

Reinforcement Learning

Lecture 3

Planning

Closed Form Solution of Value Function

Iterative Policy Evaluation

Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology
Winter Term 2023/24

Recap

A Markov Decision Process (MDP) is a 3-tuple

$$(S, A, p)$$

where

S is a state space (set of states of environment)

A is an action space (set of actions of agent)

p is called the dynamics of the MDP. Formally, it is a conditional probability distribution

$$p(s', r | s, a)$$

where

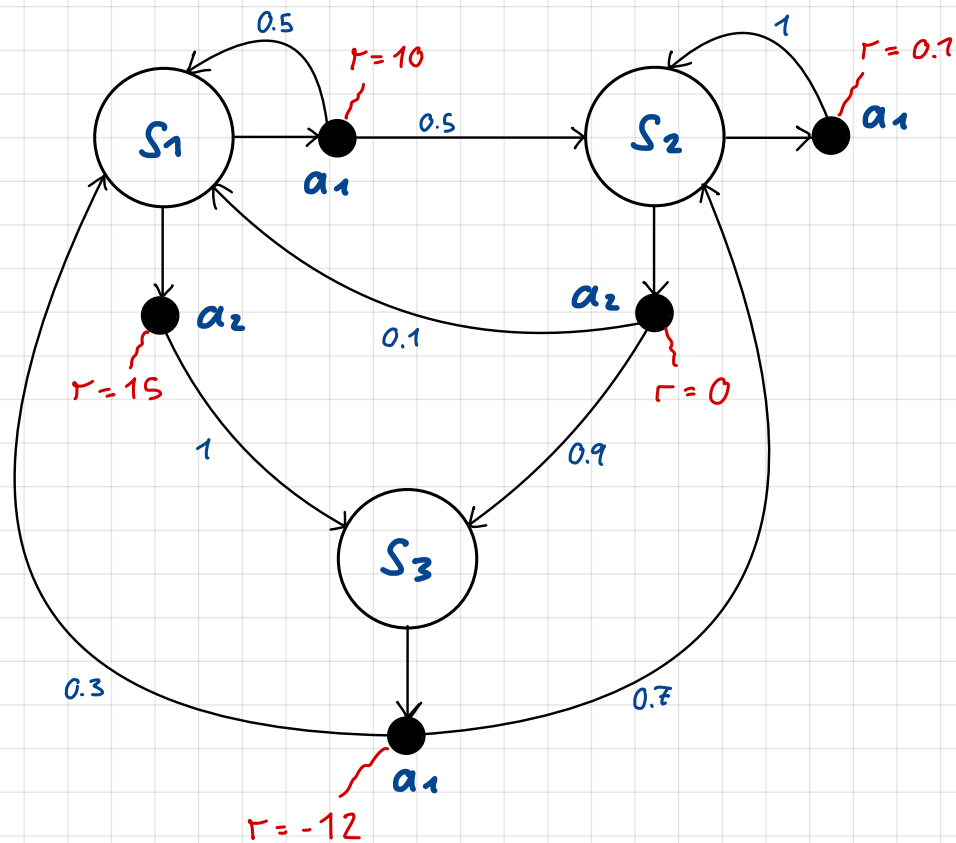
$s' \in S$ is the value of the next state S_{t+1} ,

$r \in \mathbb{R}$ is the value of the reward R_{t+1} ,

$s \in S$ is the value of the current state S_t ,

$a \in A$ is the value of the selected action A_t .

Recap



Transition Diagram



1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

Grid World

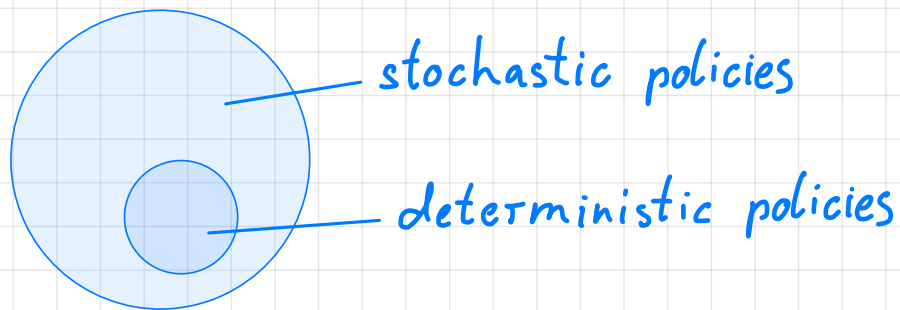
Recap

Policy (agent's behavior)

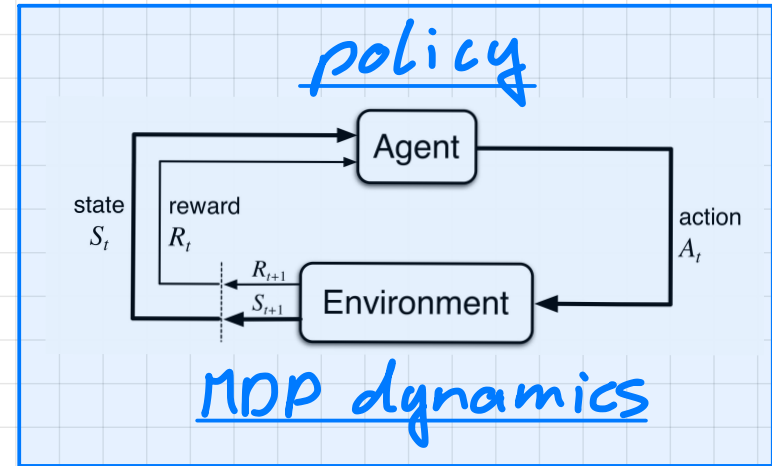
$$\hat{\pi}(a|s)$$

(conditional distribution over actions given current state)

Subsumes deterministic policies



$$\hat{\pi}_{sto}(a|s) := \begin{cases} 1 & \text{if } \hat{\pi}_{det}(s) = a \\ 0 & \text{otherwise} \end{cases}$$



Recap

The discounted return at time t is defined as

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad 0 \leq \gamma \leq 1$$

The value function $V_{\pi}: \mathcal{S} \rightarrow \mathbb{R}$ is defined as

$$V_{\pi}(s) := \mathbb{E}_{\pi}[G_t \mid S_t = s] = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

"The expected discounted return when following policy π , and starting from $s \in \mathcal{S}$ (t is any arbitrary time step)."

Bellman Expectation Equation

$$\underline{V_{\pi}(s)} = \mathbb{E}_{A_t, R_{t+1}, S_{t+1}} \left[R_{t+1} + \gamma \underline{V_{\pi}(S_{t+1})} \mid S_t = s \right]$$



Image: wikipedia.org

Recap

$$\underline{v_{\pi}(s)} = \mathbb{E}_{A_t, R_{t+1}, S_{t+1}} \left[R_{t+1} + \gamma \underline{v_{\pi}(S_{t+1})} \mid S_t = s \right]$$

$$\begin{aligned} v_{\pi}(s) &= \sum_a \pi(a|s) \left[\sum_{s', r} p(s', r | s, a) (r + \gamma v_{\pi}(s')) \right] \\ &\vdots \\ &= r(s) + \gamma \sum_{s'} p(s'|s) v_{\pi}(s') \end{aligned}$$

where

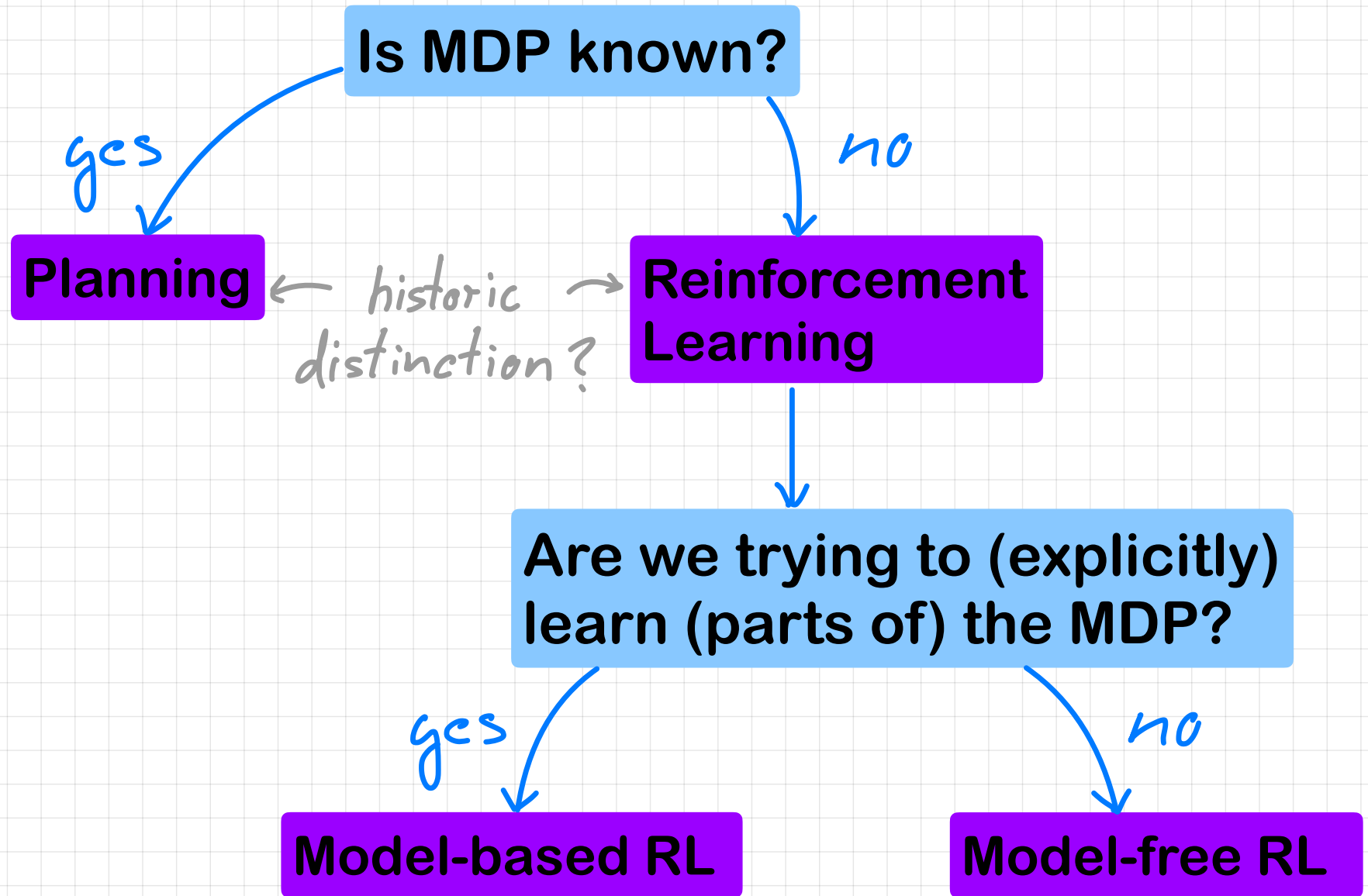
starting at s expected reward starting at s with a

$$r(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) r = \sum_a \pi(a|s) r(s, a)$$

combined state
transition

$$p(s'|s) = \sum_a \sum_r \underbrace{\pi(a|s)}_{\text{policy}} \underbrace{p(s', r | s, a)}_{\text{dynamics}} \quad p(s', a, r | s) \text{ (chain rule)}$$

Planning, Model-based RL, Model-free RL



Fundamental Problems in Planning & RL

Policy Evaluation, Prediction – How good is my policy?

given a policy π , compute the value function v_π

Control – What is the best policy?

How to improve a given policy π , or even find an optimal policy?

The two problems are related, and the algorithms to solve them are similar

Planning: Policy Evaluation, Prediction

Solving the Bellman Equation(s)

- note that we actually have one Bellman equation per state:

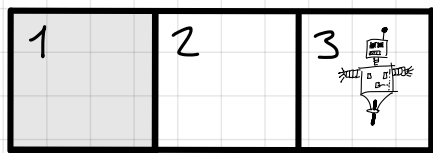
$$V_{\pi}(s_1) = r(s_1) + \gamma \sum_{s'} p(s'|s_1) V_{\pi}(s')$$

$$V_{\pi}(s_2) = r(s_2) + \gamma \sum_{s'} p(s'|s_2) V_{\pi}(s')$$

\vdots

$$V_{\pi}(s_N) = r(s_N) + \gamma \sum_{s'} p(s'|s_N) V_{\pi}(s') \quad (\text{assuming } N \text{ states})$$

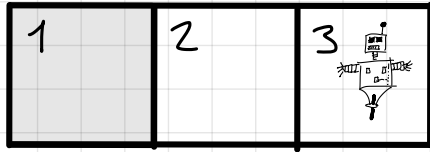
- assume a very simple grid world



- states $\mathcal{S} = \{1, 2, 3\}$, state 1 is terminal
- actions $\mathcal{A} = \{\leftarrow, \rightarrow\}$, move agent (deterministic)
- reward -1 non-terminal states, 0 for terminal

- assume a random walk policy $\pi(a|s) = 0.5$ for all s, a

Solving the Bellman Equation(s) cont'd



- states $\mathcal{S} = \{1, 2, 3\}$, state 1 is terminal
- actions $\mathcal{A} = \{\leftarrow, \rightarrow\}$, move agent (deterministic)
- reward -1 non-terminal states, 0 for terminal

combined state transition as matrix

$|\mathcal{S}| \times |\mathcal{S}|$ -matrix

$$P = \begin{pmatrix} p(1|1) & p(2|1) & p(3|1) \\ p(1|2) & p(2|2) & p(3|2) \\ p(1|3) & p(2|3) & p(3|3) \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0.5 & 0 & 0.5 \\ 0 & 0.5 & 0.5 \end{pmatrix}$$

stochastic matrix

(non-negative,
rows sum to 1)

$$\begin{aligned} \text{e.g. } p(1|2) &= \tilde{\pi}(\leftarrow|2) p(1|2, \leftarrow) + \tilde{\pi}(\rightarrow|2) p(1|2, \rightarrow) \\ &= 0.5 * 1 + 0.5 * 0 = 0.5 \end{aligned}$$

expected immediate reward as vector

$$r = \begin{pmatrix} 0 \\ -1 \\ -1 \end{pmatrix}$$

$|\mathcal{S}|$ -dim vector

Solving the Bellman Equation(s) cont'd

$$P = \begin{pmatrix} 1 & 0 & 0 \\ 0.5 & 0 & 0.5 \\ 0 & 0.5 & 0.5 \end{pmatrix} \quad r = \begin{pmatrix} 0 \\ -1 \\ -1 \end{pmatrix}$$

Bellman equation for state 1

$$\begin{aligned} v_{\pi}(1) &= r(1) + \gamma \left[p(1|1) v_{\pi}(1) + p(2|1) v_{\pi}(2) + p(3|1) v_{\pi}(3) \right] \\ &= 0 + \gamma \left[1 v_{\pi}(1) + 0 v_{\pi}(2) + 0 v_{\pi}(3) \right] \\ &= 0 + \gamma v_{\pi}(1) \end{aligned}$$

$$(1 - \gamma) v_{\pi}(1) = 0$$

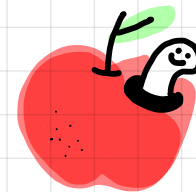
Thus, if $\gamma < 1$ (discounting) $v_{\pi}(1) = 0$

For $\gamma = 1$, Bellman equation is vacuous here...

Bellman equation is not the definition of v_{π} !

Even for $\gamma = 1$ we can argue that $v_{\pi}(1) = 0$, since

$$v_{\pi}(1) = \mathbb{E}_{\pi} [G_t | s=1] = 0 \text{ as } 1 \text{ is a terminal state}$$



Solving the Bellman Equation(s) cont'd

Bellman equations for states 2 and 3

$$v_{\pi}(2) = r(2) + \gamma \left[\cancel{p(1|2) v_{\pi}(1)} + \cancel{p(2|2) v_{\pi}(2)} + p(3|2) v_{\pi}(3) \right]$$

$$v_{\pi}(3) = r(3) + \gamma \left[\cancel{p(1|3) v_{\pi}(1)} + p(2|3) v_{\pi}(2) + p(3|3) v_{\pi}(3) \right]$$

$= 0$

$$v_{\pi}(2) = -1 + \gamma \cdot 0.5 v_{\pi}(3)$$

$$v_{\pi}(3) = -1 + \gamma \cdot 0.5 v_{\pi}(2) + \gamma \cdot 0.5 v_{\pi}(3)$$

2 equations

2 unknowns



Assume $\gamma = 1$

$$v_{\pi}(3) = -1 + 0.5 \left(\overbrace{-1 + 0.5 v_{\pi}(3)}^{v_{\pi}(2)} \right) + 0.5 v_{\pi}(3)$$

$$= -1.5 + 0.75 v_{\pi}(3)$$

$$v_{\pi}(3) = \frac{-1.5}{0.25} = -6$$

$$v_{\pi}(2) = -1 + 0.5 v_{\pi}(3) = -4$$

Analytic Solution of Bellman Equation

$$p(s'|s) = \sum_{r,a} p(s', r|s, a) \tilde{\pi}(a|s)$$

$$\underline{V_{\pi}}(s) = r(s) + \gamma \sum_{s'} p(s'|s) V_{\pi}(s')$$

Assume an arbitrary enumeration of states $\mathcal{S} = \{s_1, s_2, \dots, s_N\}$

$$\underline{V} = \begin{pmatrix} V_{\pi}(s_1) \\ V_{\pi}(s_2) \\ \vdots \\ V_{\pi}(s_N) \end{pmatrix} \quad \underline{r} = \begin{pmatrix} r(s_1) \\ r(s_2) \\ \vdots \\ r(s_N) \end{pmatrix} \quad P = \begin{pmatrix} p(s_1|s_1) & \dots & p(s_N|s_1) \\ \vdots & \ddots & \vdots \\ p(s_1|s_N) & \dots & p(s_N|s_N) \end{pmatrix}$$

We can write the Bellman equations as $\underline{V} = \underline{r} + \gamma P \underline{V}$

Thus,

$$\underline{V} - \gamma P \underline{V} = \underline{r}$$

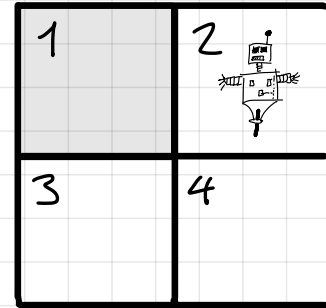
$$(I - \gamma P) \underline{V} = \underline{r}$$

$$\underline{V_{\pi}} = (I - \gamma P)^{-1} \underline{r}$$

Example: 4 State Grid World

$$P = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.25 & 0.5 & 0 & 0.25 \\ 0.25 & 0 & 0.5 & 0.25 \\ 0 & 0.25 & 0.25 & 0.5 \end{pmatrix} \quad r = \begin{pmatrix} 0 \\ -1 \\ -1 \\ -1 \end{pmatrix} \quad \pi(a|s) = 0.25$$

$$\gamma = 0.9999$$



```
>>> import numpy as np
>>> P = np.array([[1., 0., 0., 0.], [0.25, 0.5, 0., 0.25], [0.25, 0., 0.5, 0.25], [0., 0.25, 0.25, 0.5]])
>>> P
array([[1., 0., 0., 0.],
       [0.25, 0.5, 0., 0.25],
       [0.25, 0., 0.5, 0.25],
       [0., 0.25, 0.25, 0.5]])
>>> r = np.array([0., -1., -1., -1.])
>>> r
array([ 0., -1., -1., -1.])
>>> gamma = 0.9999
>>> np.linalg.inv(np.eye(4) - gamma*P) @ r
array([-2.25605909e-16, -5.99660198e+00, -5.99660198e+00, -7.99520280e+00])
>>>
```

$$\underline{V} = (\underline{I} - \gamma P)^{-1} \underline{r} \approx \begin{pmatrix} 0 \\ -5.99 \\ -5.99 \\ -7.99 \end{pmatrix}$$

double check

\underline{V} must satisfy Bellman equation

$$\underline{V} = \underline{r} + \gamma P \underline{V}$$

1	2
0	-5.99
3	4
-5.99	-7.99

Uniqueness of Solution

Bellman Equation (three different notations)

$$V_{\pi}(s) = E_{A_t, R_{t+1}, S_{t+1}} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s]$$

$$V_{\pi}(s) = r(s) + \sum_{s'} p(s'|s) \gamma V_{\pi}(s')$$

$$\underline{V} = \underline{r} + \gamma P \underline{V}$$

Is there always a (unique) solution?

- For $\gamma=1$, $(I-\gamma P)$ is singular and the solution is not unique
- For $\gamma < 1$, $(I-\gamma P)$ is always invertible — solution unique
- Thus, discounting makes the problem well-posed

Planning: Iterative Policy Evaluation

Iterative Policy Evaluation

- solving the Bellman equation via matrix inversion does not scale to large MDPs (cubic complexity)
- also not applicable to later problems, e.g. control
- iterative method?
- reconsider the Bellman equation:

$$V_{\pi}(s) = r(s) + \gamma \sum_{s'} p(s'|s) V_{\pi}(s')$$

- this looks like an update rule!

Iterative Policy Evaluation

- initialize $V(s)$ arbitrarily (e.g. all 0)
- repeat
 - for $s \in \mathcal{S}$ do $V_{\text{new}}(s) \leftarrow r(s) + \gamma \sum_{s'} p(s'|s) V(s')$
 - if $\forall s: V_{\text{new}}(s) \approx V(s) \rightarrow$ break
 - $V \leftarrow V_{\text{new}}$

Iterative Policy Evaluation

$$\gamma = 0.9999$$

iteration

V

0

1	2
0	0
3	4
0	0

$$P = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0.25 & 0.5 & 0 & 0.25 \\ 0.25 & 0 & 0.5 & 0.25 \\ 0 & 0.25 & 0.25 & 0.5 \end{pmatrix} \quad r = \begin{pmatrix} 0 \\ -1 \\ -1 \\ -1 \end{pmatrix}$$

1	2
3	4

1

1	2
0	-1
3	4
-1	-1

$$\underline{V} = \underline{r} + \gamma P \underline{V}$$

2

1	2
0	-2.4
3	4
-2.4	-2.9

Recall:

$$\underline{V} = (I - \gamma P)^{-1} \underline{r} \approx \begin{pmatrix} 0 \\ -5.99 \\ -5.99 \\ -7.99 \end{pmatrix}$$

⋮

⋮

10

1	2
0	-4.8
3	4
-4.8	-6.3

⋮

20

1	2
0	-5.8
3	4
-5.8	-7.6

⋮

50

1	2
0	-5.99
3	4
-5.99	-7.99

Does iterative policy evaluation always converge to V_{π} ?

Iterative Policy Evaluation, Larger Gridworld

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

- actions \leftrightarrow move agent deterministically
- reward -1 whenever on non-terminal (white) cell
- discount factor $\gamma = 0.999$
- policy: $\pi(a|s) = 0.25$ (random walk)

After K iterations:

$K=0$

0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

$K=1$

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

$K=2$

0	-1.7	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

...

$K=10$

0	-6.1	-8.3	-8.9
-6.1	-7.7	-8.4	-8.3
-8.3	-8.4	-7.7	-6.1
-8.9	-8.3	-6.1	0

...

$K=199$

0	-13.8	-19.6	-21.6
-13.8	-17.7	-19.6	-19.6
-19.6	-19.6	-17.7	-13.8
-21.6	-19.6	-13.8	0

$K=200$

0	-13.8	-19.6	-21.6
-13.8	-17.7	-19.6	-19.6
-19.6	-19.6	-17.7	-13.8
-21.6	-19.6	-13.8	0

$\approx V_{\pi}$

Some Mathematical Tools

Lipschitz function

Let V be any vector space and $T: V \rightarrow V$ be a function from V to V . If there is a constant L such that

$$\|T(v) - T(v')\| \leq L \|v - v'\| \quad \forall v, v' \in V$$

then T is called L -Lipschitz, where $\|\cdot\|$ is any norm.

Contraction (Contraction mapping)

Let V be any vector space. An L -Lipschitz function $T: V \rightarrow V$ with $L < 1$ is called a contraction (L -contraction).

Contraction

For any two distinct points u, v :

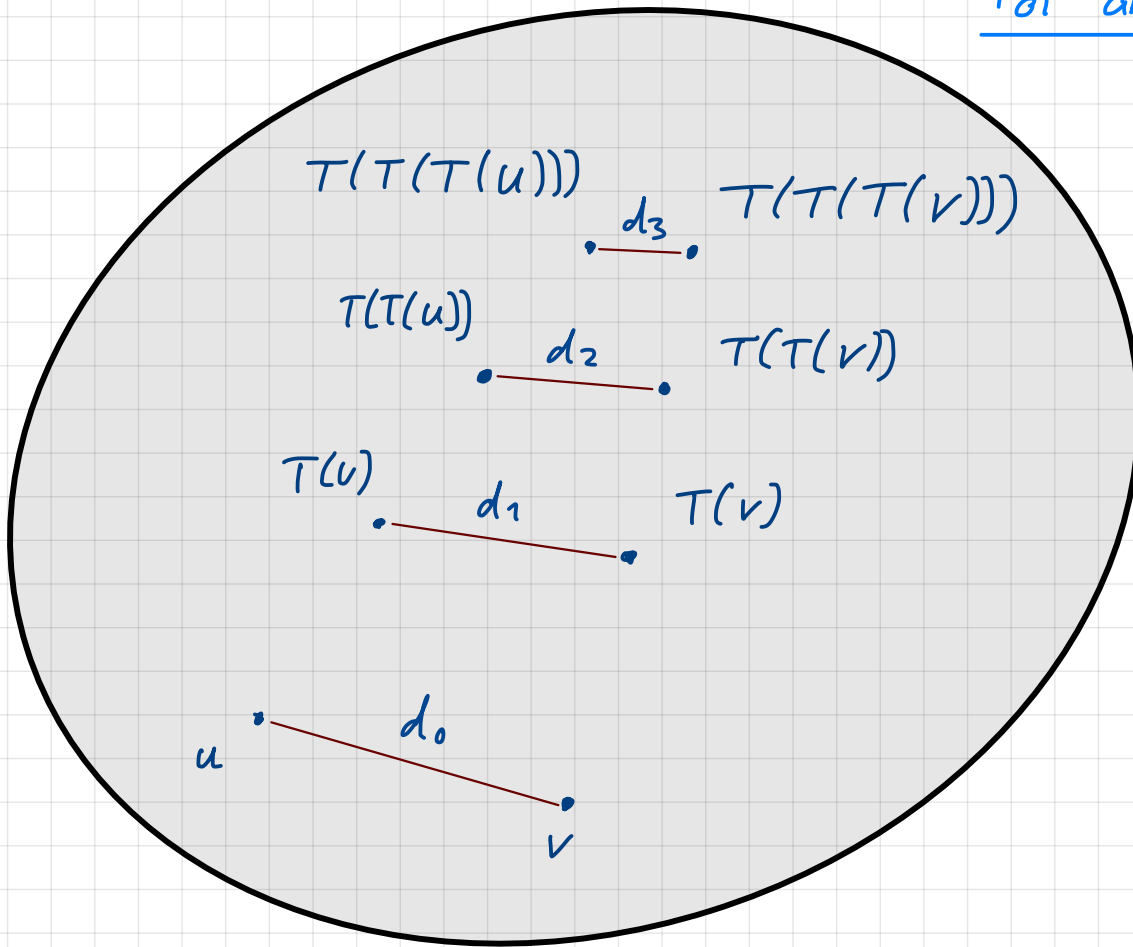
$$d_0 = \|u - v\|$$

$$d_1 = \|T(u) - T(v)\|$$

$$d_2 = \|T(T(u)) - T(T(v))\|$$

\vdots

$$d_0 > d_1 > d_2 > d_3 > \dots$$



Banach's Fixed Point Theorem

Fixed point

Let V be some vector space and let $T: V \rightarrow V$ be a function. A vector $v \in V$ is called a fixed point of T if $v = T(v)$, i.e. if T maps v onto itself.

Let V be some vector space and let $T: V \rightarrow V$ be a contraction. Then T has a unique fixed point $v^* \in V$. Moreover, for any $v_0 \in V$, the sequence defined via

$$v_n = T(v_{n-1})$$

converges to v^* , i.e. $v_n \rightarrow v^*$ as $n \rightarrow \infty$.

Details

Norm: We tacitly assumed some norm $\|\cdot\|$ of vector space V .
A common choice is the max norm (uniform norm, infinity norm)

$$\|\underline{v}\|_{\infty} = \max_i |v_i|$$

Complete space: The fixed-point theorem requires V to be complete, i.e. that every Cauchy sequence converges to some element in V .

For example, the real numbers \mathbb{R} are complete.

The rational numbers \mathbb{Q} are not complete; e.g., there are sequences in \mathbb{Q} which converge to a irrational number (e.g. π).

A complete vector space with a norm is a Banach space.

Thus, the Banach fixed-point theorem says that iterating a contraction T in a Banach space always converges to a unique fixed-point.

Iterative Policy Evaluation

For $\gamma < 1$, Iterative Policy Evaluation always converges to v_{π} ,
i.e. for any MDP, any π and any initialisation of v_{π} .

- for $\gamma < 1$, the Bellman Equation is a contraction.
(see HW1 in KU)
- thus, iterating it will converge to a unique fixed-point
- v_{π} is a fixed-point, thus Iterative Policy Evaluation converges to v_{π}

Contraction arguments like this are frequently used in RL theory!

Summary

Planning: Policy Evaluation, Prediction

- We have learned about two methods for evaluating a policy π
 - closed form

$$\underline{v}_{\pi} = (I - \gamma P)^{-1} \underline{r}$$

- iterative policy evaluation

init \underline{v}_0

iterate $\underline{v}_n = \underline{r} + \gamma P \underline{v}_{n-1}$

- Policy evaluation just tells us how good π is, for each s .
- But how to actually learn π ?
- Or, given some π , how to improve it? (next lecture)