

Kapitel 10

Dynamische Programmierung I (deterministische Probleme)

10.1 Einleitung; Bellman-Prinzip und Bellman-Gleichung

Neben Variationsrechnung und Kontrolltheorie (nach Pontryagin) steht mit der *dynamischen Programmierung* (nach Bellman) ein auf den ersten Blick vollkommen andersartiges Konzept zur Behandlung dynamischer Optimierungsprobleme zur Verfügung. Obwohl die dynamische Programmierung gerade für *stochastische* Probleme (wo also Unsicherheit über die zukünftige Entwicklung der Zustandsvariable und damit auch die zukünftigen Erträge besteht) besser geeignet ist als die Kontrolltheorie, beschränken wir uns in diesem Kapitel auf deterministische Probleme. Da verschiedene Aspekte der dynam. Programmierung, wie z.B. ihr rekursiver Charakter, *closed-loop controls* usw., am deutlichsten bei Problemen mit einer diskreten Einteilung des Zeitintervalls $[0, T]$ in Perioden $t = 0, 1, \dots, T-1$ hervortreten, beginnen wir mit *zeitdiskreten* Problemen der Form:

$$\max_{u_0, \dots, u_{T-1} \in \mathcal{U}} \sum_{t=0}^T f_t(x_t, u_t), \quad x_0 \text{ gegeben, } x_{t+1} = g_t(x_t, u_t) \quad (t = 0, \dots, T-1)$$

statt der bisherigen zeitstetigen Probleme (die später im Kapitel behandelt werden):

$$\max_{u(t) \in \mathcal{U} \text{ für } t \in [0, T]} \int_0^T f(t, x(t), u(t)) dt + S(T, x(T)), \quad x(0) = x_0, \quad \dot{x}(t) = g(t, x(t), u(t))$$

(Durchgängig bezeichnet x die Zustandsvariable, u die Kontrollvariable; bei der dynamischen Programmierung wird standardmäßig ein freies x_T bzw. $x(T)$ unterstellt).

Grundlegend für die dynamische Programmierung ist das **Bellman-Prinzip** (BELLMAN (1957)):

An optimal policy has the property that, whatever the initial state and initial decision are, the remaining decisions must constitute an optimal policy with regard to the state resulting from the first decision.

Das heißt: Eine Lösung des „Gesamtpblems“ stellt ab jedem Zeitpunkt $t \in \{1, \dots, T\}$ eine Lösung des „Restproblems“ dar, das ausgeht vom zum Zeitpunkt t erreichten Zustand $x = x_t$. Gegeben ein Zeitpunkt t und ein zu diesem Zeitpunkt erreichter Zustand x , spielt es für die zukünftigen optimalen Entscheidungen $u_s, s \geq t$ also keine Rolle, wie man in den Zustand x gelangt ist. Dieses Prinzip gilt gleichermaßen für zeit-diskrete wie zeit-stetige Probleme. (Es gilt

auch für stochastische Probleme, sofern der stochastische Prozess, dem x_t folgt, die *Markov-Eigenschaft* hat). Das legt es nahe (bzw. erlaubt es bei stochastischen Problemen), das „vom Zeitpunkt t aus bestehende Restproblem“ – und insbesondere dessen Optimalwert V_t – als *Funktion des Startzustands x* zu betrachten. Betrachtet man nun die Restprobleme in zwei aufeinanderfolgenden Zeitpunkten t und $t+1$ (bzw. t und $t+dt$ bei zeitstetigen Problemen), so kann man den Wert der Zielfunktion des Restproblems zur Zeit t zerlegen in den unmittelbaren Ertrag $f_t(x, u)$ (bzw. $f(t, x, u) dt$) und den Zielfunktionswert des Restproblems zum Zeitpunkt $t+1$ (bzw. $t+dt$). Die vom Zeitpunkt t aus zu treffenden Entscheidungen kann man analog unterteilen in die aktuelle Entscheidung $u = u_t$ und die zukünftigen Entscheidungen $u_s, s > t$. Nach dem Bellman-Prinzip spielt für die zukünftigen Entscheidungen aber nur der Ausgangszustand \tilde{x} eine Rolle, der sich nach der aktuellen Entscheidung ergibt. Ist u die unmittelbare Entscheidung, so ist $\tilde{x} = g_t(x, u)$ und der Optimalwert des zukünftigen Problems lautet $V_{t+1}(g_t(x, u))$. Wenn die Funktion $V_{t+1}(\tilde{x})$ bekannt ist, bestimmt sich die unmittelbare Entscheidung u also durch Maximierung von $f_t(x, u) + V_{t+1}(g_t(x, u))$, und das Ergebnis dieser Maximierung ist gerade $V_t(x)$. Es besteht also folgende Beziehung zwischen den Funktionen V_t und V_{t+1} , die wir als **Bellman-Gleichung** (hier: eines zeitdiskreten deterministischen Problems) bezeichnen:¹

$$V_t(x) = \max_{u \in \mathcal{U}} \{f_t(x, u) + V_{t+1}(g_t(x, u))\} \quad [\text{Bellman-Gl. im zeitdiskreten determin. Fall}] \quad (*)$$

Ausgehend von $V_T(x) = \max_{u \in \mathcal{U}} f_T(x, u)$ kann man mit der Beziehung $(*)$ *rekursiv* (in zeitlich absteigender Reihenfolge) $V_{T-1}(x), V_{T-2}(x), \dots, V_0(x)$ bestimmen und erhält damit schließlich den Optimalwert $V(x) = V_0(x)$ des Ausgangsproblems als Funktion des Startzustands $x = x_0$.

Die Optimalwertfunktionen der Restprobleme, $V_{T-n}(x)$, an sich sind oft gar nicht so interessant. Was man mit der Bellman-Rekursion aber auch bekommt, ist eine Darstellung der optimalen Kontrolle $u = u_t$ zur Zeit t als Funktion des Zustands $x = x_t$ zu diesem Zeitpunkt: $u_t = u_t^*(x)$. (Dies enthält die eigentliche Bedeutung des Begriffs „*optimal policy*“: Man sagt, welche Entscheidung u optimalerweise zu treffen ist, wenn man den Zustand x des Systems kennt). Man spricht dann auch von **closed-loop controls** oder **feedback controls**. Will man die **open-loop controls** $u = u_t$ bzw. $u = u(t)$, d.h. die optimale Kontrolle als Funktion der Zeit t (nicht als Funktion des Zustands x), wie sie in der Kontrolltheorie ermittelt werden, aus den closed loop controls bestimmen, so ist eine erneute Rekursion, diesmal in zeitlich aufsteigender Reihenfolge, erforderlich: Ausgehend von x_0 bestimmt man $u_0 = u_0^*(x_0)$, dann $x_1 = g_0(x_0, u_0)$, dann $u_1 = u_1^*(x_1)$ usw.

Die Bellman-Gleichung erfasst sehr kompakt das Entscheidungsproblem, vor dem beispielsweise ein Unternehmen zum Zeitpunkt t steht, dessen Erträge (pro Periode) in den folgenden Perioden $s = t, t+1, \dots, T$ als Funktion $f_s(x, u)$ des Zustands x des Unternehmens (beispielsweise Kapitalstock) und der Entscheidungsvariable u (z.B. Höhe der Investition pro Periode) durch $f_s(x, u)$ gegeben sind: Ist x der Zustand zum Zeitpunkt t , so muss die unmittelbare Entscheidung u die Summe aus unmittelbarem Ertrag, das ist $f_t(x, u)$, und dem Wert des Unternehmens, der sich nach der Entscheidung ergeben würde, das ist $V_{t+1}(g_t(x, u))$, maximieren. Beim „Wert des Unternehmens, der sich nach der Entscheidung ergeben würde“ muss das UN nämlich davon ausgehen, dass es auch ab der Periode $t+1$ optimal entscheidet – und dies ist in V_{t+1} erfasst. Dass V_{t+1} im Zustand $x_{t+1}(x, u)$ ausgewertet werden muss, berücksichtigt, dass die unmittelbare Entscheidung $u = u_t$ einen Effekt auf den Zustand x_{t+1} zu Beginn der nächsten Periode hat.

¹Analoge Beziehungen gelten auch für zeitstetige und stochastische Probleme; bei zeitstetigen Problemen erhält man i.w. eine partielle Differentialgleichung für die Funktion $V(t, x)$, die dann als Hamilton-Jacobi-Bellman-Gleichung bezeichnet wird. Die Bellman-Gleichung des Problems stellt jeweils Dreh- und Angelpunkt der dynamischen Programmierung dar und wird auch als **fundamentale Gleichung für Optimalität** bezeichnet.

10.2 Bellman-Gleichung für zeitdiskrete determinist. Probleme

10.2.1 Formulierung im Gegenwartswert (present value)

Wir betrachten in diesem Abschnitt folgendes dynamische Optimierungsproblem in der Zustandsvariable $\mathbf{x} \in \mathbb{R}^n$ und der Entscheidungsvariable $\mathbf{u} \in \mathbb{R}^m$ (da es prinzipiell und auch formal überhaupt keinen Unterschied macht, lassen wir gleich mehrdimensionale Zustands- und Kontrollvariablen zu, obwohl in den einfachsten Beispielen immer $n = m = 1$ sein wird):

$$\max_{\mathbf{u}_t \in \mathcal{U}_t, t=0, \dots, T} \sum_{t=0}^T f_t(\mathbf{x}_t, \mathbf{u}_t) \quad \text{unt. d. NB} \quad \begin{cases} \mathbf{x}_0 \text{ gegeben,} \\ \mathbf{x}_{t+1} = \mathbf{g}_t(\mathbf{x}_t, \mathbf{u}_t) \quad \text{für } t = 0, \dots, T-1 \end{cases} \quad (1)$$

Hierbei sind, für $t = 0, \dots, T-1$, $f_t(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$ und $\mathbf{g}_t(\mathbf{x}, \mathbf{u})$ gegebene Funktionen, die den „Momentanertrag in der Periode von t zu $t+1$ “ (f_t) bzw. die „Bewegungsgleichung der Zustandsvariable“ (\mathbf{g}_t) beschreiben. Man beachte, dass \mathbf{g}_t hier nicht den Zuwachs $\Delta \mathbf{x}_t = \mathbf{x}_{t+1} - \mathbf{x}_t$, sondern direkt \mathbf{x}_{t+1} wiedergibt. Die Funktion $f_T(\mathbf{x}, \mathbf{u})$ entspricht dem früheren „Terminalwert“ $S(T, \mathbf{x})$, der hier auch von der Kontrollgröße \mathbf{u} zum Zeitpunkt T abhängen kann. Anders als in der Kontrolltheorie gilt das Folgende auch dann ohne Modifikationen, wenn die „Kontrollbereiche“ $\mathcal{U}_t \subset \mathbb{R}^m$ vom Zustand \mathbf{x} zur Zeit t abhängen (d.h. die Kontrollvariablenrestriktion könnte auch $\mathbf{u}_t \in \mathcal{U}_t(\mathbf{x}_t)$ lauten). Man beachte, dass wir das Problem als Funktion des Anfangswertes \mathbf{x} der Zustandsvariable betrachten.

Lösungskonzept: Gesucht sind Funktionen $\mathbf{u}_t^*(\mathbf{x})$, so dass man mit der Rekursion $\mathbf{u}_t = \mathbf{u}_t^*(\mathbf{x}_t)$, $\mathbf{x}_{t+1} = \mathbf{g}_t(\mathbf{x}_t, \mathbf{u}_t)$ das Problem (1) für jeden Anfangswert \mathbf{x}_0 löst. Wir nennen dies eine **Markov-Politik-Lösung** des Problems (obwohl der Begriff bei determinist. Problemen unangebracht erscheint; der wichtige Punkt bei einer ‘Markov-policy’ ist, dass die aktuelle Entscheidung \mathbf{u} auf Basis des aktuellen Systemzustands \mathbf{x} getroffen werden kann – die Vergangenheit von \mathbf{x} spielt dafür keine Rolle).

Zur Ermittlung der optimalen Politiken, d.h. der Funktionen $\mathbf{u}_t^*(\mathbf{x})$, betrachtet die dynamische Programmierung den Optimalwert des Restproblems, das sich vom Zeitpunkt t an noch stellt, als Funktion von dessen Initialzustand $\mathbf{x} = \mathbf{x}_t$:

$$V_t(\mathbf{x}) := \max_{\mathbf{u}_s \in \mathcal{U}_s, s=t, \dots, T} \sum_{s=t}^T f_s(\mathbf{x}_s, \mathbf{u}_s) \quad \text{unt. d. NB} \quad \begin{cases} \mathbf{x}_t = \mathbf{x} \text{ und für } s = t, \dots, T-1 : \\ \mathbf{x}_{s+1} = \mathbf{g}_s(\mathbf{x}_s, \mathbf{u}_s) \end{cases} \quad (2)$$

Das Bellman-Prinzip besagt gerade, dass man mit der Markov-Politik-Lösung $\mathbf{u}_t^*(\mathbf{x}), \dots, \mathbf{u}_T^*(\mathbf{x})$ des Problems (2) die optimalen Politiken des gesamten, vom Zeitpunkt 0 an betrachteten Problems (1), im restlichen Zeitraum t, \dots, T bekommt. Damit lässt sich zeigen:

Satz 10.1 (Bellman-Gleichung als notwend. und hinreich. Bed. für Optimalität)

a) Sofern die Optimalwertfunktionen $V_t(\mathbf{x})$ aus (2) existieren (d.h. $V_t(\mathbf{x}) < \infty \forall t, \mathbf{x}$), erfüllt jede Markov-Politik-Lösung $\mathbf{u}_t^*(\mathbf{x})$ ($t = 0, \dots, T$) des Optimierungsproblems (1) die Bedingungen

$$\begin{aligned} \mathbf{u}_T^*(\mathbf{x}) &= \arg \max_{\mathbf{u} \in \mathcal{U}_T(\mathbf{x})} f_T(\mathbf{x}, \mathbf{u}) \\ \mathbf{u}_t^*(\mathbf{x}) &= \arg \max_{\mathbf{u} \in \mathcal{U}_t(\mathbf{x})} \{f_t(\mathbf{x}, \mathbf{u}) + V_{t+1}(\mathbf{g}_t(\mathbf{x}, \mathbf{u}))\}, \quad t = T-1, T-2, \dots, 0 \end{aligned} \quad (\text{ARG})$$

und die Optimalwertfunktionen $V_t(\mathbf{x})$ erfüllen die Bellman-Gleichung

$$\begin{aligned} V_T(\mathbf{x}) &= \max_{\mathbf{u} \in \mathcal{U}_T(\mathbf{x})} f_T(\mathbf{x}, \mathbf{u}) \\ V_t(\mathbf{x}) &= \max_{\mathbf{u} \in \mathcal{U}_t(\mathbf{x})} \{f_t(\mathbf{x}, \mathbf{u}) + V_{t+1}(\mathbf{g}_t(\mathbf{x}, \mathbf{u}))\}, \quad t = T-1, T-2, \dots, 0 \end{aligned} \quad (\text{BGL})$$

b) Wenn Funktionen $V_t(\mathbf{x})$ ($< \infty \forall t, \mathbf{x}$) existieren, die die Bellman-Gleichungen (BGL) erfüllen, und durch (ARG) eine Markov-Politik definiert wird, dann ist diese Politik eine Lösung des Optimierungsproblems (1).

Die Bellman-Gl. (BGL) reduziert das dynamische Optimierungsproblem auf eine **Sequenz statischer Optimierungsprobleme**, die in zeitlich negativer Richtung zu durchlaufen ist.

Beispiel 1:

$$\max_{u_t \in \mathbb{R}} \sum_{t=0}^T (x_t - u_t^2) \quad \text{unt. d. NB} \quad \begin{cases} x_0 \text{ gegeben,} \\ x_{t+1} = x_t + u_t \quad \text{für } t = 0, \dots, T-1 \end{cases}$$

Hier ist $f_t(x, u) = x - u^2$ und $g_t(x, u) = x + u$ für $t = 0, \dots, T-1$ sowie $f_T(x, u) = x^2 - u^2$. Wir bestimmen zunächst $V_T(x)$, $V_{T-1}(x)$, $V_{T-2}(x)$ und $V_{T-3}(x)$. Dabei wird sich das Bildungsgesetz für $V_{T-n}(x)$ abzeichnen, das sich dann durch vollständige Induktion über n nachweisen lässt:

$V_T(x)$: Da $f_T(x, u) = x - u^2$, bestimmt sich $V_T(x)$ als

$$V_T(x) = \max_{u \in \mathbb{R}} \{f_T(x, u)\} \stackrel{(\text{hier})}{=} \max_{u \in \mathbb{R}} \{x - u^2\}$$

Das Maximum über $u \in \mathbb{R}$ von $h(u) := x - u^2$ (mit x als Parameter) wird offensichtlich in $u = 0$ angenommen. Also ist $u_T^*(x) = 0$ und

$$V_T(x) = x$$

$V_{T-1}(x)$: Gemäß Bellman-Gl. ergibt sich $V_{T-1}(x)$ als:

$$\begin{aligned} V_{T-1}(x) &= \max_{u \in \mathbb{R}} \{f_{T-1}(x, u) + V_T(g_{T-1}(x, u))\} \\ &\stackrel{(\text{hier})}{=} \max_{u \in \mathbb{R}} \{x - u^2 + x + u\} = \max_{u \in \mathbb{R}} \{2x - u^2 + u\} \end{aligned}$$

Das Maximum über $u \in \mathbb{R}$ der Parabel $h(u) := 2x - u^2 + u$ (mit x als Parameter) wird dort angenommen, wo $h'(u) = -2u + 1 = 0$ ist, also $u = \frac{1}{2}$ (und somit $u_{T-1}^*(x) = \frac{1}{2}$). Somit ist

$$V_{T-1}(x) = 2x - \left(\frac{1}{2}\right)^2 + \frac{1}{2} = 2x + \frac{1}{4}$$

$V_{T-2}(x)$: Gemäß Bellman-Gl. ergibt sich $V_{T-2}(x)$ als:

$$\begin{aligned} V_{T-2}(x) &= \max_{u \in \mathbb{R}} \{f_{T-2}(x, u) + V_{T-1}(g_{T-2}(x, u))\} \\ &\stackrel{(\text{hier})}{=} \max_{u \in \mathbb{R}} \{x - u^2 + 2(x + u) + \frac{1}{4}\} = \max_{u \in \mathbb{R}} \{3x - u^2 + 2u + \frac{1}{4}\} \end{aligned}$$

Das Maximum über $u \in \mathbb{R}$ der Parabel $h(u) := 3x - u^2 + 2u + \frac{1}{4}$ (mit x als Parameter) wird dort angenommen, wo $h'(u) = -2u + 2 = 0$ ist, also $u = 1$ (und somit $u_{T-2}^*(x) = 1$). Also ist

$$V_{T-2}(x) = 3x - 1^2 + 2 \cdot 1 + \frac{1}{4} = 3x + \frac{5}{4}$$

$V_{T-3}(x)$: Gemäß Bellman-Gl. ergibt sich $V_{T-3}(x)$ als:

$$\begin{aligned} V_{T-3}(x) &= \max_{u \in \mathbb{R}} \{f_{T-3}(x, u) + V_{T-2}(g_{T-3}(x, u))\} \\ &\stackrel{(\text{hier})}{=} \max_{u \in \mathbb{R}} \{x - u^2 + 3(x + u) + \frac{5}{4}\} = \max_{u \in \mathbb{R}} \{4x - u^2 + 3u + \frac{5}{4}\} \end{aligned}$$

Das Maximum über $u \in \mathbb{R}$ der Parabel $h(u) := 4x - u^2 + 3u + \frac{5}{4}$ (mit x als Parameter) wird dort angenommen, wo $h'(u) = -2u + 3 = 0$ ist, also $u = \frac{3}{2}$ (und somit $u_{T-3}^*(x) = \frac{3}{2}$). Also ist

$$V_{T-3}(x) = 4x - \left(\frac{3}{2}\right)^2 + 3 \cdot \frac{3}{2} + \frac{5}{4} = 4x - \frac{9}{4} + \frac{9}{2} + \frac{5}{4} = 4x + \frac{9}{4} + \frac{5}{4} = 4x + \frac{14}{4}$$