Jules Hedges, Riu Rodríguez Sakamoto

Mathematically Structured Programming group, Department of Computer and Information Sciences, University of Strathclyde

riu.rodriguez-sakamoto@strath.ac.uk

5th International Conference on Applied Category Theory, Glasgow (Scotland)

July 21st, 2022

Riu Rguez Sakamoto

Dynamic Programming

Introduction

Decision Processes

Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Value Iteration

Q-learning

Related work

Conclusion

References

Towards Reinforcement Learning algorithms

- *n*-armed bandit problem: Exploration vs exploitation (greedy, *ε*-greedy)
- Contextual bandits: Many states, immediate rewards
- Reinforcement Learning problem: Many states, delayed rewards

Decision Processes

Riu Rguez Sakamoto

Dynamic Programmin

Introduction

Decision Processes

Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learning

Related wor

References

state space X action space A transition $f: X \times A \rightarrow X$ utility $U: X \times A \rightarrow \mathbb{R}$

Objective: Which actions to choose where, to optimize long-run reward?

policy $\pi: X \to A$

Long-run reward: Discounted¹ sum with infinite horizon

$$x_0 \xrightarrow{f(x_0,\pi(x_0))} x_1 \xrightarrow{f(x_1,\pi(x_1))} x_2 \longrightarrow \cdots$$

 $V_{\pi}(x_0) = \sum_{k=0}^{\infty} \beta^k U(x_k,\pi(x_k))$

¹Discount factor $\beta \in (0, 1)$

Dynamic Programming

Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processe

Dynamic Programming

Value iteration with optics

Optics

- Optics for Valu Iteration
- Q-learning
- Related work

Conclusion

References

Assign value functions to policies.

Start with some policy $\pi : X \to A$, and some value function $V : X \to \mathbb{R}$. Bellman update steps:

- Value improvement: V'(x) = U(x, π(x)) + βV(f(x, π(x))) (policy evaluation)
- Policy improvement: $\pi'(x) = \arg \max_{a \in A} U(x, a) + \beta V(f(x, a))$

Bellman optimality condition: Fixpoint of the update function.

- Value improvement: $V(x) = U(x, \pi(x)) + \beta V(f(x, \pi(x)))$
- Policy improvement: $\pi(x) = \arg \max_{a \in A} U(x, a) + \beta V(f(x, a))$

Riu Rguez Sakamoto

Dynamic Programmir

Introduction Decision Processes

Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Va Iteration

Q-learning

Related work

Conclusion

References

Dynamic Programming

Classical DP algorithms:

- Policy Iteration: ((V impr.)*; π impr.)*
- Value Iteration: (V impr.; π impr.)*
 - Fusion of both steps:

$$V(x) = \max_{a \in A} U(x, a) + \beta V(f(x, a))$$



Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processe Dynamic Programming

MDPs

Value iteration with optics

Optics for

O lassa'a

Deleteral

Related won

Conclusion

References

Markov Decision Processes

MDP: Some function is stochastic: Kleisli morphism of Δ (e.g. Gridworld with uncertainty):

state space X action space A transition $f : X \times A \rightarrow \Delta X$ utility $U : X \times A \rightarrow \Delta \mathbb{R}$

policy $\pi: X \to \Delta A$



Riu Rguez Sakamoto

Dynamic Programmi

Introduction Decision Processes Dynamic

Programi MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learnin

Related work

Conclusion

References

Continuous spaces w/o stochasticity

Continuous state or action space, e.g. savings model in economics:

- $X = \mathbb{R}$: Assets (savings)
- $A = \mathbb{R}$: Consumption
- $(1 + \gamma)$: Gross rate of return
- *i*: Interest (constant)

• Policy: $\pi : x \mapsto c$ • Transition: $f : (x, c) \mapsto (1 + \gamma)(x - c + i)$

• Utility: $U: (_, c) \mapsto U(c)$ (w/ Inada or transversality condition)

$$V'(x) = \max_{0 \le c \le x+i} U(c) + \beta V((1+\gamma)(x-c+i))$$

Also, the inverted pendulum control problem. Continuous spaces and stochastic maps: Stoch = Kl(Giry) or some subcategory.





Optics

Riu Rguez Sakamoto

Dynamic Programming

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

- Optics for Va Iteration Q-learning
- Related work
- Conclusion

References

$$\begin{split} \mathcal{M} \mbox{ monoidal category} \\ \downarrow \bullet : \mathcal{M} \times \mathcal{C} \to \mathcal{C} \mbox{ action of } \mathcal{M} \mbox{ on } \mathcal{C} \\ \mathcal{C} \mbox{ category} \end{split}$$

Then C is a \mathcal{M} -actegory. Given \mathcal{M} monoidal, and two \mathcal{M} -actegories C, \mathcal{D} , the category of mixed optics $\text{Optic}_{\mathcal{C},\mathcal{D}}$ has

objects
$$\begin{pmatrix} X \\ X' \end{pmatrix} \in \mathcal{C}$$

• morphisms are coends

$$\operatorname{Optic}_{\mathcal{C},\mathcal{D}}\left(\binom{X}{X'},\binom{Y}{Y'}\right) = \int^{M:\mathcal{M}} \underbrace{\mathcal{C}(X, M \bullet Y)}_{f} \times \underbrace{\mathcal{D}(M \bullet Y', X')}_{f'}$$

i.e. equivalence classes of (M, f, f').

Optics

Riu Rguez Sakamoto

Dynamic Programming

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learning

Related wor

Conclusion

References

Let $\mathcal{C} = \mathcal{M} = \text{Mark} = \text{Kl}(\Delta)$, and $\mathcal{D} = (\text{Conv}, \boxtimes, \{*\}) = \text{EM}(\Delta)$, where \boxtimes is similar to tensor product in *R*-Mod. Conv is monoidal closed under \boxtimes , where the function set [X, Y] has pointwise convex structure.

$$\mathsf{Optic}_{\mathsf{Mark},\mathsf{Conv}}\left(\binom{X}{X'},\binom{Y}{Y'}\right) = \int^{M:\mathsf{Mark}} \mathsf{Mark}(X, M \otimes Y) \times \mathsf{Conv}(\Delta M, [Y', X'])$$

String diagram syntax (informal for mixed case):







Optics for Value Iteration



Riu Rguez Sakamoto

Dynamic Programming

Introduction Decision Proces

Programming MDPs

Value iteration with optics

Optics

Optics for Value Iteration

Q-learning

Related work

Conclusion

References

simulate: look into the future Х Х update our values \mathbb{R} \mathbb{R} (a) Lenses are iteration steps Х f π Х

U

VB

πъ

 \mathbb{R}

Riu Rguez

Optics for Value Iteration

Sakamoto Dynamic Programmin

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Value Iteration

Q-learning

Related work

References



• V improv: $(\pi, V) \mapsto (\pi, \overline{\pi} \ ; \lambda \ V)$



• π improv: $(\pi, V) \mapsto (x \mapsto \arg \max_{a \in A} (\lambda \circ V)(x, a), V)$

Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Value Iteration

Q-learning Related work

References

Optics for Value Iteration Mixed optics for MDPs

$$\lambda: \begin{pmatrix} \mathbf{X} \times \mathbf{A} \\ \mathbb{R} \end{pmatrix} \to \begin{pmatrix} \mathbf{X} \\ \mathbb{R} \end{pmatrix}$$

where
$$\lambda = (X \otimes A, g, g')$$

 (\mathbf{V}, \mathbf{A})

 $g: X \otimes A \to X \otimes A \otimes A$ in Mark

 $\langle \mathbf{v} \rangle$

$$g':\Delta(X\otimes A)
ightarrow [\mathbb{R},\mathbb{R}]$$
 in Conv $g'(lpha)(r)=\mathbb{E}U(lpha)+eta r$



 $\alpha : \Delta(X \otimes A)$, joint distribution on states and actions.

Riu Rguez Sakamoto

Optics for Value

Iteration

Optics for Value Iteration

Mixed optics for continuous-space MDPs

Savings problem:
$$V'(x) = \max_{0 \le c \le x+i} U(c) + \beta V((1+\gamma)(x-c+i))$$

- V improv: $V'(x) = U(\pi(x)) + \beta V(f(x, \pi(x)))$
 - π improv:

$$\pi'(x) = \operatorname*{argmax}_{0 \le c \le x+i} U(c) + \beta V(f(x,c))$$

Recall:

- Policy: $\pi: x \mapsto c$
- Transition: $f:(x,c)\mapsto (1+\gamma)(x-c+i)$
- Utility: $U: (-, c) \mapsto U(c)$ Stoch: Markov category of measurable spaces and stochastic maps (Gauss: Euclidean spaces and affine functions with noise $Y = MX + \xi$)



Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learning

Related work Conclusion

References

From Value Iteration to Q-learning

Shift from learning state values $V: X \to \mathbb{R}$ to state-action values $q: X \times A \to \mathbb{R}$.

• Value Iteration, in $A^X \times \operatorname{Optic}\left(\binom{X}{\mathbb{R}},\binom{I}{I}\right)$:

$$(\pi,V)\mapsto (x\mapsto rg\max_a(\lambda \ arsim V)(x,a), ar \pi \ arsim \lambda \ arsim V)$$

• State-Action Value Iteration, in $A^X \times \text{Optic}\left(\binom{X \times A}{\mathbb{R}}, \binom{I}{I}\right)$:

$$(\pi,q)\mapsto (x\mapsto rg\max_a q(x,a),\lambda\, {}_{
m S}\, {}_{
m S}\, q)$$



.

Riu Rguez Sakamoto

Dynamic Programmir

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learning

Related work

Conclusion

References



MDPs in Myers' open dyn. systems²: Functor lenses of $BiKl(C \times -, \Delta(\mathbb{R} \times -))^{op}$ Different interfaces:



(a) Interface is a value function (ours)

 \rightarrow (



(b) Interface is a policy (Myers')

$$\binom{X}{\Delta(X\times\mathbb{R})}\to \binom{O}{I}$$

²Jaz Myers, "Double Categories of Open Dynamical Systems (Extended Abstract)".

Riu Rguez Sakamoto

Dynamic Programmii

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics Optics for Val Iteration

Q-learnin

Related wo

Conclusion

References

Summary and future work

- Optics can be tailored to study different MDPs (deterministic, stochastic, discrete-space, continuous-space).
- Q-learning feels ad hoc: use categorical cybernetics to express it with optics

_	One state $(X = 1)$	Many states	Imperfect states
Without agency	-	Markov	Hidden Markov
(fixed π)		chain	Model (HMM)
With agency $(ext{learned} \ \pi)$	" Bandit problem" w/ delayed re- wards	MDP	Partially Ob- servable MDP (POMDP)

- Known/unknown environment
- Problem structure: Adaptive Control, Dynamic Programming, Reinforcement Learning.

Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processes

Dynamic Programmin MDPs

Value iteratior with optics

Optics

Optics for Valı Iteration

Q-learning

Related wo

Conclusion

References

Value iteration is optic composition

Jules Hedges, Riu Rodríguez Sakamoto

Mathematically Structured Programming group, Department of Computer and Information Sciences, University of Strathclyde

riu.rodriguez-sakamoto@strath.ac.uk

5th International Conference on Applied Category Theory, Glasgow (Scotland)

July 21st, 2022

Riu Rguez Sakamoto

Dynamic Programming

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Iteration

Q-learning

Related wor

Conclusion

References

Dynamic Programming theorems

- Policy improvement theorem: Given deterministic policies $\pi,\pi':X o A$,

$$egin{aligned} q_\pi(x,\pi'(x)) &\geq V_\pi(x) \; orall x \in X \ \end{aligned}$$
 implies $V_{\pi'}(x) &\geq V_\pi(x) \; orall x \in X \end{aligned}$

Also for stochastic policies $\pi,\pi':X o\Delta A$, by defining

$$q_{\pi}(s,\pi'(s)) = \sum_{\mathsf{a}} \pi'(\mathsf{a} \mid s) q_{\pi}(s,\mathsf{a})$$

• Convergence proofs require metric. Optics enriched in metric spaces?

Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

Optics

Optics for Valu Iteration

Q-learning

Related wo

Conclusion

References

MDP example: Inverted pendulum

• System of non-linear differential equations³:

$$(M + m)\ddot{y} + mL\ddot{\theta}\cos\theta - mL\dot{\theta}^{2}\sin\theta = a$$
$$mL\ddot{y}\cos\theta + mL^{2}\ddot{\theta} - mLg\sin\theta = 0$$

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & -\frac{mg}{M} & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{(M+m)g}{ML} & 0 \end{pmatrix} \qquad B = \begin{pmatrix} 0 \\ \frac{1}{M} \\ 0 \\ -\frac{1}{ML} \end{pmatrix} \qquad x = \begin{pmatrix} y \\ \dot{y} \\ \dot{\theta} \\ \dot{\theta} \end{pmatrix}$$

• Optimal control approach:

$$J(x,a) = \int_0^\infty C(x(t),a(t)) \mathrm{d}t$$

³Friedland, Control Systems Design, Example 2E.



Riu Rguez

Conv

Sakamoto Dynamic

Introduction Decision Processes Dynamic Programming MDPs

Value iteration with optics

- Optics
- Optics for Valu Iteration
- Q-learning
- Related work

Conclusion

References

- (Conv, \boxtimes , {*}) monoidal product⁴ of two convex spaces $A \boxtimes B$:
 - Free convex structure on $A \times B$, then take the smallest congruence relation on the set s.t.

$$\sum_{i} \alpha_i(a_i, b) \equiv (\sum_{i} \alpha_i a_i, b) \text{ and } \sum_{i} \alpha_i(a, b_i) \equiv (a, \sum_{i} \alpha_i b_i)$$

- Classifies bi-affine functions f : A × B → C, like how A ⊗ B classifies bilinear functions f : A × B → C in R-Mod.
- Monoidal closure: Convex structure on B^A = Conv(A, B) is pointwise: Given f, g ∈ Conv(A, B), (f +_α g)(a) = f(a) +_α g(a).

⁴Stirtz, "Categorical probability theory", Sec 2.2.

Riu Rguez Sakamoto

Dynamic Programmin

Introduction Decision Processe Dynamic Programming MDPs

Value iteratior with optics

Optics Optics for Val

Q-learning

Related worl

Conclusion

References

 Friedland, Bernard. Control Systems Design. An Introduction To State-Space Methods. McGraw-Hill Companies, 1985, p. 513. ISBN: 9780070224414.

 Jaz Myers, David. "Double Categories of Open Dynamical Systems (Extended Abstract)". en. In: *Electronic Proceedings in Theoretical Computer Science* 333 (Feb. 2021), pp. 154–167. ISSN: 2075-2180. DOI: 10.4204/EPTCS.333.11. URL: http://arxiv.org/abs/2005.05956v2 (visited on 12/12/2021).

► Stirtz, Kirk. "Categorical probability theory". arXiv:1406.6030. 2015.