

Policy and Value Iteration

Last Time

Last Time

- How is a **Markov decision process** defined?

Last Time

- How is a **Markov decision process** defined?
- What is a **policy**?

Last Time

- How is a **Markov decision process** defined?
- What is a **policy**?
- How do we **evaluate** policies?

Last Time

- How is a **Markov decision process** defined?
- What is a **policy**?
- How do we **evaluate** policies?

(MDP notebook)

Guiding Questions

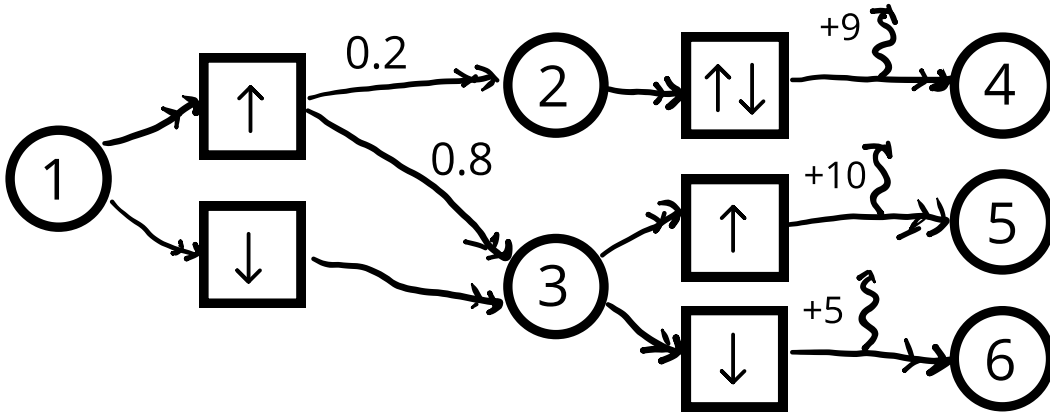
Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

Value-Based Policy Evaluation

MDP Example: Up-Down Problem

From *Dynamic Programming and the Calculus of Variations*, 1965



Algorithm: Bellman Backup

Given: MDP $(S, A, R, T, S_T, \gamma)$

1. $U^*(s) \leftarrow 0 \quad \forall s \in S_T$
2. Repeat until $U^*(s)$ known for all states:
 1. Choose s where U^* is known for all children
 2. Calculate $U^*(s)$
3. Extract $\pi^*(s) = \operatorname{argmax} Q^*(s, a)$

Break: DIA Run

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$
6. return π

Policy Iteration

Algorithm: Policy Iteration

Given: MDP (S, A, R, T, γ, b)

1. initialize π, π' (differently)
2. while $\pi \neq \pi'$
3. $\pi \leftarrow \pi'$
4. $U^\pi \leftarrow (I - \gamma T^\pi)^{-1} R^\pi$
5. $\pi'(s) \leftarrow \operatorname{argmax}_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a) U^\pi(s')) \quad \forall s \in S$
6. return π

(Policy iteration notebook)

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_{\infty} > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$
5. return U'

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_\infty > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$
5. return U'

- Returned U' will be close to U^* !

Value Iteration

Algorithm: Value Iteration

Given: MDP (S, A, R, T, γ, b) , tolerance ϵ

1. initialize U, U' (differently)
2. while $\|U - U'\|_\infty > \epsilon$
3. $U \leftarrow U'$
4. $U'(s) \leftarrow \max_{a \in A} (R(s, a) + \gamma \sum_{s' \in S} T(s'|s, a)U(s')) \quad \forall s \in S$
5. return U'

- Returned U' will be close to U^* !
- π^* is easy to extract: $\pi^*(s) = \arg \max (R(s, a) + \gamma E[U^*(s)])$

Bellman's Equations

Guiding Questions

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

Guiding Questions

- How do we reason about the **future consequences** of actions in an MDP?
- What are the basic **algorithms for solving MDPs**?

"In any small change he will have to consider only these quantitative indices (or "values") in which all the relevant information is concentrated; and by adjusting the quantities one by one, he can appropriately rearrange his dispositions without having to solve the whole puzzle ab initio, or without needing at any stage to survey it at once in all its ramifications."

-- F. A. Hayek, "The use of knowledge in society", 1945