

Einführung MDP - Policy Iteration

Arne Huckemann

16.11.2022

1 Motivation und Grundmodell

2 Erweiterung zur Policy

3 Optimalität

4 Policy Iteration

Motivation

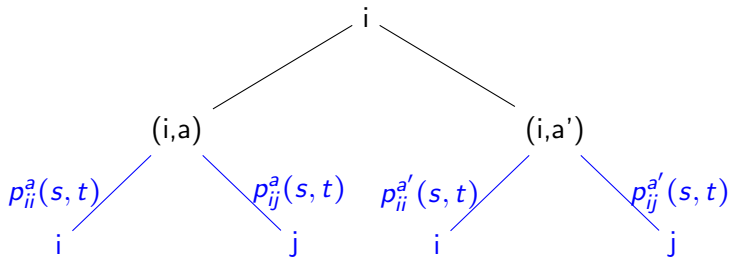
- Markov Decision Processes (MDP) aus 50/60er;
Howard, Ronald A. (1960)
- Möglichst gute Entscheidungen finden

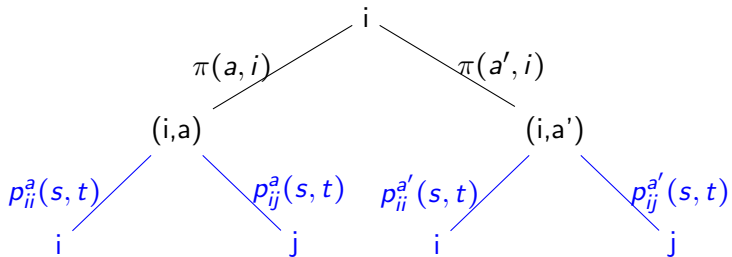
Motivation

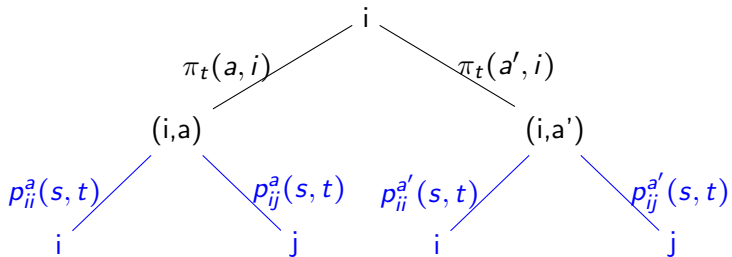
- Markov Decision Processes (MDP) aus 50/60er;
Howard, Ronald A. (1960)
- Möglichst gute Entscheidungen finden
- Anwendung
 - Robotik
 - **Epidemiologie**

Motivation

- Markov Decision Processes (MDP) aus 50/60er;
Howard, Ronald A. (1960)
- Möglichst gute Entscheidungen finden
- Anwendung
→ Robotik
→ **Epidemiologie**
- Wie findet man die optimale Lösung?
- Guo, X.; Hernández-Lerma, O. (2009) Continuous-Time Markov Decision Processes. Stochastic Modelling and Applied Probability, Springer







Annahme 1.1

- Zustands-Menge/Raum: \mathcal{Z} abzählbar, mit \mathbb{N} indiziert;
 $(\mathcal{Z}, \mathcal{P}(\mathcal{Z}))$

Annahme 1.1

- Zustands-Menge/Raum: \mathcal{Z} abzählbar, mit \mathbb{N} indiziert;
 $(\mathcal{Z}, \mathcal{P}(\mathcal{Z}))$
- Zeit: Kontinuierliche Zeit: $t \in [0, \infty) =: T$
- W-Raum: $(\Omega, \mathcal{F}, \mathbb{P})$

Annahme 1.1

- Zustands-Menge/Raum: \mathcal{Z} abzählbar, mit \mathbb{N} indiziert;
 $(\mathcal{Z}, \mathcal{P}(\mathcal{Z}))$
- Zeit: Kontinuierliche Zeit: $t \in [0, \infty) =: T$
- W-Raum: $(\Omega, \mathcal{F}, \mathbb{P})$

Definition 1.2 (Stochastischer Prozess)

- Familie von Zufallsvariablen
- $x: \Omega \times T \rightarrow \mathcal{Z}, (\omega, t) \mapsto x(\omega, t) = j$
 $\rightarrow \mathcal{F} - \mathcal{P}(\mathcal{Z})$ messbar

Annahme 1.1

- Zustands-Menge/Raum: \mathcal{Z} abzählbar, mit \mathbb{N} indiziert;
 $(\mathcal{Z}, \mathcal{P}(\mathcal{Z}))$
- Zeit: Kontinuierliche Zeit: $t \in [0, \infty) =: T$
- W-Raum: $(\Omega, \mathcal{F}, \mathbb{P})$

Definition 1.2 (Stochastischer Prozess)

- Familie von Zufallsvariablen
- $x: \Omega \times T \rightarrow \mathcal{Z}, (\omega, t) \mapsto x(\omega, t) = j$
 $\rightarrow \mathcal{F} - \mathcal{P}(\mathcal{Z})$ messbar
- $x: T \rightarrow \mathcal{Z}, t \mapsto x(t) = j, (\text{lassen } \omega \text{ weg})$

Definition 1.3 (Markov Eigenschaft)

Ein stochastischer Prozess heißt Markov Prozess, falls

$\forall n \in \mathbb{N} : 0 \leq s_1 \leq \dots \leq s_n \leq s \leq t < \infty; i_1, \dots, i_n, i, j \in \mathcal{Z} :$

$\mathbb{P}(x(t) = j \mid x(s) = i, x(s_n) = i_n, \dots, x(s_1) = i_1) = \mathbb{P}(x(t) = j \mid x(s) = i)$

Definition 1.3 (Markov Eigenschaft)

Ein stochastischer Prozess heißt Markov Prozess, falls

$\forall n \in \mathbb{N} : 0 \leq s_1 \leq \dots \leq s_n \leq s \leq t < \infty; i_1, \dots, i_n, i, j \in \mathcal{Z} :$

$\mathbb{P}(x(t) = j \mid x(s) = i, x(s_n) = i_n, \dots, x(s_1) = i_1) = \mathbb{P}(x(t) = j \mid x(s) = i)$

■ Übergangswahrscheinlichkeit:

$$p_{ij}(s, t) := \mathbb{P}(x(t) = j \mid x(s) = i) \quad \forall s \leq t \text{ und } i, j \in \mathcal{Z}$$

Definition 1.3 (Markov Eigenschaft)

Ein stochastischer Prozess heißt Markov Prozess, falls

$$\forall n \in \mathbb{N} : 0 \leq s_1 \leq \dots \leq s_n \leq s \leq t < \infty; i_1, \dots, i_n, i, j \in \mathcal{Z} :$$

$$\mathbb{P}(x(t) = j \mid x(s) = i, x(s_n) = i_n, \dots, x(s_1) = i_1) = \mathbb{P}(x(t) = j \mid x(s) = i)$$

- Übergangswahrscheinlichkeit:

$$p_{ij}(s, t) := \mathbb{P}(x(t) = j \mid x(s) = i) \quad \forall s \leq t \text{ und } i, j \in \mathcal{Z}$$

Definition 1.4

- Stabil, wenn

$$\lim_{\Delta \rightarrow 0^+} p_{ij}(t, t + \Delta) = \delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

Bemerkung 1.5

$$\sum_{j \in \mathcal{Z}} p_{ij}(s, t) = 1 \quad \forall s \leq t \text{ und } i, j \in \mathcal{Z}$$

Bemerkung 1.5

$$\sum_{j \in \mathcal{Z}} p_{ij}(s, t) = 1 \quad \forall s \leq t \text{ und } i, j \in \mathcal{Z}$$

Definition 1.6

Übergangsmatrix: $P(s, t) := (p_{ij}(s, t))_{(i,j) \in \mathcal{Z} \times \mathcal{Z}} \quad \forall s \leq t < \infty$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \forall s \leq u \leq t$$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \quad \forall s \leq u \leq t$$

Beweis:

$$\mathbb{P}(x(t) = j \mid x(s) = i) = \mathbb{P}(\cup_{k \in \mathcal{Z}} \{x(t) = j, x(u) = k\} \mid x(s) = i)$$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \quad \forall s \leq u \leq t$$

Beweis:

$$\mathbb{P}(x(t) = j \mid x(s) = i) = \mathbb{P}(\cup_{k \in \mathcal{Z}} \{x(t) = j, x(u) = k\} \mid x(s) = i)$$

$$\stackrel{\sigma \text{Add.}}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j, x(u) = k \mid x(s) = i) \quad \mathbb{P}(A, B \mid C) = \mathbb{P}(A \mid B, C)\mathbb{P}(B \mid C)$$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \quad \forall s \leq u \leq t$$

Beweis:

$$\mathbb{P}(x(t) = j \mid x(s) = i) = \mathbb{P}(\cup_{k \in \mathcal{Z}} \{x(t) = j, x(u) = k\} \mid x(s) = i)$$

$$\stackrel{\sigma Add.}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j, x(u) = k \mid x(s) = i) \quad \mathbb{P}(A, B \mid C) = \mathbb{P}(A \mid B, C)\mathbb{P}(B \mid C)$$

$$\stackrel{Bed. WK}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j \mid x(u) = k, x(s) = i) \mathbb{P}(x(u) = k \mid x(s) = i)$$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \quad \forall s \leq u \leq t$$

Beweis:

$$\mathbb{P}(x(t) = j \mid x(s) = i) = \mathbb{P}(\cup_{k \in \mathcal{Z}} \{x(t) = j, x(u) = k\} \mid x(s) = i)$$

$$\stackrel{\sigma Add.}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j, x(u) = k \mid x(s) = i) \quad \mathbb{P}(A, B \mid C) = \mathbb{P}(A \mid B, C)\mathbb{P}(B \mid C)$$

$$\stackrel{Bed. WK}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j \mid x(u) = k, x(s) = i) \mathbb{P}(x(u) = k \mid x(s) = i)$$

$$\stackrel{Markov}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j \mid x(u) = k) \mathbb{P}(x(u) = k \mid x(s) = i)$$

Satz 1.7 (Chapman-Kolmogorov Equation)

$$P(s, t) = P(s, u)P(u, t), \quad \forall s \leq u \leq t$$

Beweis:

$$\mathbb{P}(x(t) = j \mid x(s) = i) = \mathbb{P}(\cup_{k \in \mathcal{Z}} \{x(t) = j, x(u) = k\} \mid x(s) = i)$$

$$\stackrel{\sigma Add.}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j, x(u) = k \mid x(s) = i) \quad \mathbb{P}(A, B \mid C) = \mathbb{P}(A \mid B, C)\mathbb{P}(B \mid C)$$

$$\stackrel{Bed. WK}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j \mid x(u) = k, x(s) = i) \mathbb{P}(x(u) = k \mid x(s) = i)$$

$$\stackrel{Markov}{=} \sum_{k \in \mathcal{Z}} \mathbb{P}(x(t) = j \mid x(u) = k) \mathbb{P}(x(u) = k \mid x(s) = i)$$

$$\stackrel{Def.}{=} \sum_{k \in \mathcal{Z}} p_{kj}(u, t) p_{ik}(s, u)$$

Satz und Definition 1.8

Für einen stabilen Markov Prozess gilt $\forall 0 \leq s$:

$$\infty > \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - \delta_{ij}}{t - s} =: q_{ij}(s)$$

Satz und Definition 1.8

Für einen stabilen Markov Prozess gilt $\forall 0 \leq s$:

$$\infty > \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - \delta_{ij}}{t - s} =: q_{ij}(s)$$

Beweis: Nur für Fall $i = j$:

Annahme: $p_{ij}(s, t) \neq 1$

Satz und Definition 1.8

Für einen stabilen Markov Prozess gilt $\forall 0 \leq s$:

$$\infty > \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - \delta_{ij}}{t - s} =: q_{ij}(s)$$

Beweis: Nur für Fall $i = j$:

Annahme: $p_{ii}(s, t) \neq 1$

1) $f(s, t) := -\log(p_{ii}(s, t))$

Satz und Definition 1.8

Für einen stabilen Markov Prozess gilt $\forall 0 \leq s$:

$$\infty > \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - \delta_{ij}}{t - s} =: q_{ij}(s)$$

Beweis: Nur für Fall $i = j$:

Annahme: $p_{ii}(s, t) \neq 1$

$$1) f(s, t) := -\log(p_{ii}(s, t))$$

$$p_{ii}(s, t) = \sum_{k \in \mathcal{Z}} p_{ik}(s, u) p_{ki}(u, t) \geq p_{ii}(s, u) p_{ii}(u, t)$$

Satz und Definition 1.8

Für einen stabilen Markov Prozess gilt $\forall 0 \leq s$:

$$\infty > \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - \delta_{ij}}{t - s} =: q_{ij}(s)$$

Beweis: Nur für Fall $i = j$:

Annahme: $p_{ii}(s, t) \neq 1$

$$1) f(s, t) := -\log(p_{ii}(s, t))$$

$$p_{ii}(s, t) = \sum_{k \in \mathcal{Z}} p_{ik}(s, u) p_{ki}(u, t) \geq p_{ii}(s, u) p_{ii}(u, t)$$

$$\Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$1) f(s, t) := -\log(p_{ij}(s, t))$$

$$2) p_{ij}(s, t) \geq p_{ij}(s, u)p_{ij}(u, t) \Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$1) f(s, t) := -\log(p_{ij}(s, t))$$

$$2) p_{ij}(s, t) \geq p_{ij}(s, u)p_{ij}(u, t) \Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$\Rightarrow^{[2]} \sup_{t \rightarrow s^+} \frac{f(s, t)}{t-s} = \lim_{t \rightarrow s^+} \frac{f(s, t)}{t-s}$$

$$1) f(s, t) := -\log(p_{ij}(s, t))$$

$$2) p_{ij}(s, t) \geq p_{ij}(s, u)p_{ij}(u, t) \Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$\Rightarrow^{[2]} \sup_{t \rightarrow s^+} \frac{f(s, t)}{t-s} = \lim_{t \rightarrow s^+} \frac{f(s, t)}{t-s}$$

Damit Existenz gezeigt. Nun Eindeutigkeit:

$$\lim_{t \rightarrow s^+} \frac{p_{ij}(s, t) - 1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)} - 1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)} - 1}{f(s, t)} \frac{f(s, t)}{t-s}$$

$$1) f(s, t) := -\log(p_{ij}(s, t))$$

$$2) p_{ij}(s, t) \geq p_{ij}(s, u)p_{ij}(u, t) \Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$\Rightarrow^{[2]} \sup_{t \rightarrow s^+} \frac{f(s, t)}{t-s} = \lim_{t \rightarrow s^+} \frac{f(s, t)}{t-s}$$

Damit Existenz gezeigt. Nun Eindeutigkeit:

$$\lim_{t \rightarrow s^+} \frac{p_{ij}(s, t)-1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)}-1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)}-1}{f(s, t)} \frac{f(s, t)}{t-s}$$

$$\text{Da } \frac{e^{-f(s, t)}-1}{f(s, t)} = \frac{(1-f(s, t)+f(s, t)^2/2-\dots)-1}{f(s, t)} = \frac{f(s, t)}{f(s, t)} \underbrace{(-1 + f(s, t)/2 + \dots)}_{\rightarrow 0, t \rightarrow s^+}$$

$$1) f(s, t) := -\log(p_{ij}(s, t))$$

$$2) p_{ij}(s, t) \geq p_{ij}(s, u)p_{ij}(u, t) \Rightarrow f(s, t) \leq f(s, u) + f(u, t)$$

$$\Rightarrow^{[2]} \sup_{t \rightarrow s^+} \frac{f(s, t)}{t-s} = \lim_{t \rightarrow s^+} \frac{f(s, t)}{t-s}$$

Damit Existenz gezeigt. Nun Eindeutigkeit:

$$\lim_{t \rightarrow s^+} \frac{p_{ij}(s, t)-1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)}-1}{t-s} = \lim_{t \rightarrow s^+} \frac{e^{-f(s, t)}-1}{f(s, t)} \frac{f(s, t)}{t-s}$$

$$\text{Da } \frac{e^{-f(s, t)}-1}{f(s, t)} = \frac{(1-f(s, t)+f(s, t)^2/2-\dots)-1}{f(s, t)} = \frac{f(s, t)}{f(s, t)} \underbrace{(-1 + f(s, t)/2 + \dots)}_{\rightarrow 0, t \rightarrow s^+}$$

$$\Rightarrow -\lim_{t \rightarrow s^+} \frac{f(s, t)}{t-s} = -\sup_{t \rightarrow s^+} \frac{f(s, t)}{t-s} =: q_{ij}(s)$$

Korollar 1.9

Außerdem gilt für die Transition Rates:

- i) Nicht-Homogen: $q_{ii}(s) \leq 0$ und $q_{ij}(s) \geq 0$
- ii) Konservativ: $\sum_{j \in \mathcal{Z}} q_{ij}(s) = 0$

Korollar 1.9

Außerdem gilt für die Transition Rates:

- i) Nicht-Homogen: $q_{ii}(s) \leq 0$ und $q_{ij}(s) \geq 0$
- ii) Konservativ: $\sum_{j \in \mathcal{Z}} q_{ij}(s) = 0$

Beweis:

- i) $p_{ii}(s, t) - 1 \leq 0$ und $p_{ij}(s, t) \geq 0$

Korollar 1.9

Außerdem gilt für die Transition Rates:

- i) Nicht-Homogen: $q_{ii}(s) \leq 0$ und $q_{ij}(s) \geq 0$
- ii) Konservativ: $\sum_{j \in \mathcal{Z}} q_{ij}(s) = 0$

Beweis:

- i) $p_{ii}(s, t) - 1 \leq 0$ und $p_{ij}(s, t) \geq 0$
- ii) $p_{ii}(s, t) = 1 - \sum_{i \neq j} p_{ij}(s, t)$

$$\begin{aligned}\Rightarrow q_{ii}(s) &:= \lim_{t \rightarrow s^+} \frac{p_{ii}(s, t) - 1}{t - s} = \lim_{t \rightarrow s^+} \frac{1 - \sum_{i \neq j} p_{ij}(s, t) - 1}{t - s} \\ &= - \sum_{i \neq j} \lim_{t \rightarrow s^+} \frac{p_{ij}(s, t)}{t - s} = - \sum_{i \neq j} q_{ij}(s, t)\end{aligned}$$

Satz 1.10 (Umkehrung)

$t \geq 0 : Q(t) := (q_{ij}(t))_{i,j \in \mathcal{Z}}$ mit nicht-homogenen, konservativen, messbaren und integrierbaren Einträgen auf \mathbb{R} .

Falls für $L_1, L_2 > 0$ und $w \in [1, \infty)^{\mathcal{Z}}$:

$$w^T Q(t) \leq L_1 w^T \text{ und } -\text{diag}(Q(t)) \leq L_2 w$$

dann gibt es ein eindeutiges $P(s, t)$ das durch:

■ $q_{ij}(t) := \lim_{\Delta \rightarrow 0^+} \frac{p_{ij}(t, t+\Delta) - \delta_{ij}}{\Delta}$ bestimmt ist und es gilt

Satz 1.10 (Umkehrung)

$t \geq 0$: $Q(t) := (q_{ij}(t))_{i,j \in \mathcal{Z}}$ mit nicht-homogenen, konservativen, messbaren und integrierbaren Einträgen auf \mathbb{R} .

Falls für $L_1, L_2 > 0$ und $w \in [1, \infty)^{\mathcal{Z}}$:

$$w^T Q(t) \leq L_1 w^T \text{ und } -\text{diag}(Q(t)) \leq L_2 w$$

dann gibt es ein eindeutiges $P(s, t)$ das durch:

- $q_{ij}(t) := \lim_{\Delta \rightarrow 0^+} \frac{p_{ij}(t, t+\Delta) - \delta_{ij}}{\Delta}$ bestimmt ist und es gilt
- Kolmogorov's Backward und Forward:
 - $\frac{\partial}{\partial t} P(s, t) = P(s, t) Q(t)$
 - $\frac{\partial}{\partial s} P(s, t) = -Q(s) P(s, t)$

Beispiel Epidemiologie (1/4)

- $\mathcal{Z} := \mathbb{N}_0$
- Population = $N \in \mathbb{N}$

Beispiel Epidemiologie (1/4)

- $\mathcal{Z} := \mathbb{N}_0$
- Population = $N \in \mathbb{N}$

$$q_{ij}(s) = \begin{cases} \lambda_i + \hat{\lambda}_i & j = i + 1 \\ -(\lambda_i + \hat{\lambda}_i + \mu_i) & j = i \\ \mu_i & j = i - 1 \\ 0 & \text{sonst} \end{cases}$$

- λ_i Ansteckungsrate innerhalb der Population
- $\hat{\lambda}_i$ Ansteckungsrate außerhalb der Population
- μ_i Genesungsrate

Beispiel Epidemiologie 2/4 (Überprüfen der Annahmen)

$$a_i := \lambda_i + \hat{\lambda}_i \text{ und } b_i := \mu_i$$

$$w^T Q(t) \leq L_1 w^T$$

$$Q(t) = \begin{pmatrix} -(a_1 + b_1) & a_1 & 0 & 0 & \dots \\ b_2 & -(a_2 + b_2) & a_2 & 0 & \dots \\ 0 & b_3 & -(a_3 + b_3) & a_3 & \dots \\ 0 & 0 & b_4 & -(a_4 + b_4) & \dots \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

Beispiel Epidemiologie 2/4 (Überprüfen der Annahmen)

$$a_i := \lambda_i + \hat{\lambda}_i \text{ und } b_i := \mu_i$$

$$w^T Q(t) \leq L_1 w^T$$

$$Q(t) = \begin{pmatrix} -(a_1 + b_1) & a_1 & 0 & 0 & \dots \\ b_2 & -(a_2 + b_2) & a_2 & 0 & \dots \\ 0 & b_3 & -(a_3 + b_3) & a_3 & \dots \\ 0 & 0 & b_4 & -(a_4 + b_4) & \dots \\ \dots & \dots & \dots & \dots & \dots \end{pmatrix}$$

$$w = (1, 1, \dots)^T$$
$$L_1 := \|w^T Q(t)\|_\infty$$

Beispiel Epidemiologie 3/4 (Überprüfen der Annahmen)

$$-diag(Q(t)) \leq L_2 w$$

$$-diag(Q(t)) = \begin{pmatrix} a_1 + b_1 \\ a_2 + b_2 \\ \dots \end{pmatrix} \leq L_2 \begin{pmatrix} w_1 \\ w_2 \\ \dots \end{pmatrix}$$

Beispiel Epidemiologie 3/4 (Überprüfen der Annahmen)

$$-diag(Q(t)) \leq L_2 w$$

$$-diag(Q(t)) = \begin{pmatrix} a_1 + b_1 \\ a_2 + b_2 \\ \dots \end{pmatrix} \leq L_2 \begin{pmatrix} w_1 \\ w_2 \\ \dots \end{pmatrix}$$

Und $L_2 = \|diag(Q(t))\|_\infty$

Definition 2.1

- $(\mathcal{A}, \sigma(\mathcal{A}))$ Action Space
- $\mathcal{A}(i) \in \sigma(\mathcal{A})$ mögliche Actions in $i \in \mathcal{Z}$

Definition 2.1

- $(\mathcal{A}, \sigma(\mathcal{A}))$ Action Space
- $\mathcal{A}(i) \in \sigma(\mathcal{A})$ mögliche Actions in $i \in \mathcal{Z}$

Definition 2.2 (Markov Kern)

Für jedes i ist $\pi(\cdot, i)$ Markovkern:

- $\pi(\cdot, i) : [0, \infty) \times \sigma(\mathcal{A}(i)) \rightarrow [0, 1]$

Definition 2.1

- $(\mathcal{A}, \sigma(\mathcal{A}))$ Action Space
- $\mathcal{A}(i) \in \sigma(\mathcal{A})$ mögliche Actions in $i \in \mathcal{Z}$

Definition 2.2 (Markov Kern)

Für jedes i ist $\pi(\cdot, i)$ Markovkern:

- $\pi(\cdot, i) : [0, \infty) \times \sigma(\mathcal{A}(i)) \rightarrow [0, 1]$
 1. $\pi(A, i)$ ist $\mathcal{B}([0, \infty)) - \mathcal{B}([0, 1])$ messbar $\forall A \in \sigma(\mathcal{A}(i)), i \in \mathcal{Z}$

Definition 2.1

- $(\mathcal{A}, \sigma(\mathcal{A}))$ Action Space
- $\mathcal{A}(i) \in \sigma(\mathcal{A})$ mögliche Actions in $i \in \mathcal{Z}$

Definition 2.2 (Markov Kern)

Für jedes i ist $\pi(\cdot, i)$ Markovkern:

- $\pi(\cdot, i) : [0, \infty) \times \sigma(\mathcal{A}(i)) \rightarrow [0, 1]$
 1. $\pi(A, i)$ ist $\mathcal{B}([0, \infty)) - \mathcal{B}([0, 1])$ messbar $\forall A \in \sigma(\mathcal{A}(i)), i \in \mathcal{Z}$
 2. $\pi_t(\cdot, i)$ ist W-Maß auf $\sigma(\mathcal{A}(i)) \forall t \in [0, \infty)$ und $i \in \mathcal{Z}$

Definition 2.1

- $(\mathcal{A}, \sigma(\mathcal{A}))$ Action Space
- $\mathcal{A}(i) \in \sigma(\mathcal{A})$ mögliche Actions in $i \in \mathcal{Z}$

Definition 2.2 (Markov Kern)

Für jedes i ist $\pi(\cdot, i)$ Markovkern:

- $\pi(\cdot, i) : [0, \infty) \times \sigma(\mathcal{A}(i)) \rightarrow [0, 1]$
 1. $\pi(A, i)$ ist $\mathcal{B}([0, \infty)) - \mathcal{B}([0, 1])$ messbar $\forall A \in \sigma(\mathcal{A}(i)), i \in \mathcal{Z}$
 2. $\pi_t(\cdot, i)$ ist W-Maß auf $\sigma(\mathcal{A}(i)) \forall t \in [0, \infty)$ und $i \in \mathcal{Z}$

Policy: $\pi := (\pi_t)_{t \geq 0} \in \Pi$

Definition 2.3

- Gegebene Transition Rates: $q(j \mid i, a)$ messbar und integrierbar in $a \in \mathcal{A}(i)$
- Durchschnittliche Transition Rate von i nach j der Policy π :

$$q_{ij}^{\pi}(t) := \int_{\mathcal{A}(i)} q(j \mid i, a) d\pi_t(a, i)$$

Definition 2.3

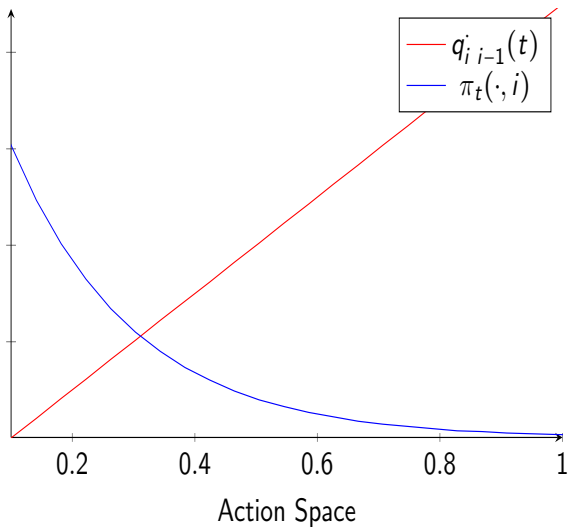
- Gegebene Transition Rates: $q(j \mid i, a)$ messbar und integrierbar in $a \in \mathcal{A}(i)$
- Durchschnittliche Transition Rate von i nach j der Policy π :

$$q_{ij}^{\pi}(t) := \int_{\mathcal{A}(i)} q(j \mid i, a) d\pi_t(a, i)$$

Definition 2.4

- Kosten: $c(i \mid a) \geq 0$ messbar und integrierbar
- Durchschnittliche Kosten der Policy π im Zustand i :

$$c_i^{\pi}(t) := \int_{\mathcal{A}(i)} c(i \mid a) d\pi_t(a, i)$$



Beispiel Epidemiologie (4/4)

- $A := [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}] = A(i) \quad \forall i \in \mathcal{Z}$
- $\mathbf{a} = (a_1, a_2) \in A$

Beispiel Epidemiologie (4/4)

- $A := [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}] = A(i) \quad \forall i \in \mathcal{Z}$

- $\mathbf{a} = (a_1, a_2) \in A$

→ $a_1 \in [\underline{a}, \bar{a}]$ Level der Quarantäne $\Rightarrow \hat{\lambda}_i := \hat{\lambda}_i(a_1)$

Beispiel Epidemiologie (4/4)

- $A := [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}] = A(i) \quad \forall i \in \mathcal{Z}$

- $\mathbf{a} = (a_1, a_2) \in A$

→ $a_1 \in [\underline{a}, \bar{a}]$ Level der Quarantäne $\Rightarrow \hat{\lambda}_i := \hat{\lambda}_i(a_1)$

→ $a_2 \in [\underline{b}, \bar{b}]$ Level der medizinischen Behandlung $\Rightarrow \lambda_i := \lambda_i(a_2)$

Beispiel Epidemiologie (4/4)

- $A := [\underline{a}, \bar{a}] \times [\underline{b}, \bar{b}] = A(i) \quad \forall i \in \mathcal{Z}$

- $a = (a_1, a_2) \in A$

→ $a_1 \in [\underline{a}, \bar{a}]$ Level der Quarantäne $\Rightarrow \hat{\lambda}_i := \hat{\lambda}_i(a_1)$

→ $a_2 \in [\underline{b}, \bar{b}]$ Level der medizinischen Behandlung $\Rightarrow \lambda_i := \lambda_i(a_2)$

$$c(i \mid a) := h_0(i) + h_1(a_1) + h_2(i, a_2); \quad \forall i \in \mathcal{Z}, a = (a_1, a_2) \in A$$

Definition 3.1

Wert der Policy mit Start in i ($x(0) = i$) und $\alpha > 0$:

$$\begin{aligned} J_i^\pi &:= \mathbb{E}_\pi^i \left[\int_0^\infty e^{-\alpha t} c_{x(t)}^\pi(t) dt \right] \\ &= \int_0^\infty e^{-\alpha t} \mathbb{E}_\pi^i \left[c_{x(t)}^\pi(t) \right] dt \end{aligned}$$

$$\blacksquare \mathbb{E}_\pi^i \left[c_{x(t)}^\pi(t) \right] = \sum_{j \in \mathcal{Z}} c_j^\pi(t) p_{ij}(0, t)$$

Definition 3.2

π^* heißt optimale Policy, falls

$$\pi^* \in \operatorname{argmin}_{\pi \in \Pi} J_i^\pi$$

Definition 4.1

Sei $f: \mathcal{Z} \rightarrow \mathcal{A}$, $i \mapsto f(i) \in \mathcal{A}(i)$

Neue stationäre deterministische Policy:

$$\pi(A, i) = \delta_{f(i)}(A), \quad A \in \sigma(\mathcal{A}(i))$$

Definition 4.1

Sei $f: \mathcal{Z} \rightarrow \mathcal{A}$, $i \mapsto f(i) \in \mathcal{A}(i)$

Neue stationäre deterministische Policy:

$$\pi(A, i) = \delta_{f(i)}(A), \quad A \in \sigma(\mathcal{A}(i))$$

Definition 4.2

$\forall a \sim \pi(\cdot, i) :$

$$q_{ij}^{\pi}(t) = \int_{\mathcal{A}(i)} q(j \mid i, a) d\pi(a, i) =: q_{ij}^{f(i)}$$

Definition 4.1

Sei $f: \mathcal{Z} \rightarrow \mathcal{A}$, $i \mapsto f(i) \in \mathcal{A}(i)$

Neue stationäre deterministische Policy:

$$\pi(A, i) = \delta_{f(i)}(A), \quad A \in \sigma(\mathcal{A}(i))$$

Definition 4.2

$\forall a \sim \pi(\cdot, i) :$

$$q_{ij}^{\pi}(t) = \int_{\mathcal{A}(i)} q(j \mid i, a) d\pi(a, i) =: q_{ij}^{f(i)}$$

$$c_i^{\pi}(t) = \int_{\mathcal{A}(i)} c(i \mid a) d\pi_t(a, i) =: c_i^{f(i)}$$

Definition 4.1

Sei $f: \mathcal{Z} \rightarrow \mathcal{A}$, $i \mapsto f(i) \in \mathcal{A}(i)$

Neue stationäre deterministische Policy:

$$\pi(A, i) = \delta_{f(i)}(A), \quad A \in \sigma(\mathcal{A}(i))$$

Definition 4.2

$\forall a \sim \pi(\cdot, i) :$

$$q_{ij}^{\pi}(t) = \int_{\mathcal{A}(i)} q(j \mid i, a) d\pi(a, i) =: q_{ij}^{f(i)}$$

$$c_i^{\pi}(t) = \int_{\mathcal{A}(i)} c(i \mid a) d\pi_t(a, i) =: c_i^{f(i)}$$

Transition Matrix $P(s, t)$ existiert und bestimmt $x(t)$ auf \mathcal{Z}

Satz 4.3

$$\begin{aligned} J_i^f &= \int_0^\infty e^{-\alpha t} \mathbb{E}_\pi^i \left[c_{x(t)}^{f(i)} \right] dt = \int_0^\infty e^{-\alpha t} \sum_{j \in \mathcal{Z}} c_j^{f(i)} p_{ij}(0, t) dt \\ &= \frac{1}{\alpha - q_{ii}^{f(i)}} \left(c_i^{f(i)} + \sum_{i \neq k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)} \right) \end{aligned}$$

Satz 4.3

$$\begin{aligned} J_i^f &= \int_0^\infty e^{-\alpha t} \mathbb{E}_\pi^i \left[c_{x(t)}^{f(i)} \right] dt = \int_0^\infty e^{-\alpha t} \sum_{j \in \mathcal{Z}} c_j^{f(i)} p_{ij}(0, t) dt \\ &= \frac{1}{\alpha - q_{ii}^{f(i)}} \left(c_i^{f(i)} + \sum_{i \neq k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)} \right) \end{aligned}$$

$$\begin{aligned} &= \overbrace{\sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}} \\ \Leftrightarrow J_i^f &= \frac{1}{\alpha - q_{ii}^{f(i)}} \left(\sum_{i \neq k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)} + J_i^f q_{ii}^{f(i)} - J_i^f q_{ii}^{f(i)} \right) \end{aligned}$$

Satz 4.3

$$\begin{aligned} J_i^f &= \int_0^\infty e^{-\alpha t} \mathbb{E}_\pi^i \left[c_{x(t)}^{f(i)} \right] dt = \int_0^\infty e^{-\alpha t} \sum_{j \in \mathcal{Z}} c_j^{f(i)} p_{ij}(0, t) dt \\ &= \frac{1}{\alpha - q_{ii}^{f(i)}} \left(c_i^{f(i)} + \sum_{i \neq k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)} \right) \end{aligned}$$

$$\Leftrightarrow J_i^f = \frac{1}{\alpha - q_{ii}^{f(i)}} \left(\overbrace{\sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}} + J_i^f q_{ii}^{f(i)} - J_i^f q_{ii}^{f(i)} \right)$$

$$\Leftrightarrow (\alpha - q_{ii}^{f(i)}) J_i^f = c_i^{f(i)} + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)} - J_i^f q_{ii}^{f(i)}$$

$$\alpha J_i^f = c_i^{f(i)} + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}$$

$$\alpha J_i^f = c_i^{f(i)} + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}$$

Satz 4.4 (Policy Improvement)

■ $D_f(i, a) := c_i^a + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}$ für $i \in \mathcal{Z}$

$$\alpha J_i^f = c_i^{f(i)} + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}$$

Satz 4.4 (Policy Improvement)

- $D_f(i, a) := c_i^a + \sum_{k \in \mathcal{Z}} J_k^f q_{ik}^{f(i)}$ für $i \in \mathcal{Z}$

- f' Policy Improvement zu f , falls

$$\rightarrow f'(i) := \begin{cases} f(i) & , D_f(i, a) \geq \alpha J_i^f \quad \forall a \in A(i) \\ a' & , \text{irgend ein } a' \in A(i) \text{ mit } D_f(i, a') < \alpha J_i^f \end{cases}$$

$$\Rightarrow J_i^f \geq J_i^{f'}$$

Beweis: Direkt per Konstruktion

Policy Iteration Algorithmus

1. Wähle $k=0$ und zufällige/beliebige Policy $f^{(k=0)}$ und Genauigkeitsmaß $\epsilon > 0$

Policy Iteration Algorithmus

1. Wähle $k=0$ und zufällige/beliebige Policy $f^{(k=0)}$ und Genauigkeitsmaß $\epsilon > 0$
2. Iteriere:
 - i) Wert $J_i^{f^{(k)}}$ ermitteln
 - ii) Policy Improvement: $f^{(k+1)} = f'^{(k)}$
 - iii) Falls $|J_i^{f^{(k+1)}} - J_i^{f^{(k)}}| < \epsilon \ \forall i \in \mathcal{Z} \Rightarrow$ fertig.
Sonst zurück zu i) und setze $k = k + 1$

Policy Iteration Algorithmus

1. Wähle $k=0$ und zufällige/beliebige Policy $f^{(k=0)}$ und Genauigkeitsmaß $\epsilon > 0$
2. Iteriere:
 - i) Wert $J_i^{f^{(k)}}$ ermitteln
 - ii) Policy Improvement: $f^{(k+1)} = f'^{(k)}$
 - iii) Falls $|J_i^{f^{(k+1)}} - J_i^{f^{(k)}}| < \epsilon \ \forall i \in \mathcal{Z} \Rightarrow$ fertig.
Sonst zurück zu i) und setze $k = k + 1$

Satz 4.5

Dieser Algorithmus konvergiert zu einem lokalen Minimum

Literatur

- [1] Guo, X.; Hernandez-Lerma, O. (2009) Continuous-Time Markov Decision Processes. Stochastic Modelling and Applied Probability, Springer
- [2] Liuer Ye und Xianping Guo und Onésimo Hernández-Lerma, (2008) Existence and Regularity of a Nonhomogeneous Transition Matrix under Measurability Conditions, J Theor Probab
- [3] Kolmogoroff, A. (1930) Über die analytischen Methoden in der Wahrscheinlichkeitsrechnung. Math. Ann.
- [4] Howard, Ronald A. (1960) Dynamic Programming and Markov Processes, The M.I.T. Press
- [5] Marek Kuczma, (2009) An Introduction to the Theory of Functional Equations and Inequalities, Second Edition