



Machine Intelligence 1

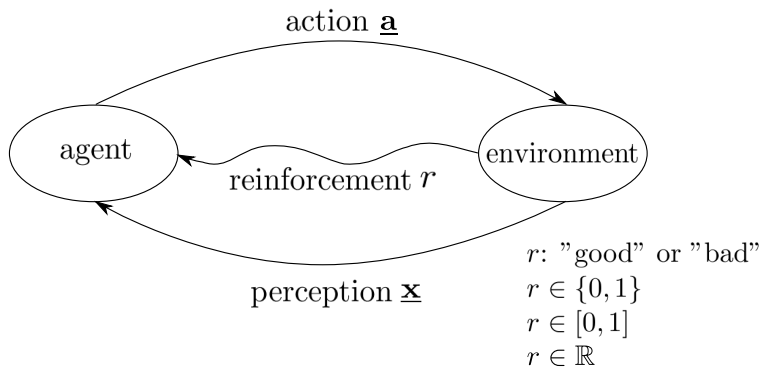
4.1 Reinforcement Learning – Evaluation

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

WS 2016/2017

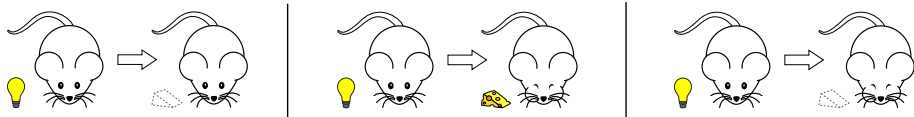
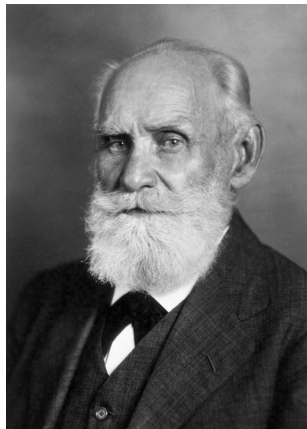
Reinforcement learning



4.1.1 Conditioning

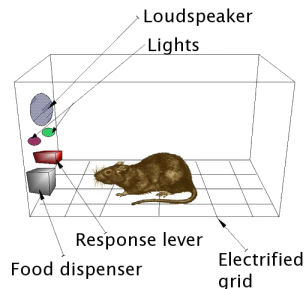
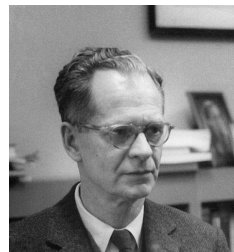
Classical conditioning

- Ivan Pavlov (1849–1936)
- 💡: conditioned stimulus (neutral)
- 🧀: unconditioned stimulus (rewarding)
- experience reinforces *involuntary response*
- animal learns to *expect* reward



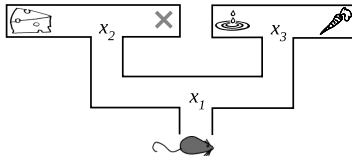
Operant conditioning

- B.F. Skinner (1904–1990)
- animal has to act voluntarily
- actions are rewarded or punished
- experience reinforces *voluntary behavior*
- animal learns how to *achieve* reward



Future rewards

- not all decisions are immediately rewarded
 - decision in **state** x_1 is crucial, but not rewarded
- some decisions require foresight
 - future reward of decision in x_1 depends on decisions in x_2 and x_3
- animal must *delay* the reinforcement of behavior



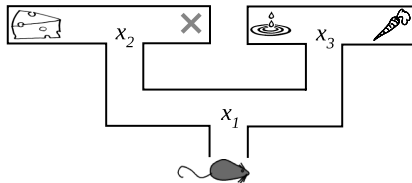
(see Dayan and Niv, 2008; Dayan, 2008)

4.1.2 Markov Decision Processes

Markov decision processes

A Markov decision process (MDP) consist of

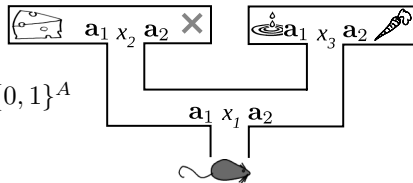
- a set of **states** $\underline{x} \in \mathcal{X}$,
 - e.g. $\mathcal{X} := \{\underline{x}_1, \dots, \underline{x}_S\} \subset \{0, 1\}^S$
with 1-out-of- S encoding



Markov decision processes

A Markov decision process (MDP) consist of

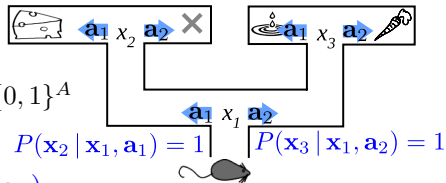
- a set of **states** $\underline{x} \in \mathcal{X}$,
 - e.g. $\mathcal{X} := \{\underline{x}_1, \dots, \underline{x}_S\} \subset \{0, 1\}^S$
with 1-out-of- S encoding
- a set of **actions** $\underline{a} \in \mathcal{A}$,
 - e.g. $\mathcal{A} := \{\underline{a}_1, \dots, \underline{a}_A\} \subset \{0, 1\}^A$
with 1-out-of- A encoding



Markov decision processes

A Markov decision process (MDP) consist of

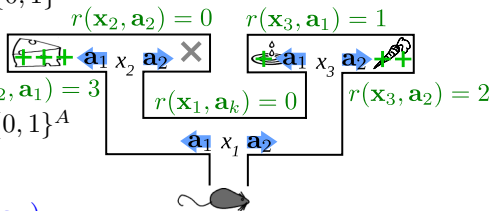
- a set of **states** $\underline{x} \in \mathcal{X}$,
 - e.g. $\mathcal{X} := \{\underline{x}_1, \dots, \underline{x}_S\} \subset \{0, 1\}^S$
with 1-out-of- S encoding
- a set of **actions** $\underline{a} \in \mathcal{A}$,
 - e.g. $\mathcal{A} := \{\underline{a}_1, \dots, \underline{a}_A\} \subset \{0, 1\}^A$
with 1-out-of- A encoding
- a **transition model** $P(\underline{x}_j | \underline{x}_i, \underline{a}_k)$
 - probability to end up in \underline{x}_j after choosing \underline{a}_k in \underline{x}_i
 - stationary distribution (Markov property)



Markov decision processes

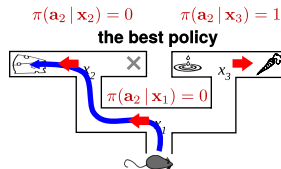
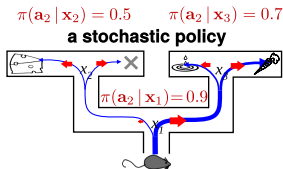
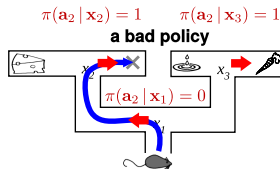
A Markov decision process (MDP) consist of

- a set of **states** $\underline{x} \in \mathcal{X}$,
 - e.g. $\mathcal{X} := \{\underline{x}_1, \dots, \underline{x}_S\} \subset \{0, 1\}^S$
with 1-out-of- S encoding
- a set of **actions** $\underline{a} \in \mathcal{A}$,
 - e.g. $\mathcal{A} := \{\underline{a}_1, \dots, \underline{a}_A\} \subset \{0, 1\}^A$
with 1-out-of- A encoding
- a **transition model** $P(\underline{x}_j | \underline{x}_i, \underline{a}_k)$
 - probability to end up in \underline{x}_j after choosing \underline{a}_k in \underline{x}_i
 - stationary distribution (Markov property)
- a bounded **reward function** $r(\underline{x}_i, \underline{a}_k)$
 - denotes the *average immediate reward* for choosing \underline{a}_k in \underline{x}_i
 - extension with randomized rewards possible



Policy

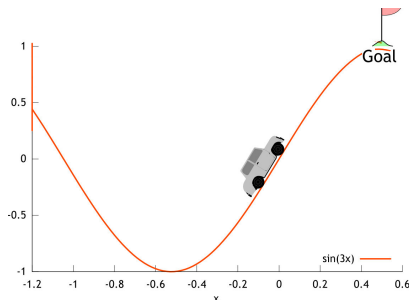
- the agent's behavior is expressed by a **policy** $\pi(\underline{a}_k | \underline{x}_i)$
 - the probability that the agent chooses \underline{a}_k in \underline{x}_i



- the goal of RL is to find the “optimal policy” π^*

Example: mountain car

- a car in a valley between mountains
 - \mathcal{X} : position and velocity
- agent drives the car
 - \mathcal{A} : forward, backward, nothing (i.e., accelerate the car by $+a$, $-a$ and 0)
- dynamics are given by physics
 - transition model P simulated
 - gravitation but no friction
- goal: reach right hilltop
 - reward $r=0$, except $r=1$ at goal
- but car is underpowered
 - policy π must first pick up speed



Markov chains

- a Markov chain of length p

- is a sequence of states and actions

$$\{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}\}_{t=0}^p \subset \mathcal{X} \times \mathcal{A}$$

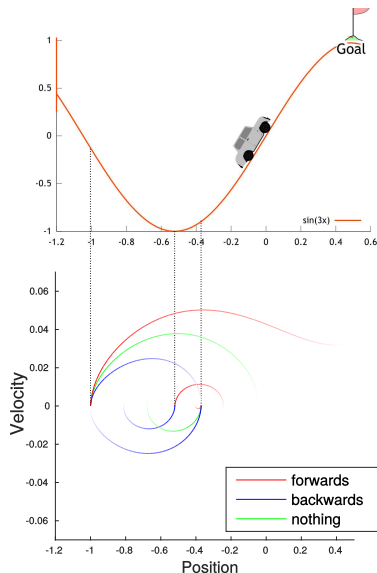
- actions $\underline{\mathbf{a}}^{(t)}$ are drawn from **policy**:

$$\underline{\mathbf{a}}^{(t)} \sim \pi(\cdot | \underline{\mathbf{x}}^{(t)})$$

- successive states $\underline{\mathbf{x}}^{(t+1)}$ are drawn from **transition model**:

$$\underline{\mathbf{x}}^{(t+1)} \sim P(\cdot | \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)})$$

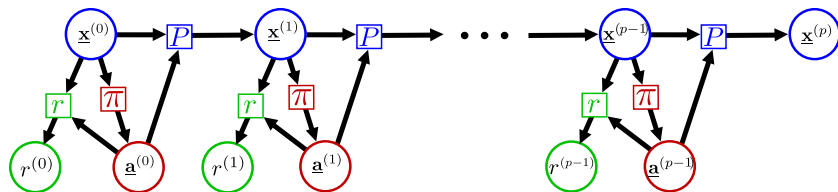
- given an MDP, a Markov chain depends on initial $\underline{\mathbf{x}}^{(0)}$ and **policy** π



Markov chain distribution

- Markov chains are sets of random variables
 - depend on initial state $\underline{x}^{(0)}$ and policy π
- joint distribution of states in a Markov chain factorizes

$$P(\underline{x}^{(0)}, \dots, \underline{x}^{(p)}) = P(\underline{x}^{(0)}) \prod_{t=0}^{p-1} \sum_{\mathbf{a}_k}^A \pi(\mathbf{a}_k | \underline{x}^{(t)}) P(\underline{x}^{(t+1)} | \underline{x}^{(t)}, \mathbf{a}_k)$$



(see blackboard)

4.1.3 Policy Evaluation

Value function

- a **value function** measures the quality of a policy π in state $\underline{\mathbf{x}}^{(0)}$
 - $V^\pi(\underline{\mathbf{x}}^{(0)})$ is the *expected* *reward*

$$V^\pi(\underline{\mathbf{x}}^{(0)}) = \mathbb{E} \left[r(\underline{\mathbf{x}}^{(0)}, \underline{\mathbf{a}}^{(0)}) \mid \underline{\mathbf{a}}^{(0)} \sim \pi(\cdot \mid \underline{\mathbf{x}}^{(0)}) \right]$$

- *average* over *selected actions*

Value function

- a **value function** measures the quality of a policy π in state $\underline{\mathbf{x}}^{(0)}$
 - $V^\pi(\underline{\mathbf{x}}^{(0)})$ is the *expected finite sum of rewards*

$$V^\pi(\underline{\mathbf{x}}^{(0)}) = \mathbb{E} \left[\sum_{t=0}^H r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \mid \begin{array}{l} \underline{\mathbf{a}}^{(t)} \sim \pi(\cdot | \underline{\mathbf{x}}^{(t)}) \\ \underline{\mathbf{x}}^{(t+1)} \sim P(\cdot | \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \end{array} \right]$$

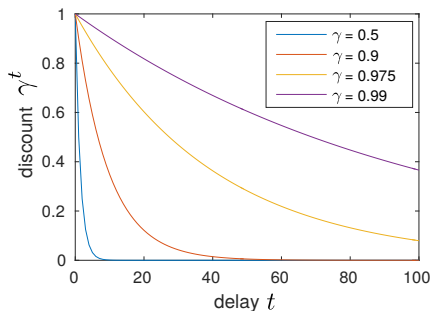
- *average* over Markov chains
- Markov chains start at $\underline{\mathbf{x}}^{(0)}$
and follow the transition
model P and the *policy* π

Value function

- a **value function** measures the quality of a policy π in state $\underline{\mathbf{x}}^{(0)}$
 - $V^\pi(\underline{\mathbf{x}}^{(0)})$ is the *expected infinite sum of discounted future rewards*

$$V^\pi(\underline{\mathbf{x}}^{(0)}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \mid \begin{array}{l} \underline{\mathbf{a}}^{(t)} \sim \pi(\cdot \mid \underline{\mathbf{x}}^{(t)}) \\ \underline{\mathbf{x}}^{(t+1)} \sim P(\cdot \mid \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \end{array} \right], \quad \gamma \in [0, 1].$$

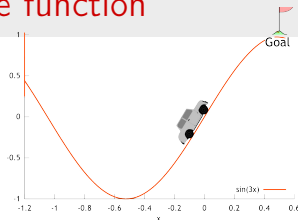
- *average* over Markov chains
- Markov chains start at $\underline{\mathbf{x}}^{(0)}$ and follow the transition model P and the **policy** π
- **discount factor** γ : preference for short- vs. long-term goals



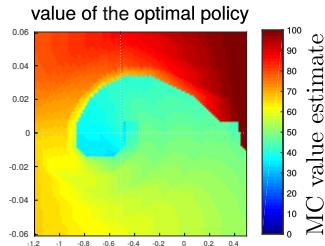
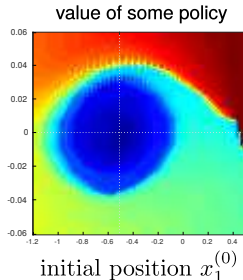
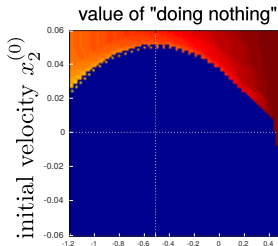
Monte Carlo (MC) estimation of the value function

- finite approximation of infinite Markov chains

- rewards weighted by $\gamma^H < \epsilon$ are neglected
- value is the discounted reward averaged over n Markov chains of length H
- n must be sufficiently large



- requires simulator to draw n chains from the same initial state $\underline{x}^{(0)}$
- every state must be evaluated often \leadsto not sample efficient



The Bellman equation (1)

$$\begin{aligned}
 V^\pi(\underline{\mathbf{x}}_i) &= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \mid \begin{array}{l} \underline{\mathbf{x}}^{(0)} := \underline{\mathbf{x}}_i \\ \underline{\mathbf{a}}^{(t)} \sim \pi(\cdot \mid \underline{\mathbf{x}}^{(t)}) \\ \underline{\mathbf{x}}^{(t+1)} \sim P(\cdot \mid \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \end{array} \right] \\
 &= \mathbb{E} \left[r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}^{(0)}) \mid \underline{\mathbf{a}}^{(0)} \sim \pi(\cdot \mid \underline{\mathbf{x}}_i) \right] \\
 &\quad + \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^t r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \mid \begin{array}{l} \underline{\mathbf{x}}^{(0)} := \underline{\mathbf{x}}_i \\ \underline{\mathbf{a}}^{(t)} \sim \pi(\cdot \mid \underline{\mathbf{x}}^{(t)}) \\ \underline{\mathbf{x}}^{(t+1)} \sim P(\cdot \mid \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \end{array} \right] \\
 &= \mathbb{E} \left[r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}^{(0)}) \mid \underline{\mathbf{a}}^{(0)} \sim \pi(\cdot \mid \underline{\mathbf{x}}_i) \right] + \gamma \mathbb{E} \left[V^\pi(\underline{\mathbf{x}}^{(1)}) \mid \begin{array}{l} \underline{\mathbf{a}}^{(0)} \sim \pi(\cdot \mid \underline{\mathbf{x}}_i) \\ \underline{\mathbf{x}}^{(1)} \sim P(\cdot \mid \underline{\mathbf{x}}_i, \underline{\mathbf{a}}^{(0)}) \end{array} \right] \\
 &= \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k \mid \underline{\mathbf{x}}_i) \left(r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) + \gamma \sum_{j=1}^S P(\underline{\mathbf{x}}_j \mid \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) V^\pi(\underline{\mathbf{x}}_j) \right)
 \end{aligned}$$



Richard E. Bellman
(1920–1984)

$\underline{\mathbf{x}}_i \in \{0, 1\}^S$: 1-out-of- S coded state i

The Bellman equation (2)

$$\begin{aligned}
 V^\pi(\underline{\mathbf{x}}_i) &= \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) \left(r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) + \gamma \sum_{j=1}^S P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) V^\pi(\underline{\mathbf{x}}_j) \right) \\
 &= \underbrace{\sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k)}_{\text{"controlled" reward function } r_i^\pi} + \underbrace{\gamma \sum_{j=1}^S \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k)}_{\text{"controlled" transition model } P_{ij}^\pi} V^\pi(\underline{\mathbf{x}}_j)
 \end{aligned}$$

$$\begin{aligned}
 \underline{\mathbf{v}}^\pi &= \underline{\mathbf{r}}^\pi + \gamma \underline{\mathbf{P}}^\pi \underline{\mathbf{v}}^\pi, \quad \text{with } \left\{ \begin{array}{l} r_i^\pi := \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) \\ P_{ij}^\pi := \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) \end{array} \right. \\
 &=: \hat{B}^\pi[\underline{\mathbf{v}}^\pi]
 \end{aligned}$$

"controlled" models $\underline{\mathbf{r}}^\pi \in \mathbb{R}^S$ and $\underline{\mathbf{P}}^\pi \in \mathbb{R}^{S \times S}$

$\underline{\mathbf{x}}_i \in \{0, 1\}^S$: 1-out-of- S coded state i ,

$\underline{\mathbf{v}}^\pi \in \mathbb{R}^S$: vector containing all values V^π

4.1.4 Model-based Approaches

The analytic solution of the Bellman equation

Bellman operator \hat{B}^π for discrete state values

$$\hat{B}^\pi[\tilde{\mathbf{v}}] = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \tilde{\mathbf{v}}, \quad \forall \tilde{\mathbf{v}} \in \mathbb{R}^S$$

- Bellman operator \hat{B}^π of **policy** π uses “controlled” models
 - of the reward function $\mathbf{r}^\pi \in \mathbb{R}^S$
 - and transition model $\mathbf{P}^\pi \in \mathbb{R}^{S \times S}$

- \hat{B}^π has an analytic solution of the value function $\mathbf{v}^\pi \in \mathbb{R}^S$

$$\mathbf{v}^\pi = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \mathbf{v}^\pi \leadsto (\mathbf{I} - \gamma \mathbf{P}^\pi) \mathbf{v}^\pi = \mathbf{r}^\pi \leadsto \mathbf{v}^\pi = (\mathbf{I} - \gamma \mathbf{P}^\pi)^{-1} \mathbf{r}^\pi$$

- matrix $(\mathbf{I} - \gamma \mathbf{P}^\pi) \in \mathbb{R}^{S \times S}$ is always invertible
 - $|\lambda_k| \leq 1$ for all eigenvalues λ_k of transition matrices \mathbf{P}^π
 - discount factor $\gamma < 1$

(see e.g. Bertsekas, 2007, for details)

Model-based value iteration

- the value function $\underline{\mathbf{v}}^\pi$ is the **fixed-point** of the Bellman operator \hat{B}^π

$$\underline{\mathbf{v}}^\pi = \hat{B}^\pi[\underline{\mathbf{v}}^\pi] = \underline{\mathbf{r}}^\pi + \gamma \underline{\mathbf{P}}^\pi \underline{\mathbf{v}}^\pi$$

- **value iteration**: repeated application of the Bellman operator

$$\tilde{\underline{\mathbf{v}}}^{\pi(t+1)} = \underline{\mathbf{r}}^\pi + \gamma \underline{\mathbf{P}}^\pi \tilde{\underline{\mathbf{v}}}^{\pi(t)}$$

- is value iteration convergent, i.e. $\lim_{t \rightarrow \infty} \tilde{\underline{\mathbf{v}}}^{\pi(t)} = \underline{\mathbf{v}}^\pi$?

Convergence of value iteration

Contraction mapping (in supremum norm)

A function $\hat{B} : \mathbb{R}^S \rightarrow \mathbb{R}^S$ is called a *contraction mapping* with Lipschitz constant $\lambda < 1$ if $\max_j |(\hat{B}[\tilde{\mathbf{v}}] - \hat{B}[\tilde{\mathbf{w}}])_j| \leq \lambda \max_j |\tilde{v}_j - \tilde{w}_j|, \forall \tilde{\mathbf{v}}, \tilde{\mathbf{w}} \in \mathbb{R}^S$.

■ application to the Bellman operator $\hat{B}^\pi[\tilde{\mathbf{v}}] = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \tilde{\mathbf{v}}$

$$\begin{aligned} \max_j |\hat{B}^\pi[\tilde{\mathbf{v}}]_j - \hat{B}^\pi[\tilde{\mathbf{w}}]_j| &= \max_j |r_j^\pi + \gamma (\mathbf{P}^\pi \tilde{\mathbf{v}})_j - r_j^\pi - \gamma (\mathbf{P}^\pi \tilde{\mathbf{w}})_j| \\ &\stackrel{(i)}{\leq} \max_j \gamma (\mathbf{P}^\pi |\tilde{\mathbf{v}} - \tilde{\mathbf{w}}|)_j \stackrel{(ii)}{\leq} \gamma \max_j |\tilde{v}_j - \tilde{w}_j| \end{aligned}$$

$$(i) \quad \left| \sum_{i=1}^S P_{ji}^\pi x_i \right| \leq \sum_{i=1}^S P_{ji}^\pi |x_i|, \quad \forall \mathbf{x} \in \mathbb{R}^S \quad (\text{Jensen's inequality})$$

$$(ii) \quad \sum_{i=1}^S P_{ji}^\pi |x_i| \leq \sum_{i=1}^S P_{ji}^\pi \max_{1 \leq k \leq S} |x_k| = \max_{1 \leq k \leq S} |x_k| \quad \left(\sum_{i=1}^S P_{ji}^\pi = 1 \right)$$

Convergence of value iteration

Contraction mapping (in supremum norm)

A function $\hat{B} : \mathbb{R}^S \rightarrow \mathbb{R}^S$ is called a *contraction mapping* with Lipschitz constant $\lambda < 1$ if $\max_j |(\hat{B}[\tilde{\mathbf{v}}] - \hat{B}[\tilde{\mathbf{w}}])_j| \leq \lambda \max_j |\tilde{v}_j - \tilde{w}_j|, \forall \tilde{\mathbf{v}}, \tilde{\mathbf{w}} \in \mathbb{R}^S$.

- \hat{B}^π is a contraction mapping with Lipschitz constant γ

$$\tilde{\mathbf{v}}^{\pi(t+1)} = \hat{B}^\pi[\tilde{\mathbf{v}}^{\pi(t)}] = \mathbf{r}^\pi + \gamma \mathbf{P}^\pi \tilde{\mathbf{v}}^{\pi(t)} \quad (\text{value iteration})$$

- $\max_j |(\hat{B}^\pi[\mathbf{v}^{\pi(t)}])_j - v_j^\pi| \leq \gamma \max_j |v_j^{\pi(t)} - v_j^\pi|$ for *any* $\mathbf{v}^{\pi(t)} \in \mathbb{R}^S$

\Rightarrow value iteration converges in the limit to unique fix-point \mathbf{v}^π

- number of iterations until convergence $\sim -\frac{1}{\log(\gamma)}$
- analytic solution is faster for large γ

4.1.5 Model-free Approaches: Online Value Estimation

Inductive value estimation

- agent must learn through interaction with the environment
 - “controlled” models $\underline{\mathbf{r}}^\pi$ and $\underline{\mathbf{P}}^\pi$ are not available

$$V^\pi(\underline{\mathbf{x}}_i) = \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) \left(r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) + \gamma \sum_{j=1}^S P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) V^\pi(\underline{\mathbf{x}}_j) \right)$$

- estimate value function *inductively* from one long Markov chain
 - actions are drawn according to the policy $\underline{\mathbf{a}}^{(t)} \sim \pi(\cdot | \underline{\mathbf{x}}^{(t)})$
 - which lead to transitions $\underline{\mathbf{x}}^{(t+1)} \sim P(\cdot | \underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)})$
 - and yield rewards $r_t := r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)})$

Temporal difference (TD) learning

- online estimation named after the difference in values (TD-error ΔV_t)

$$\tilde{V}_{t+1}^{\pi}(\underline{\mathbf{x}}^{(t)}) = \tilde{V}_t^{\pi}(\underline{\mathbf{x}}^{(t)}) + \eta \underbrace{\left(\textcolor{green}{r}_t + \gamma \tilde{V}_t^{\pi}(\textcolor{blue}{\underline{\mathbf{x}}}^{(t+1)}) - \tilde{V}_t^{\pi}(\underline{\mathbf{x}}^{(t)}) \right)}_{\text{TD-error } \Delta V_t}$$

- TD learning performs value iteration *on average*

- for the average over all Markov chains that pass $\underline{\mathbf{x}}_i$ at time t holds:

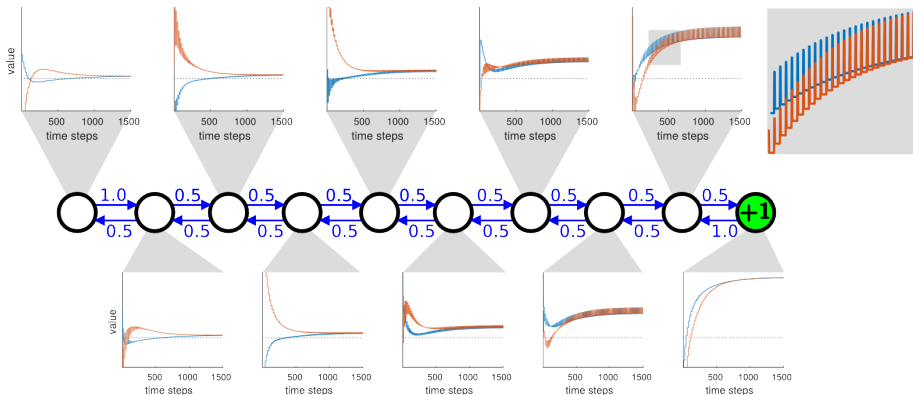
$$\underbrace{\mathbb{E}[\tilde{V}_{t+1}^{\pi}(\underline{\mathbf{x}}^{(t)})]}_{\tilde{v}_i^{\pi(t+1)}} = (1 - \eta) \underbrace{\mathbb{E}[\tilde{V}_t^{\pi}(\underline{\mathbf{x}}^{(t)})]}_{\tilde{v}_i^{\pi(t)}} + \eta \underbrace{\left(\mathbb{E}[\textcolor{green}{r}_t] + \gamma \mathbb{E}[\tilde{V}_t^{\pi}(\textcolor{blue}{\underline{\mathbf{x}}}^{(t+1)})] \right)}_{(\textcolor{green}{\mathbf{r}}^{\pi} + \gamma \textcolor{blue}{\mathbf{P}}^{\pi} \tilde{\mathbf{v}}^{\pi(t)})_i}$$

- *asynchronous online estimate* of $\hat{B}^{\pi}[\tilde{\mathbf{v}}^{\pi(t)}] = \textcolor{green}{\mathbf{r}}^{\pi} + \gamma \textcolor{blue}{\mathbf{P}}^{\pi} \tilde{\mathbf{v}}^{\pi(t)}$
 - asynchronous update of one state at a time
 - estimates Bellman operator \hat{B}^{π} by online average

(see Sutton and Barto, 1998)

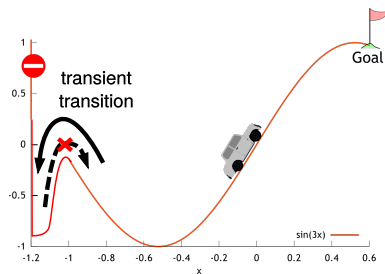
Convergence of TD learning

- **example:** Markov chain running back and forth on 10 states
 - two randomly initialized value functions (red/blue)
 - deterministic transitions with stochastic policy
 - rightmost state is rewarded, $\gamma = 0.95$, $\eta = 0.5$
- TD learning **contracts** different initializations, but does **not converge**



Requirements for contraction

- TD learning contracts
 - for an infinite Markov chain,
 - which visits *all* states infinitely often
- no **transient** transitions allowed
 - transitions must be reversible
 - “you cannot learn from death”
- **positive recurrence**: a non-zero probability to return in finite time



Ergodic Markov chains

Ergodicity

A Markov chain is **ergodic** if it is **positively recurrent** (non-zero probability to leave any state and eventually return to it) and **aperiodic** (returns to the same state can occur at irregular times).

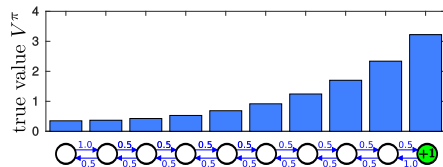
■ **steady state distribution** $P_{ss}(\underline{x}) > 0$ exists and visits all states \underline{x}

⇒ TD learning is a contraction mapping for *ergodic* Markov chains

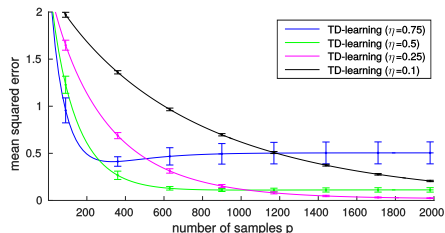
Influence of learning rate η

$$\begin{aligned}\tilde{V}_{t+1}^{\pi}(\underline{\mathbf{x}}^{(t)}) &= \tilde{V}_t^{\pi}(\underline{\mathbf{x}}^{(t)}) + \eta \Delta V_t \\ \Delta V_t &= \textcolor{green}{r}_t + \gamma \tilde{V}_t^{\pi}(\textcolor{blue}{\underline{\mathbf{x}}}^{(t+1)}) - \tilde{V}_t^{\pi}(\underline{\mathbf{x}}^{(t)})\end{aligned}$$

- stochastic transitions/rewards
 $\leadsto \tilde{V}_t^{\pi}$ may not converge
- TD learning let \tilde{V}_t^{π} fluctuate around the true value function V^{π}
- influence of the learning rate η
 - large η : fast learning, large variance
 - small η : slow learning, small variance
 - decaying η_t are not practical as ΔV_t are (initially) not stationary



- 10 states Markov chain
- regular movement back and forth
- rightmost state rewarded, $\gamma = 0.95$



4.1.6 Model-free Approaches: Eligibility Traces & $TD(\lambda)$

Value propagation in TD learning

$$\begin{aligned}\tilde{V}_{t+1}^{\pi}(\mathbf{x}^{(t)}) &= \tilde{V}_t^{\pi}(\mathbf{x}^{(t)}) + \eta \Delta V_t \\ \Delta V_t &= r_t + \gamma \tilde{V}_t^{\pi}(\mathbf{x}^{(t+1)}) - \tilde{V}_t^{\pi}(\mathbf{x}^{(t)})\end{aligned}$$

- TD learning propagates values one step into the past

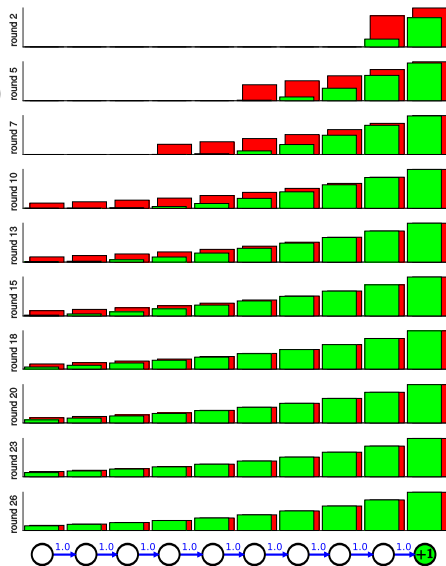
- many steps to convergence

- deterministic example:

- 10 states, 1 action
- only forward transitions
- reward in last state
- $\gamma = 0.9$; $\eta = 1$ or $\eta = 0.5$

- value propagation requires

- exactly 10 rounds ($\eta = 1$)
- roughly 26 rounds ($\eta = 0.5$)



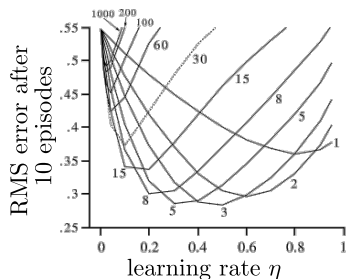
n -step temporal difference learning

- accumulation of observed rewards

$$\begin{aligned}
 R_t^{(1)} &= r_t + \gamma \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+1)}) \\
 R_t^{(2)} &= r_t + \gamma r_{t+1} + \gamma^2 \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+2)}) \\
 &\vdots \\
 R_t^{(n)} &= \sum_{\tau=0}^{n-1} \gamma^\tau r_{t+\tau} + \gamma^n \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+n)})
 \end{aligned}$$

- online estimation similar to TD learning

$$\tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) \leftarrow \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) + \eta \left(R_t^{(n)} - \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) \right)$$



RMS averaged over 100 random-walks on a 19-state chain, rewarded at one end

(Sutton and Barto, 1998)

Discounted average

$$\tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) \leftarrow \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) + \eta \left(R_t^{(n)} - \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t)}) \right)$$

- there is an optimal combination of η and n , however,
 - agent must memorize the last n steps
 - values are updated with a delay of n steps
- trick: consider a discounted average of $R_t^{(n)}$

$$R_t^\lambda = (1 - \lambda) \sum_{k=0}^{\infty} \lambda^k R_t^{(k+1)}$$

Eligibility traces & TD(λ)

- the **eligibility trace** $\underline{e}^{(t)} \in \mathbb{R}^S$ stores the past of state \underline{x}_i

$$e_i^{(t)} = \sum_{k=0}^t (\gamma \lambda)^{t-k} \delta_{ik}, \quad \delta_{ik} = \underline{x}_i^\top \underline{x}^{(k)} \quad \forall \underline{x}_i \in \mathcal{X},$$

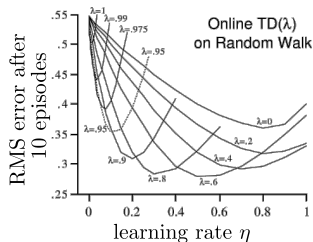
- The **TD(λ) method**:

$$\begin{aligned} \tilde{V}_{t+1}^\pi(\underline{x}_i) &= \tilde{V}_t^\pi(\underline{x}_i) + \eta e_i \left(\overbrace{r_t + \gamma \tilde{V}_t^\pi(\underline{x}^{(t+1)}) - \tilde{V}_t^\pi(\underline{x}^{(t)})}^{\text{TD-error } \Delta V_t} \right) \\ \underline{e}^{(t+1)} &= \gamma \lambda \underline{e}^{(t)} + \underline{x}^{(t+1)} \end{aligned}$$

- TD(0)**: TD learning as defined before

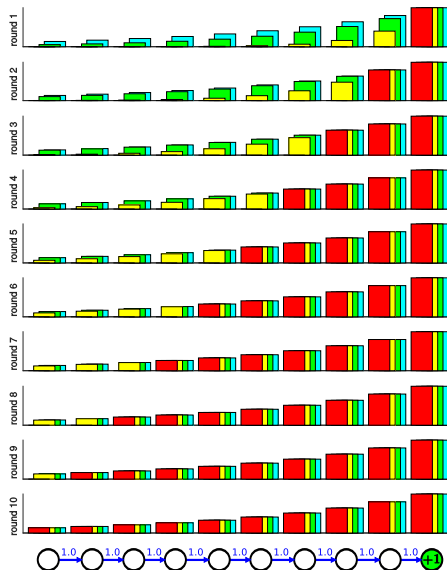
RMS averaged over 100 random-walks
on a 19-state chain, rewarded at one end

(Sutton and Barto, 1998)



Value propagation in TD(λ)

- deterministic example:
 - 10 states, 1 action
 - only forward transitions
 - reward in last state
 - $\gamma = 1, \eta = 1$
- value propagation finishes
 - after 1 round with $\lambda = 1$
 - after 4 rounds with $\lambda = 0.9$
 - after 7 rounds with $\lambda = 0.5$
 - after 10 rounds with $\lambda = 0$



4.1.7 Model-free approaches: Batch Value Estimation

Reminder: the Bellman equation



Richard E. Bellman
(1920–1984)

$$V^\pi(\underline{\mathbf{x}}_i) = \underbrace{\sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k)}_{\text{"controlled" reward function } \underline{\mathbf{r}}_i^\pi} + \gamma \underbrace{\sum_{j=1}^S \sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k)}_{\text{"controlled" transition model } \underline{\mathbf{P}}_{ij}^\pi} V^\pi(\underline{\mathbf{x}}_j)$$

$$\underline{\mathbf{v}}^\pi = \hat{B}^\pi[\underline{\mathbf{v}}^\pi] = \underline{\mathbf{r}}^\pi + \gamma \underline{\mathbf{P}}^\pi \underline{\mathbf{v}}^\pi$$

$\underline{\mathbf{x}}_i \in \{0, 1\}^S$: 1-out-of- S coded state i ,

$\underline{\mathbf{r}}^\pi \in \mathbb{R}^S$ “controlled” reward function,

$\underline{\mathbf{v}}^\pi \in \mathbb{R}^S$: vector containing all values V^π

$\underline{\mathbf{P}}^\pi \in \mathbb{R}^{S \times S}$ “controlled” transition model

Batch approximation of the Bellman operator (1)

- approximate $\tilde{V}_{t+1}^\pi \approx \hat{B}^\pi[\tilde{V}_t^\pi]$ using samples from an
ergodic Markov chain $\{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}\}_{t=0}^p$, executing policy π

$$\hat{B}^\pi[\tilde{V}_t^\pi](\underline{\mathbf{x}}_i) = \underbrace{\sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) \left(r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) + \gamma \overbrace{\sum_{j=1}^S P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) \tilde{V}_t^\pi(\underline{\mathbf{x}}_j)}^{\text{average } \{\tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+1)}) | \underline{\mathbf{x}}^{(t)} = \underline{\mathbf{x}}_i\}} \right)}_{\text{approximate by averaging over } \{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)} | \underline{\mathbf{x}}^{(t)} = \underline{\mathbf{x}}_i, \underline{\mathbf{a}}^{(t)} \sim \pi\}}$$

Batch approximation of the Bellman operator (1)

- approximate $\tilde{V}_{t+1}^\pi \approx \hat{B}^\pi[\tilde{V}_t^\pi]$ using samples from an
ergodic Markov chain $\{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}\}_{t=0}^p$, executing policy π

$$\begin{aligned} \hat{B}^\pi[\tilde{V}_t^\pi](\underline{\mathbf{x}}_i) &= \underbrace{\sum_{k=1}^A \pi(\underline{\mathbf{a}}_k | \underline{\mathbf{x}}_i) \left(r(\underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) + \gamma \underbrace{\sum_{j=1}^S P(\underline{\mathbf{x}}_j | \underline{\mathbf{x}}_i, \underline{\mathbf{a}}_k) \tilde{V}_t^\pi(\underline{\mathbf{x}}_j)}_{\text{average } \{\tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+1)}) | \underline{\mathbf{x}}^{(t)} = \underline{\mathbf{x}}_i\}} \right)}_{\text{approximate by averaging over } \{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)} | \underline{\mathbf{x}}^{(t)} = \underline{\mathbf{x}}_i, \underline{\mathbf{a}}^{(t)} \sim \pi\}} \\ &\approx \underbrace{\frac{1}{\sum_{\tau=0}^{p-1} \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(\tau)}}}_{\text{normalization}} \sum_{t=0}^{p-1} \underbrace{\underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(t)}}_{\text{selection}} \left(r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) + \gamma \tilde{V}_t^\pi(\underline{\mathbf{x}}^{(t+1)}) \right) \end{aligned}$$

- selection $\underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(t)}$ applies update only for states $\underline{\mathbf{x}}_i = \underline{\mathbf{x}}^{(t)}$
- normalization $\sum_{\tau=0}^{p-1} \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(\tau)}$ counts how often $\underline{\mathbf{x}}_i$ appears in batch

Batch approximation of the Bellman operator (2)

- approximate $\tilde{V}_{t+1}^\pi \approx \hat{B}^\pi[\tilde{V}_t^\pi]$ using samples from an **ergodic Markov chain** $\{\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}\}_{t=0}^p$, **executing policy π**

$$\begin{aligned}\hat{B}^\pi[\tilde{V}^\pi](\underline{\mathbf{x}}_i) &\approx \frac{1}{\sum_{\tau=0}^{p-1} \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(\tau)}} \sum_{t=0}^{p-1} \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(t)} \left(r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) + \gamma \tilde{V}^\pi(\underline{\mathbf{x}}^{(t+1)}) \right) \\ &= \underline{\mathbf{x}}_i^\top \underbrace{\left(\sum_{\tau=0}^{p-1} \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(\tau)} \right)^{-1}}_{\mathbf{C}_{ii}} \underbrace{\left(\sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \right)}_{\underline{\mathbf{b}}} + \gamma \underbrace{\sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} \tilde{V}^\pi(\underline{\mathbf{x}}^{(t+1)})}_{\underline{\mathbf{D}}^\pi \tilde{\mathbf{v}}^\pi}\end{aligned}$$

$$\hat{B}^\pi[\tilde{\mathbf{v}}^\pi] \approx \underline{\mathbf{C}}^{-1}(\underline{\mathbf{b}} + \gamma \underline{\mathbf{D}}^\pi \tilde{\mathbf{v}}^\pi)$$

$\underline{\mathbf{C}} = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} (\underline{\mathbf{x}}^{(t)})^\top \in \mathbb{R}^{S \times S}$
diagonal normalization matrix

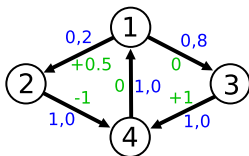
$\underline{\mathbf{D}}^\pi = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} (\underline{\mathbf{x}}^{(t+1)})^\top \in \mathbb{R}^{S \times S}$
matrix to count **transitions**

$\underline{\mathbf{b}} = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} r(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)}) \in \mathbb{R}^S$
vector with the sum of **rewards**

Example batch approximation

- the approximated Bellman operator:

$$\hat{B}^{\pi}[\tilde{\mathbf{v}}^{\pi}] \approx \underline{\mathbf{C}}^{-1}(\underline{\mathbf{b}} + \gamma \underline{\mathbf{D}}^{\pi} \tilde{\mathbf{v}}^{\pi})$$



- example MDP with 4 states:

- transition probabilities $\underline{\mathbf{P}}^{\pi} \approx \underline{\mathbf{C}}^{-1} \underline{\mathbf{D}}^{\pi}$
- reward for transitions $\underline{\mathbf{r}}^{\pi} \approx \underline{\mathbf{C}}^{-1} \underline{\mathbf{b}}$

$$\underline{\mathbf{C}} = \begin{bmatrix} 10 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 7 & 0 \\ 0 & 0 & 0 & 10 \end{bmatrix},$$

state visit count

$$\underline{\mathbf{D}}^{\pi} = \begin{bmatrix} 0 & 3 & 7 & 0 \\ 0 & 0 & 0 & 3 \\ 0 & 0 & 0 & 7 \\ 10 & 0 & 0 & 0 \end{bmatrix},$$

transition count

$$\underline{\mathbf{b}} = \begin{bmatrix} 1.5 \\ -3 \\ +7 \\ 0 \end{bmatrix}$$

collected rewards

$$\underline{\mathbf{C}} = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} (\underline{\mathbf{x}}^{(t)})^{\top} \in \mathbb{R}^{S \times S} \quad \underline{\mathbf{D}}^{\pi} = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} (\underline{\mathbf{x}}^{(t+1)})^{\top} \in \mathbb{R}^{S \times S} \quad \underline{\mathbf{b}} = \sum_{t=0}^{p-1} \underline{\mathbf{x}}^{(t)} r_{(\underline{\mathbf{x}}^{(t)}, \underline{\mathbf{a}}^{(t)})} \in \mathbb{R}^S$$

Solution to the approximated Bellman operator

- the approximated Bellman operator:

$$\hat{B}^{\pi}[\tilde{\mathbf{v}}^{\pi}] \approx \underline{\mathbf{C}}^{-1}(\underline{\mathbf{b}} + \gamma \underline{\mathbf{D}}^{\pi} \tilde{\mathbf{v}}^{\pi})$$

- fixed-point $\mathbf{v}^* \approx \hat{B}^{\pi}[\mathbf{v}^*]$ can be computed analytically

$$\mathbf{v}^* = (\underline{\mathbf{C}} - \gamma \underline{\mathbf{D}}^{\pi})^{-1} \underline{\mathbf{b}}$$

Solution to the approximated Bellman operator

- the approximated Bellman operator:

$$\hat{B}^{\pi}[\tilde{\mathbf{v}}^{\pi}] \approx \underline{\mathbf{C}}^{-1}(\underline{\mathbf{b}} + \gamma \underline{\mathbf{D}}^{\pi} \tilde{\mathbf{v}}^{\pi})$$

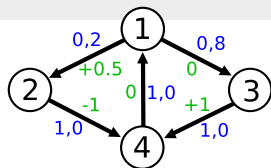
- fixed-point $\mathbf{v}^* \approx \hat{B}^{\pi}[\mathbf{v}^*]$ can be computed analytically

$$\mathbf{v}^* = (\underline{\mathbf{C}} - \gamma \underline{\mathbf{D}}^{\pi})^{-1} \underline{\mathbf{b}} = (\underline{\mathbf{I}} - \gamma \underbrace{\underline{\mathbf{C}}^{-1} \underline{\mathbf{D}}^{\pi}}_{\tilde{\mathbf{P}}^{\pi}})^{-1} \underbrace{\underline{\mathbf{C}}^{-1} \underline{\mathbf{b}}}_{\tilde{\mathbf{r}}^{\pi}} = (\underline{\mathbf{I}} - \gamma \tilde{\mathbf{P}}^{\pi})^{-1} \tilde{\mathbf{r}}^{\pi}$$

Solution to the approximated Bellman operator

- the approximated Bellman operator:

$$\hat{B}^{\pi}[\tilde{\mathbf{v}}^{\pi}] \approx \underline{\mathbf{C}}^{-1}(\underline{\mathbf{b}} + \gamma \underline{\mathbf{D}}^{\pi} \tilde{\mathbf{v}}^{\pi})$$



- fixed-point $\mathbf{v}^* \approx \hat{B}^{\pi}[\mathbf{v}^*]$ can be computed analytically

$$\mathbf{v}^* = (\underline{\mathbf{C}} - \gamma \underline{\mathbf{D}}^{\pi})^{-1} \underline{\mathbf{b}} = (\underline{\mathbf{I}} - \underbrace{\gamma \underline{\mathbf{C}}^{-1} \underline{\mathbf{D}}^{\pi}}_{\tilde{\mathbf{P}}^{\pi}})^{-1} \underbrace{\underline{\mathbf{C}}^{-1} \underline{\mathbf{b}}}_{\tilde{\mathbf{r}}^{\pi}} = (\underline{\mathbf{I}} - \gamma \tilde{\mathbf{P}}^{\pi})^{-1} \tilde{\mathbf{r}}^{\pi}$$

- equivalent to empirically estimated model-based solution

$$\tilde{\mathbf{P}}^{\pi} = \underline{\mathbf{C}}^{-1} \underline{\mathbf{D}}^{\pi} = \begin{bmatrix} 0 & .3 & .7 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}, \quad \tilde{\mathbf{r}} = \underline{\mathbf{C}}^{-1} \underline{\mathbf{b}} = \begin{bmatrix} .15 \\ -1 \\ +1 \\ 0 \end{bmatrix}$$

- in the limit convergence to V^{π} for ergodic Markov chains

- $\tilde{\mathbf{P}}^{\pi} \rightarrow \mathbf{P}^{\pi}$ and $\tilde{\mathbf{r}}^{\pi} \rightarrow \mathbf{r}^{\pi}$ if all states are visited infinitely often

Comparison of batch and online value estimation

■ different reward propagation

TD(0): one time step into the past

TD(λ): all λ -discounted steps

batch: instantaneous everywhere

■ different convergence behavior

TD(0): fluctuates around V^π

TD(λ): fluctuates around V^π

batch: converges exactly to V^π

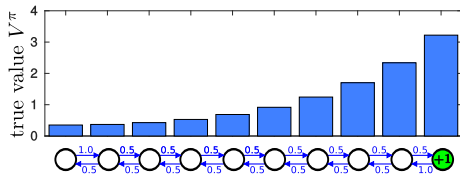
■ different complexities (computational and memory)

TD(0): $\mathcal{O}(p)$ and $\mathcal{O}(S)$

TD(λ): $\mathcal{O}(pS)$ and $\mathcal{O}(S)$

batch: $\mathcal{O}(p + S^3)$ and $\mathcal{O}(S^2)$

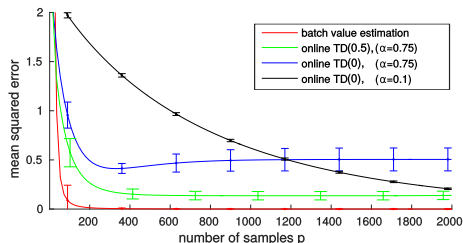
S : number of states, p : number of samples



■ Markov chain back and forth

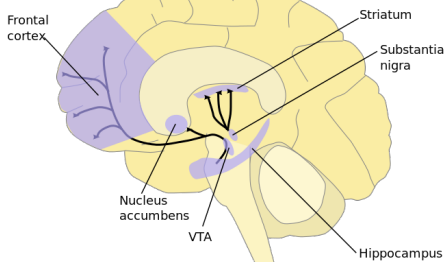
■ only last state rewarded

■ value estimated for $\gamma = 0.95$

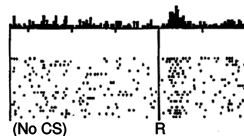


Neurological relevance of reinforcement learning

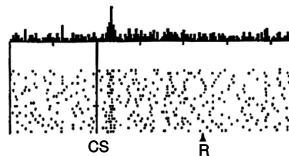
- dopamine neurons encode online TD-errors in most mammals



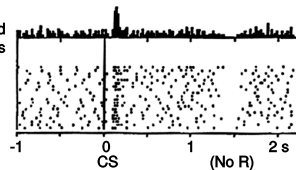
No prediction
Reward occurs



Reward predicted
Reward occurs



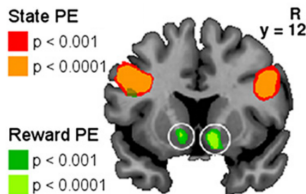
Reward predicted
No reward occurs



(Schultz et al., 1997)

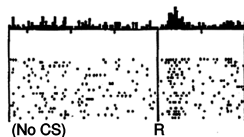
Neurological relevance of reinforcement learning

- dopamine neurons encode online TD-errors in most mammals
- model-based prediction errors were found in human *pre-frontal cortex*

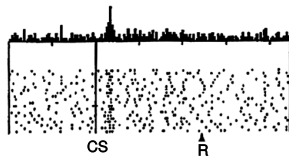


(Gläscher et al., 2010)

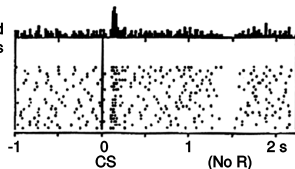
No prediction
Reward occurs



Reward predicted
Reward occurs



Reward predicted
No reward occurs



(Schultz et al., 1997)

End of Section 4.1

the following slides contain

OPTIONAL MATERIAL

The many faces of classical conditioning

FORWARD CONDITIONING



SIMULTANEOUS CONDITIONING



SECOND ORDER CONDITIONING



TEMPORAL CONDITIONING



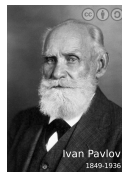
EXTINCTION



BLOCKING



INHIBITION



Contraction properties of TD learning

- asynchronous TD update at time t for all states $\underline{\mathbf{x}}_i$:

- let $v_t := \underline{\mathbf{v}}^\top \underline{\mathbf{x}}^{(t)}$ and $\mu_{it} = \underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}^{(t)}$

$$\hat{B}_t^\pi[\underline{\mathbf{v}}]_i := v_i + \underbrace{\eta \mu_{it} (\textcolor{green}{r}_t + \gamma \textcolor{blue}{v}_{t+1} - v_t)}_{\text{TD-error } \Delta v_t \text{ if } \underline{\mathbf{x}}_i = \underline{\mathbf{x}}^{(t)}}$$

- \hat{B}_t^π is *in general* a **non-expansion**

$$\max_{1 \leq i \leq S} |\hat{B}_t^\pi[\underline{\mathbf{v}}]_i - \hat{B}_t^\pi[\underline{\mathbf{w}}]_i| \leq \max_{1 \leq i \leq S} |v_i - w_i|$$

- \hat{B}_t^π is *sometimes* a **contraction mapping**

- in states $\underline{\mathbf{x}}^{(t)}$ with $|v_t - w_t| \geq \max_{i \neq t} |v_i - w_i|$

$$|\hat{B}_t^\pi[\underline{\mathbf{v}}]_t - \hat{B}_t^\pi[\underline{\mathbf{w}}]_t| \leq (1 - \eta(1 - \gamma)) |v_t - w_t|$$

Temporal difference learning with eligibility traces: TD(λ)

The TD(λ) algorithm

for $t \in \{0, \dots, p-1\}$ **do**

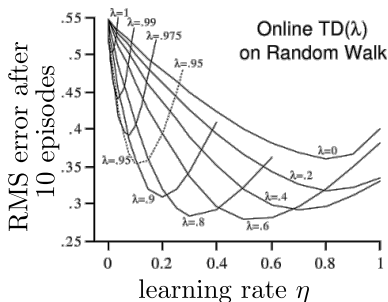
$\Delta v \leftarrow r_t + \gamma \mathbf{v}^\top \underline{\mathbf{x}}^{(t+1)} - \mathbf{v}^\top \underline{\mathbf{x}}^{(t)}$

$\mathbf{v} \leftarrow \mathbf{v} + \eta \Delta v \mathbf{e}$

$\mathbf{e} \leftarrow \gamma \lambda \mathbf{e} + \underline{\mathbf{x}}^{(t+1)}$

end

// TD-error Δv at time t
 // update all visited states
 // update eligibility trace \mathbf{e}



- TD(0): TD learning as defined before

RMS averaged over 100 random-walks
on a 19-state chain, rewarded at one end

(Sutton and Barto, 1998)

References I

- D. P. Bertsekas. *Dynamic Programming and Optimal Control*, volume 2. Athena Scientific, 3rd edition, 2007.
- Peter Dayan. The role of value systems in decision making. In *Better Than Conscious?: Decision Making, the Human Mind, and Implications for Institutions*. MIT Press, 2008. ISBN 0-262-19580-1.
- Peter Dayan and Yael Niv. Reinforcement learning: The good, the bad and the ugly. *Current Opinion in Neurobiology*, 18:185–196, 2008.
- Jan Gläscher, Nathaniel Daw, Peter Dayan, and John P. O'Doherty. States versus rewards: Dissociable neural prediction error signals underlying model-based and model-free reinforcement learning. *Neuron*, 66:585–595, 2010.
- Wolfram Schultz, Peter Dayan, and P. Read Montague. A Neural Substrate of Prediction and Reward. *Science*, 275(5306):1593–1599, March 1997. ISSN 1095-9203.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.