes use of a generalized advantage estimate and a clipped surrogate objective function, see [4.0]. It's worth noting that there are other variants of PPO, such as PPO-KL, which incorporate an adaptive Kullback-Leibler penalty term, see [4.0]. Our discussion will focus exclusively on PPO-Clip. This approach is straightforward to implement and has been shown to work well in practice. It doesn't require a second-order optimization process, making it computationally less intensive.

Generalized Advantage Estimation

PPO employs a Generalized Advantage Estimation (GAE) as its advantage function. The purpose of GAE is to significantly reduce the variance of the estimator while keeping the bias introduced as low as possible [4.0, p. 136]. This goal mirrors the fundamental principles of Eligibility Traces, cf. [4.0, p. 125]. Our present discussion is fundamentally grounded in the methodologies outlined in the supplementary material provided by [4.0]. Generally, the Monte Carlo return (4.106) serves as an unbiased estimator of the expected return at a specific state. However, each reward r_k may vary due to the environment dynamics, leading to a high variance in the overall estimation. To mitigate this, an employ an n -step return

$$\begin{aligned}\hat{R}_{[k:k+n]} &= r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots + \gamma^n r_{k+n} \\ &= \sum_{j=k}^{n-1} \gamma^{j-k} r_j + \gamma^n \hat{V}^\pi(\mathbf{x}_{k+n}) ,\end{aligned}\tag{4.149}$$

where we estimate the remaining return using a value function estimate $\hat{V}^\pi(\mathbf{x})$. It offers a lower-variance but slightly biased estimate by truncating the sum of returns after n steps. Particular instances such as $\hat{R}_{[k:k+1]} = r_k + \gamma \hat{V}^\pi(\mathbf{x}_{k+1})$, used in Q -learning, and $\hat{R}_{[k,\infty]} = \sum_{j=k}^N \gamma^{j-k} r_j = \hat{R}_k$, which reverts to the original Monte Carlo return, illustrate the variance-bias trade-off. An alternative for balancing bias and variance is the λ -return, computed as a weighted average of n -step returns, see [4.0, p. 289], and defined as

$$\hat{R}_k(\lambda) = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{R}_{[k:k+n]} .\tag{4.150}$$

Assuming all rewards after step N are zero, such that $\hat{R}_{[k:k+n]} = \hat{R}_{[k:N]}$ for all $n \geq N - k$, the infinite sum can be calculated according to

$$\hat{R}_k(\lambda) = (1 - \lambda) \sum_{n=1}^{N-k-1} \lambda^{n-1} \hat{R}_{[k:k+n]} + \lambda^{N-k-1} \hat{R}_{[k:N]} .\tag{4.151}$$

Here, $\lambda = 0$ yields $\hat{R}_{[k,k+1]}$, and $\lambda = 1$ provides the full Monte Carlo return $\hat{R}_{[k,\infty]}$. Intermediate values of $\lambda \in (0, 1)$ produce interpolants that can be used to balance the

bias and variance of the value estimator. Updating the value function with the λ -return leads to the so-called TD(λ) algorithm, and its application for advantage estimation results in the Generalized Advantage Estimator GAE(λ):

$$\hat{A}_k^\pi(\lambda) = \hat{R}_k(\lambda) - \hat{V}^\pi(\mathbf{x}_k) . \quad (4.152)$$

For instance, setting $\lambda = 0$ yields the expression

$$\hat{A}_k^\pi(0) = \hat{R}_{[k:k+1]} - \hat{V}^\pi(\mathbf{x}_k) = r_k + \gamma \hat{V}_\phi^\pi(\mathbf{x}_{k+1}) - \hat{V}^\pi(\mathbf{x}_k) = \delta_k , \quad (4.153)$$

which is equivalent to the TD residual (4.110). Notice that $\hat{R}_k(0) = \hat{R}_{[k:k+1]} = r_k + \gamma \hat{V}^\pi(\mathbf{x}_{k+1})$ is the 1-step estimator for $\hat{Q}^\pi(\mathbf{x}_k, \mathbf{u}_k)$. Choosing a value closer to zero introduces more bias due to the immediate approximation, resulting in lower variance, as discussed in [4.0]. Conversely, employing $\lambda = 1$ leads to high variance and nearly zero bias due to the summation of rewards, i.e.,

$$\hat{A}_k^\pi(1) = \hat{R}_{[k:N]} - \hat{V}^\pi(\mathbf{x}_k) = \sum_{j=k}^N \gamma^{j-k} r_j . \quad (4.154)$$

which is equivalent to the Monte Carlo return (4.106). To summarize, adjusting the value of λ allows control over the tradeoff between bias and variance. A smaller value, such as $\lambda = 0$, reduces variance at the expense of introducing more bias, while a larger value, like $\lambda = 1$, results in higher variance with minimal bias due to reward summation.

Clipping

PPO-clip updates policies via

$$\theta^{(i+1)} = \arg \max_{\theta} \mathbb{E}_{\mathbf{u}_k \sim \pi_{\theta^{(i)}}} \left\{ L(\mathbf{x}_k, \mathbf{u}_k, \theta^{(i)}, \theta) \right\} , \quad (4.155)$$

typically taking multiple steps of (usually mini-batch) Stochastic Gradient Ascent (SGA) to maximize the objective function. The (clipped surrogate) objective function L is given by, see [4.0],

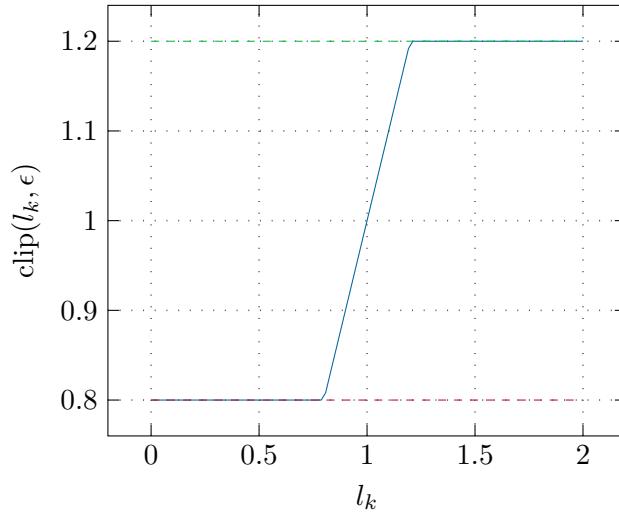
$$L(\mathbf{x}_k, \mathbf{u}_k, \theta^{(i)}, \theta) = \min(l_k A^{\pi_{\theta^{(i)}}}(\mathbf{x}_k, \mathbf{u}_k), \text{clip}(l_k, \epsilon) A^{\pi_{\theta^{(i)}}}(\mathbf{x}_k, \mathbf{u}_k)) , \quad (4.156)$$

with likelihood ratio

$$l_k = l_k(\theta^{(i)}, \theta) = \frac{\pi_\theta(\mathbf{x}_k, \mathbf{u}_k)}{\pi_{\theta^{(i)}}(\mathbf{x}_k, \mathbf{u}_k)} = \frac{\text{new policy}}{\text{old policy}} \quad (4.157)$$

and clipping operator

$$\text{clip}(l_k, \epsilon) = \begin{cases} 1 - \epsilon & \text{if } l_k < 1 - \epsilon \\ 1 + \epsilon & \text{if } l_k > 1 + \epsilon \\ l_k & \text{else} . \end{cases} \quad (4.158)$$

Figure 4.7: Clipping function (4.158) with $\epsilon = 0.2$.

The hyperparameter ϵ is a small positive constant, typically set to a value like 0.2, see Figure 4.7 for an illustration. The purpose of this ϵ -clipping is to prevent overly large policy updates, thereby maintaining stability in the learning process. This is achieved by ensuring that the current policy does not deviate excessively from the older one. The term l_k represents the probability ratio between the current and old policy. Let's interpret the implications of l_k :

- If $l_k > 1$, it signifies that for a given state \mathbf{x}_k , the corresponding action \mathbf{u}_k is more probable under the current policy than it was under the old policy.
- Conversely, if l_k is between 0 and 1, the action \mathbf{u}_k is less likely under the current policy compared to the old one.

This likelihood ratio serves as a straightforward method to gauge the divergence between the old and current policy. Algorithm 5 shows vanilla Proximal Policy Optimization with Clipping.

Algorithm 5: Proximal Policy Optimization with Clipping

```

/* Initialization */
```

1 Initialize policy parameters θ , old policy parameters θ_{old} ;

2 Initialize value function parameters ϕ ;

3 **for** $e = 1, 2, 3, \dots$ *episodes* **do**

/* Collect rollouts using the old policy $\pi^{(e)} = \pi(\theta^{(e)})$ in
environment */

4 $\mathcal{D}^{(e)} = \{\tau^{(1)}, \dots, \tau^{(E)}\}$;

5 **for** $k = 0, 1, 2, \dots, N - 1$ *time steps* **do**

/* Compute λ -returns */

6 $\hat{R}_k^{(e)}(\lambda) \leftarrow (1 - \lambda) \sum_{n=1}^{N-k-1} \lambda^{n-1} \hat{R}_{[k:k+n]} + \lambda^{N-k-1} \hat{R}_{[k:N]}$;

/* Compute generalized advantage estimates */

7 $\hat{A}_k^{(e)} \leftarrow \hat{R}_k^{(e)}(\lambda) - \hat{V}_\phi^\pi(\mathbf{x}_k)$;

8 **end**

/* Actor learning - update policy parameters via SGA in
mini-batch */

9 $l_k^{(e)} \leftarrow \frac{\pi_\theta(\mathbf{x}_k, \mathbf{u}_k)}{\pi_{\theta^{(e)}}(\mathbf{x}_k, \mathbf{u}_k)}$;

10 $\theta^{(e)} \leftarrow \arg \max_{\theta} \frac{1}{|\mathcal{D}^{(e)}|} \sum_{\tau \in \mathcal{D}^{(e)}} \sum_{k=0}^{N-1} \min(l_k^{(e)} \hat{A}_k^{(e)}, \text{clip}(\epsilon, l_k^{(e)}) \hat{A}_k^{(e)})$;

/* Critic learning - fit the value function by regression */

11 $\phi^{(e)} \leftarrow \arg \min_{\phi} \frac{1}{|\mathcal{D}^{(e)}| N} \sum_{\tau \in \mathcal{D}^{(e)}} \sum_{k=0}^{N-1} (\hat{V}_\phi^\pi(\mathbf{x}_k) - \hat{R}_k^{(e)}(\lambda))^2$;

12 **end**

4.4 Literatur

- [4.0] M. P. Deisenroth, A. Faisal, and C. S. Ong, *Mathematics for Machine Learning*. Cambridge University Press, 2020. [Online]. Available: <https://mml-book.github.io/>.
- [4.0] E. T. Jaynes, *Probability Theory: The Logic of Science*. Cambridge University Press, 2013, vol. 1.
- [4.0] D. P. Bertsekas, *Dynamic Programming and Optimal Control*. Athena Scientific, 2005, vol. 1.
- [4.0] M. L. Puterman, *Markov Decision Processes*. John Wiley & Sons, 2005.
- [4.0] C. Szepesvari, *Algorithms for Reinforcement Learning*. Morgan & Claypool, 2009.
- [4.0] R. S. Sutton and A. G. Barto, *Reinforcement Learning*. MIT Press, 2018.
- [4.0] D. P. Bertsekas, *Reinforcement Learning and Optimal Control*. Athena Scientific, 2019.
- [4.0] M. Sugiyama, *Statistical Reinforcement Learning*. CRC Press, 2015.
- [4.0] D. Silver, *Reinforcement learning: An introduction*, 2015. [Online]. Available: <http://www.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>.
- [4.0] F. L. Lewis, D. L. Varbie, and V. Syrmos, *Optimal Control*. John Wiley & Sons, 2012.
- [4.0] M. Vidyasagar, “A tutorial introduction to reinforcement learning,” 2023. [Online]. Available: <https://arxiv.org/abs/2304.00803v1>.
- [4.0] R. Stengel, *Optimal Control and Estimation*. Dover Publications, 1986.
- [4.0] J. Schulman, “Optimizing expectations: From deep reinforcement learning to stochastic computation graphs,” Ph.D. dissertation, University of California, Berkeley, 2016.
- [4.0] J. Kober, J. A. Bagnell, and J. Peters, “Reinforcement learning in robotics: A survey,” *The International Journal of Robotics Research*, pp. 1238–1274, 2013.
- [4.0] J. Peters and S. Schaal, “Reinforcement learning of motor skills with policy gradients,” *Neural Networks*, vol. 21, no. 4, pp. 682–697, 2008.
- [4.0] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, “High-dimensional continuous control using generalized advantage estimation,” in *ICLR 2015*, 2015. [Online]. Available: [arXiv:1506.02438](https://arxiv.org/abs/1506.02438).
- [4.0] T. Degrif, M. White, and R. S. Sutton, “Off-policy actor-critic,” in *Proceedings of the 29 th International Conference on Machine Learning*, Edinburgh, Scotland, UK, 2012.
- [4.0] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” 2017. [Online]. Available: [arXiv:1707.06347](https://arxiv.org/abs/1707.06347).
- [4.0] L. Graesser and W. Keng, *Foundation of Deep Reinforcement Learning*. Addison-Wesley, 2020.

- [4.0] X. Peng, P. Abbeel, S. Levine, and M. van de Panne, “Deepmimic: Example-guided deep reinforcement learning of physics-based character skills,” 2018. [Online]. Available: [arXiv:1804.02717](https://arxiv.org/abs/1804.02717).
- [4.0] C. M. Bischof, *Pattern Recognition and Machine Learning*. Springer, 2006.

A Basics of probability theory

A brief summary of basic results from probability theory is given next. For an introduction to probability theory, we refer to the [4.0] and [4.0]. We use sans-serif letters such as x, \mathbf{x} , and \mathbf{X} to represent *random variables* (scalars, vectors and matrices) and serif letters such as x, \mathbf{x} , and \mathbf{X} to represent the corresponding *deterministic variables* or observations.

A.1 Variables

Quantitative variables are categorized as either discrete or continuous.

- Categorical Variable: Categorical variables contain a finite number of categories or distinct groups. Categorical data might not have a logical order. Categorical data may not follow a logical sequence. Examples include variables like gender, type of material, and methods of payment.
- Discrete Variable: Discrete variables are numeric variables that have a countable number of values between any two values. Discrete variables are exclusively numeric. Common examples include the number of customer complaints or the count of defects or flaws.
- Continuous Variable: Continuous variables are numeric variables that have an infinite number of values between any two values. They can be in the form of numeric data or date/time values. Typical examples are the measurement of a component's length or the exact moment a payment is processed.

A.2 Random variables

A random variable is a mapping from the sample space Ξ to the real numbers \mathbb{R} . It assigns a numerical value to each outcome $\xi \in \Xi$. We denote random variables using sans-serif letters, i.e., $x : \Xi \rightarrow \mathbb{R}$, whereas the numerical values are denoted using serif letters, i.e., $x = x(\xi)$. If the range of the random variable, i.e., the set of possible values of x , is countable, we say that the random variable is discrete. Otherwise, it is continuous.

Example A.1. An example of a discrete random variable is the mapping that assigns integers from 1 to 6 to the outcomes of a die roll. For instance, $x(\square) = 1$, $x(\square) = 2$, $x(\square) = 3$, $x(\square) = 4$, $x(\square) = 5$, $x(\square) = 6$. This process is called categorical variable encoding in machine learning.

Each outcome is mapped to a real number and consequently any event $\mathcal{A} \subseteq \Xi$ corresponds to a subset $\mathcal{A}' \subseteq \mathbb{R}$ of the real line. Conversely, any subset \mathcal{A}' of the real line identifies an event in the sample space. The relationship between \mathcal{A} and \mathcal{A}' can be written as

$$\mathcal{A}' = \{x : x = x(\xi), \xi \in \mathcal{A}\} \quad \text{and} \quad \mathcal{A} = \{\xi : x(\xi) \in \mathcal{A}'\}. \quad (\text{A.1})$$

Based on the correspondence, we denote \mathcal{A}' as an event and assign the probability

$$\mathbb{P}(x(\xi) \in \mathcal{A}') = \mathbb{P}(\xi \in \mathcal{A}). \quad (\text{A.2})$$

A.3 Probability distribution

We write $x \sim p(x)$ to denote that a random variable x is distributed according to $p(x)$, with observation x . We will use the expression *probability distribution* not only for the discrete probability mass function (PMF) but also for the continuous probability density function (PDF). This is in line with most machine learning literature, cf. [4.0]. We rely on the context to distinguish the different use of the phrase *probability distribution*.

A.4 Expectation

The expectation or expected value of a deterministic function $g : \mathbb{R} \rightarrow \mathbb{R}$ of a univariate random variable $x \sim p(x)$ is given by

$$\mathbb{E}_{x \sim p(x)}\{g(x)\} = \sum_{x \in \mathcal{X}} p(x)g(x). \quad (\text{A.3})$$

In the continuous case, we have

$$\mathbb{E}_{x \sim p(x)}\{g(x)\} = \int_{\mathcal{X}} p(x)g(x)dx. \quad (\text{A.4})$$

For multivariate random variables, we define the expected value element wise, see, e.g., [4.0, p. 187] for more information.

A.5 Probability of an event

Next we show how to rewrite the probability of arbitrary events as an expectation. This involves the introduction of the indicator function of a set $\mathcal{A}' \subseteq \mathbb{R}$, defined as a function taking values 0 or 1:

$$I_{\mathcal{A}'}(x) = \begin{cases} 1, & x \in \mathcal{A}', \\ 0, & \text{else} \end{cases}. \quad (\text{A.5})$$

Setting $g(x) = I_{\mathcal{A}'}(x)$, results in

$$\mathbb{E}_{x \sim p(x)}\{I_{\mathcal{A}'}(x)\} = \mathbb{P}\{x \in \mathcal{A}'\} = \begin{cases} \sum_{x \in \mathcal{A}'} p(x) & \text{if } x \text{ is discrete} \\ \int_{\mathcal{A}'} p(x)dx & \text{if } x \text{ is continuous} \end{cases}. \quad (\text{A.6})$$

Since the indicator function $I_{\mathcal{A}'}(x)$ denotes whether x belongs to \mathcal{A}' or not, its expectation measures the average with which x falls inside \mathcal{A}' , which coincides with the probability of x being in \mathcal{A}' .

A.6 Conditional expectation

We can use conditional distributions to compute conditional expectations. More specifically, the conditional expectation of a random variable $z = g(x)$ given an observation $y = y$ is defined as

$$\mathbb{E}_{x \sim p(x)}\{g(x) | y = y\} = \begin{cases} \sum_{x \in \mathcal{X}} g(x)p(x | y) & \text{if } x \text{ is discrete} \\ \int_{\mathcal{X}} g(x)p(x | y)dx & \text{if } x \text{ is continuous} \end{cases}. \quad (\text{A.7})$$

Example A.2. For example, if \mathcal{X} and \mathcal{U} are continuous, the rollout $\tau \in \mathcal{T}$ is continuous as well. For a given policy π and function $R : \mathbb{R} \rightarrow \mathbb{R}$, then

$$\mathbb{E}_{\tau \sim p(\cdot)}\{R(\tau) | \pi\} = \int_{\mathcal{T}} p(\tau | \pi)R(\tau)d\tau = \int_{\mathcal{T}} p^\pi(\tau)R(\tau)d\tau = \mathbb{E}_{\tau \sim p^\pi(\cdot)}\{R(\tau)\}. \quad (\text{A.8})$$

Here, $p^\pi(\tau)$ is the probability distribution of rollouts obtained by executing the policy π .

A.7 Law of total expectation

The law of total expectation (also law of iterated expectations) states that if x is a random variable whose expected value $\mathbb{E}(x)$ is defined, and y is any random variable on the same probability space, then

$$\mathbb{E}(x) = \mathbb{E}(\mathbb{E}(x | y)), \quad (\text{A.9})$$

i.e., the expected value of the conditional expected value of x given y is the same as the expected value of x .

A.8 Rules of probability

The two fundamental rules of probability theory are, see [4.0, p. 187],

$$\text{Sum rule } p(x) = \begin{cases} \sum_{y \in \mathcal{Y}} p(x, y) & \text{if } y \text{ is discrete} \\ \int_{\mathcal{Y}} p(x, y)dy & \text{if } y \text{ is continuous} \end{cases} \quad (\text{A.10})$$

$$\text{Product rule } p(x, y) = p(y | x)p(x). \quad (\text{A.11})$$

Here, $p(x, y)$ represents the *joint probability*, one says "the probability of x and y ". Similarly, $p(y | x)$ denotes the *conditional probability*, it is the "the probability of y given x ". Lastly, $p(x)$ is the *marginal probability* and is simply phrased as "the probability of x ". The sum and product rule can be combined to get the *law of total probability*. It is a fundamental rule relating marginal probabilities to conditional probabilities

$$p(y) = \mathbb{E}_{x \sim p(x)}\{p(y | x)\} = \begin{cases} \sum_{x \in \mathcal{X}} p(x, y) = \sum_{x \in \mathcal{X}} p(y | x)p(x) & \text{if } x, y \text{ are discrete} \\ \int_{\mathcal{X}} p(x, y)dx = \int_{\mathcal{X}} p(y | x)p(x)dx & \text{if } x, y \text{ are continuous} \end{cases} \quad (\text{A.12})$$

A.9 The chain rule of conditional probabilities

Any joint probability distribution over many random variables can be decomposed into conditional distributions over only one variable. This observation is known as the chain rule

$$p(x_0, \dots, x_N) = p(x_0) \prod_{k=0}^{N-1} p(x_{k+1} | x_k, \dots, x_0). \quad (\text{A.13})$$

This follows directly from the definition of product rule (A.11). For example, if we apply the definition twice, we get

$$\begin{aligned} p(x_0, x_1, x_2) &= p(x_0 | x_1, x_2)p(x_1, x_2) \\ p(x_1, x_2) &= p(x_1 | x_2)p(x_2) \\ p(x_0, x_1, x_2) &= p(x_0 | x_1, x_2)p(x_1 | x_2)p(x_2) . \end{aligned} \quad (\text{A.14})$$

A.10 Bernoulli distribution

The (discrete) Bernoulli distribution will be referred to as $\mathcal{B}(x; \mu)$, where $x \in \{0, 1\}$. It represents for example the result of flipping coin. The values 1 occurs with success probability μ and $x = 0$ occurs with failure probability $1 - \mu$, respectively. Thus for random variable x , which is Bernoulli distributed, we write

$$x \sim \mathcal{B}(x; \mu) = \mu^x(1 - \mu)^{1-x} . \quad (\text{A.15})$$

A.11 Normal distribution

The (continuous) normal distribution will be referred to as $\mathcal{N}(x; \mu, \sigma^2)$, where $x \in \mathbb{R}$. Thus when a random variable x is normally distributed with mean μ and standard deviation σ , one may write

$$x \sim \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x - \mu}{\sigma}\right)^2\right) . \quad (\text{A.16})$$