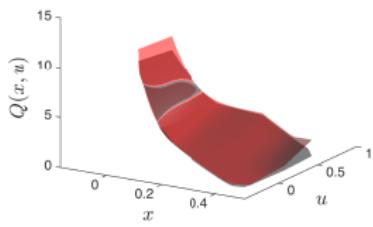
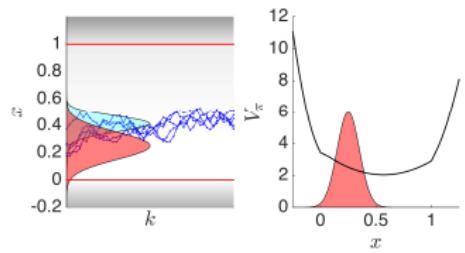
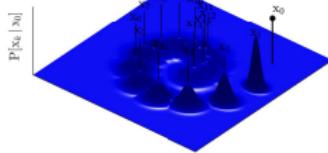


What are we going to discuss this week?

samples = 1000000



$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

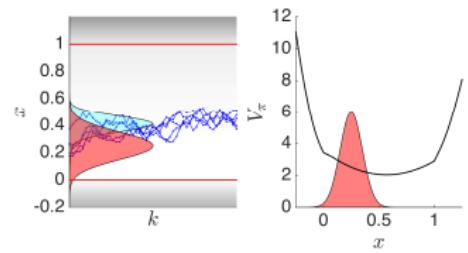
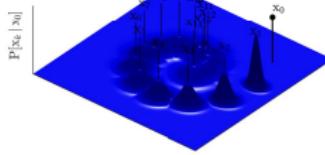
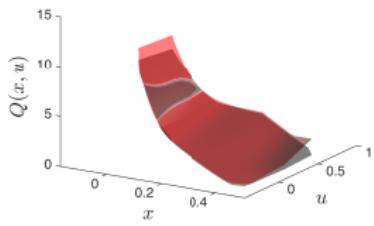


What are we going to discuss this week?

- ➊ Optimal control on (almost) anything: Markov Decision Process & Bellman equations

samples = 1000000

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

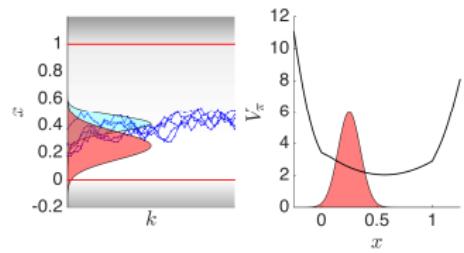
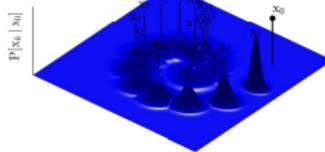
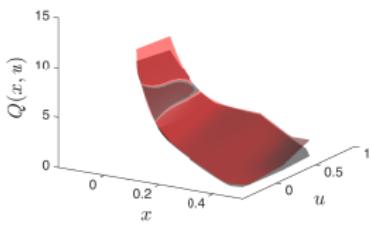


What are we going to discuss this week?

- ➊ Optimal control on (almost) anything: Markov Decision Process & Bellman equations
- ➋ Optimal control without models: Q-learning

samples = 1000000

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

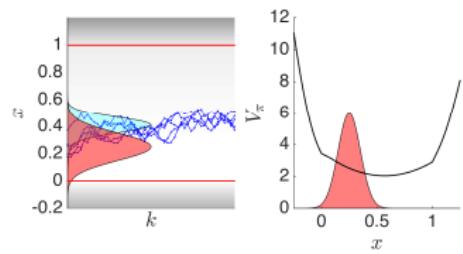
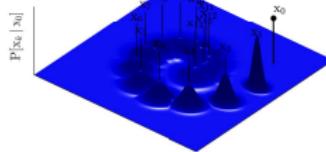
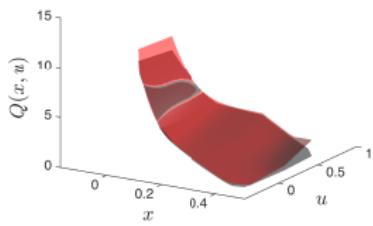


What are we going to discuss this week?

- ➊ Optimal control on (almost) anything: Markov Decision Process & Bellman equations
- ➋ Optimal control without models: Q -learning
- ➌ Optimal control without models in practice: function approximation

samples = 1000000

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

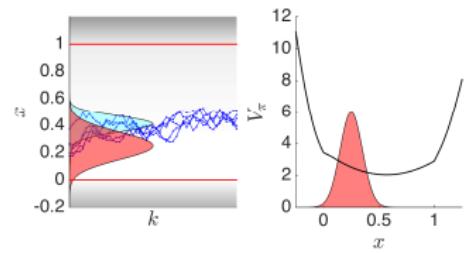
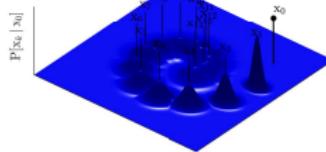
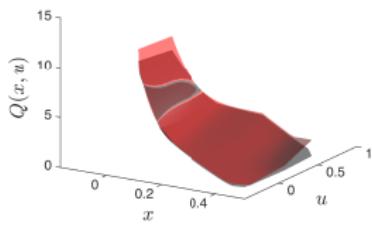


What are we going to discuss this week?

- ➊ Optimal control on (almost) anything: Markov Decision Process & Bellman equations
- ➋ Optimal control without models: Q -learning
- ➌ Optimal control without models in practice: function approximation
- ➍ Optimizing control policies: policy gradient methods

samples = 1000000

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

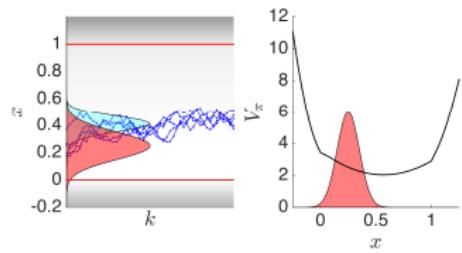
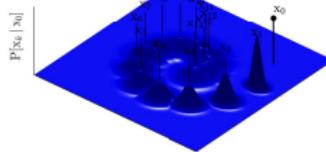
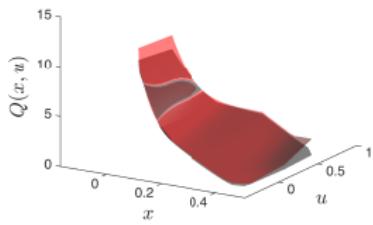


What are we going to discuss this week?

- ➊ Optimal control on (almost) anything: Markov Decision Process & Bellman equations
- ➋ Optimal control without models: Q -learning
- ➌ Optimal control without models in practice: function approximation
- ➍ Optimizing control policies: policy gradient methods
- ➎ Merging MPC and Learning: outlook on current research

samples = 1000000

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$



What are we going to discuss today?

“Classic” optimal control for given initial state s :

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{u}} \quad & \int_0^T L(\mathbf{x}, \mathbf{u}) dt \\ \text{s.t.} \quad & \mathbf{f}(\dot{\mathbf{x}}, \mathbf{x}, \mathbf{u}) = 0 \\ & \mathbf{h}(\mathbf{x}, \mathbf{u}) \leq 0 \\ & \mathbf{x}(0) = \mathbf{s} \end{aligned}$$

Can we define & treat optimality for dynamic systems in a more abstract way?

What are we going to discuss today?

"Classic" optimal control for given initial state s :

$$\begin{aligned} \min_{\mathbf{x}, \mathbf{u}} \quad & \int_0^T L(\mathbf{x}, \mathbf{u}) dt \\ \text{s.t.} \quad & \mathbf{f}(\dot{\mathbf{x}}, \mathbf{x}, \mathbf{u}) = 0 \\ & \mathbf{h}(\mathbf{x}, \mathbf{u}) \leq 0 \\ & \mathbf{x}(0) = \mathbf{s} \end{aligned}$$

Can we define & treat optimality for dynamic systems in a more abstract way?

We want to address:

- Any set/space of state & input (functional, discrete, games, etc.)
- Any kind of dynamics (stochastic, mixed-integer, etc.)
- Any kind of stage cost (discontinuous, taking infinite values, stochastic, etc.)
- Infinite-horizon problems, state-dependent termination (game, target)
- Obtain input policies as opposed to input profiles

Markov Decision Processes: the Bellman equations and Dynamic Programming

Sébastien Gros

Cybernetic, NTNU
Elec. Eng., Chalmers

TUM lectures on RL

Outline

- 1 Markov Decision Process
- 2 Value functions & Bellman equations
- 3 Function evaluation & Dynamic Programming
- 4 Quick view of Bellman contraction

Notation & assumptions

- We will work in discrete time for simplicity

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)
- Classic understanding applies but no restriction, e.g.:
 - ▶ x, u can be discrete (e.g. position and movement in a game)
 - ▶ x, u can belong to abstract vector spaces (e.g. functions)

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)
- Classic understanding applies but no restriction, e.g.:
 - ▶ x, u can be discrete (e.g. position and movement in a game)
 - ▶ x, u can belong to abstract vector spaces (e.g. functions)
- We will often use the objects x, u, x_+ together, meaning

$$\underbrace{x, u}_{\text{current state \& input}} \longrightarrow \underbrace{x_+}_{\text{next state}}$$

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)
- Classic understanding applies but no restriction, e.g.:
 - ▶ x, u can be discrete (e.g. position and movement in a game)
 - ▶ x, u can belong to abstract vector spaces (e.g. functions)
- We will often use the objects x, u, x_+ together, meaning

$$\underbrace{x, u}_{\text{current state \& input}} \longrightarrow \underbrace{x_+}_{\text{next state}}$$

- Stochastic: at current time, x_+ is typically unknown, even if x, u are known

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)
- Classic understanding applies but no restriction, e.g.:
 - ▶ x, u can be discrete (e.g. position and movement in a game)
 - ▶ x, u can belong to abstract vector spaces (e.g. functions)
- We will often use the objects x, u, x_+ together, meaning

$$\underbrace{x, u}_{\text{current state \& input}} \longrightarrow \underbrace{x_+}_{\text{next state}}$$

- Stochastic: at current time, x_+ is typically unknown, even if x, u are known
- We can consider ∞ and finite processes, termination time may be unknown. E.g.
 - ▶ Chess game: finite process terminating at an unknown time
 - ▶ Landing a plane: finite process terminating at a given time/state
 - ▶ Managing a smart building: infinite process

Notation & assumptions

- We will work in discrete time for simplicity
- State x represents the current “status” of a system
- Input u represents what we currently decide to do to the system (a.k.a. action)
- Classic understanding applies but no restriction, e.g.:
 - ▶ x, u can be discrete (e.g. position and movement in a game)
 - ▶ x, u can belong to abstract vector spaces (e.g. functions)
- We will often use the objects x, u, x_+ together, meaning

$$\underbrace{x, u}_{\text{current state \& input}} \longrightarrow \underbrace{x_+}_{\text{next state}}$$

- Stochastic: at current time, x_+ is typically unknown, even if x, u are known
- We can consider ∞ and finite processes, termination time may be unknown. E.g.
 - ▶ Chess game: finite process terminating at an unknown time
 - ▶ Landing a plane: finite process terminating at a given time/state
 - ▶ Managing a smart building: infinite process
- What do we keep from “classic” optimal control?
 - ▶ Causality
 - ▶ Performance is minimizing the sum of L

Notation

Conditional probability $\mathbb{P}[y | x]$ is a function of x and y

Probability distribution (discrete or density) of random variable y for x known, e.g.:

$$\mathbb{P}[y | x] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} \|y - \nu(x)\|^2}$$

Notation

Conditional probability $\mathbb{P}[y | x]$ is a function of x and y

Probability distribution (discrete or density) of random variable y for x known, e.g.:

$$\mathbb{P}[y | x] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} \|y - \nu(x)\|^2}$$

Expected value $\mathbb{E}[v(y)]$

For discrete distributions: $\mathbb{E}[v(y)] = \sum_{y \in D} v(y) \mathbb{P}[y]$

For densities: $\mathbb{E}[v(y)] = \int_{y \in D} v(y) \underbrace{\mathbb{P}[y]}_{=\mu_{\mathbb{P}[y]}(y)} dy$ (may need Lebesgue int.)

Notation

Conditional probability $\mathbb{P}[y | x]$ is a function of x and y

Probability distribution (discrete or density) of random variable y for x known, e.g.:

$$\mathbb{P}[y | x] = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2} \|y - \nu(x)\|^2}$$

Expected value $\mathbb{E}[v(y)]$

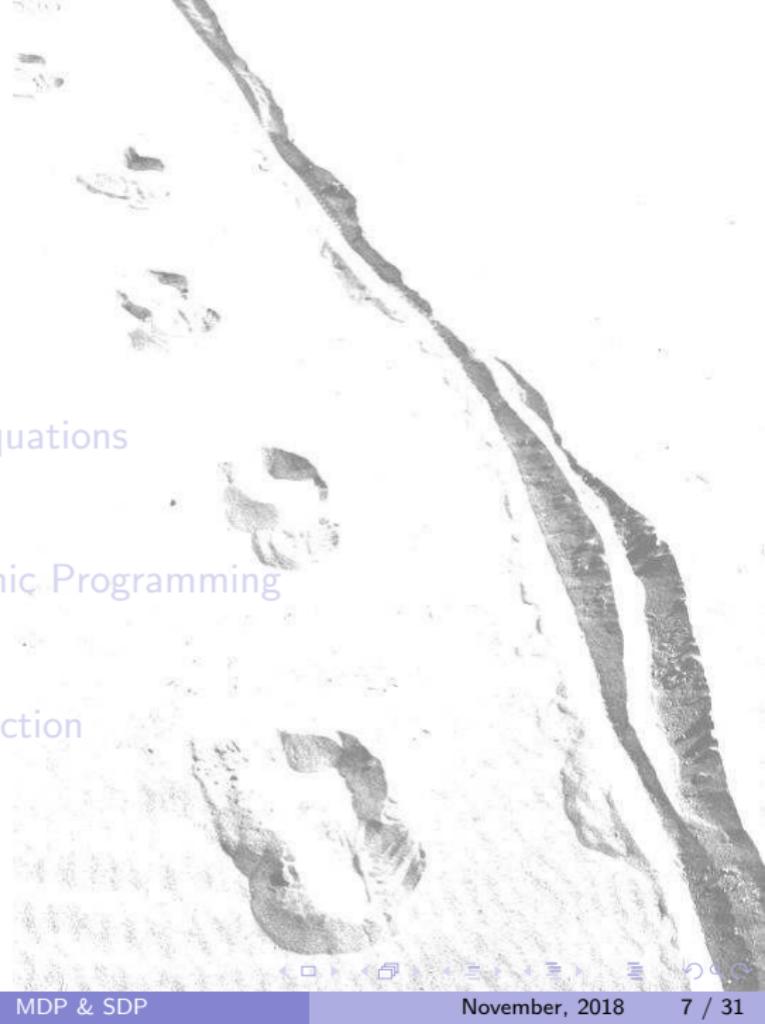
For discrete distributions: $\mathbb{E}[v(y)] = \sum_{y \in D} v(y) \mathbb{P}[y]$

For densities: $\mathbb{E}[v(y)] = \int_{y \in D} v(y) \underbrace{\mathbb{P}[y]}_{=\mu_{\mathbb{P}[y]}(y)} dy$ (may need Lebesgue int.)

Conditional expected value $\mathbb{E}[v(y) | x]$ is a function of x

For densities: $\mathbb{E}[v(y) | x] = \int_{y \in D} v(y) \mathbb{P}[y | x] dy$

Outline

- 
- 1 Markov Decision Process
 - 2 Value functions & Bellman equations
 - 3 Function evaluation & Dynamic Programming
 - 4 Quick view of Bellman contraction

State transitions - how to describe dynamics?

State transition

$$\underbrace{x, u}_{\text{current state \& input}} \longrightarrow \underbrace{x_+}_{\text{next state}}$$

State transitions - how to describe dynamics?

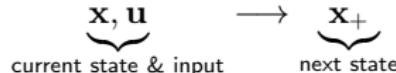
State transition

$$\underbrace{\mathbf{x}, \mathbf{u}}_{\text{current state \& input}} \longrightarrow \underbrace{\mathbf{x}_+}_{\text{next state}}$$

- In classic control, models are $\mathbf{x}_+ = \mathbf{f}(\mathbf{x}, \mathbf{u})$ for some function \mathbf{f}

State transitions - how to describe dynamics?

State transition



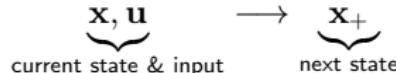
- In classic control, models are $x_+ = f(x, u)$ for some function f
- More generally state transitions are stochastic, given by conditional distributions

$$\mathbb{P}[x_+ | x, u]$$

Meaning: probability of transiting to state x_+ provided that the system is in state x and subject to input u

State transitions - how to describe dynamics?

State transition



- In classic control, models are $x_+ = f(x, u)$ for some function f
- More generally state transitions are stochastic, given by conditional distributions

$$\mathbb{P}[x_+ | x, u]$$

Meaning: probability of transiting to state x_+ provided that the system is in state x and subject to input u

- Classic model $x_+ = f(x, u)$ is a special case of stochastic transition, by defining:

$$\mathbb{P}[x_+ | x, u] = \delta(x_+ - f(x, u))$$

where δ is a Dirac function

State transitions - how to describe dynamics?

State transition



- In classic control, models are $x_+ = f(x, u)$ for some function f
- More generally state transitions are stochastic, given by conditional distributions

$$\mathbb{P}[x_+ | x, u]$$

Meaning: probability of transiting to state x_+ provided that the system is in state x and subject to input u

- Classic model $x_+ = f(x, u)$ is a special case of stochastic transition, by defining:

$$\mathbb{P}[x_+ | x, u] = \delta(x_+ - f(x, u))$$

where δ is a Dirac function

- Process noise $x_+ = f(x, u) + e$ where e.g. $e \sim \mathcal{N}(0, \Sigma)$, then

$$\mathbb{P}[x_+ | x, u] \propto e^{-\frac{1}{2}(x_+ - f(x, u))^T \Sigma (x_+ - f(x, u))}$$

Stage cost - Assessing Performance

Stage cost function: $L(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$

- $L(\mathbf{x}, \mathbf{u}) = +\infty$ for infeasible states and inputs (constraints in the cost)
- $L(\mathbf{x}, \mathbf{u})$ may not be lower-bounded

Stage cost - Assessing Performance

Stage cost function: $L(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$

- $L(\mathbf{x}, \mathbf{u}) = +\infty$ for infeasible states and inputs (constraints in the cost)
- $L(\mathbf{x}, \mathbf{u})$ may not be lower-bounded

Examples

LQR

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \frac{1}{2} \mathbf{u}^\top R \mathbf{u}$$

Stage cost - Assessing Performance

Stage cost function: $L(\mathbf{x}, \mathbf{u}) \in \mathbb{R}$

- $L(\mathbf{x}, \mathbf{u}) = +\infty$ for infeasible states and inputs (constraints in the cost)
- $L(\mathbf{x}, \mathbf{u})$ may not be lower-bounded

Examples

LQR

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \frac{1}{2} \mathbf{u}^\top R \mathbf{u}$$

MPC

$$\frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \frac{1}{2} \mathbf{u}^\top R \mathbf{u}$$

s.t. $\mathbf{x} \in \mathcal{X}, \mathbf{u} \in \mathcal{U}$

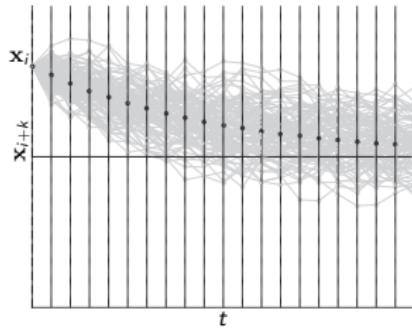
yields:

$$L(\mathbf{x}, \mathbf{u}) = \begin{cases} \frac{1}{2} \mathbf{x}^\top Q \mathbf{x} + \frac{1}{2} \mathbf{u}^\top R \mathbf{u} & \text{if } \mathbf{x} \in \mathcal{X}, \mathbf{u} \in \mathcal{U} \\ \infty & \text{otherwise} \end{cases}$$

MDP, return and discount - How to assess long-term performance?

The triplet $(\mathbb{P}[x_+ | x, u], \gamma, L)$ define a Markov Decision Process

- State and input x, u
- State transition (model) $\mathbb{P}[x_+ | x, u]$
- Stage cost L
- Discount γ



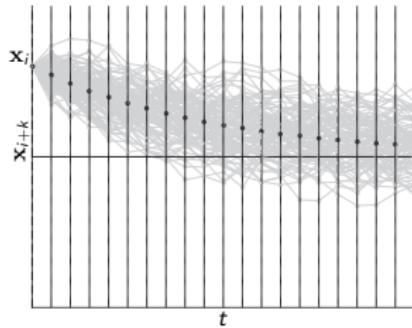
MDP, return and discount - How to assess long-term performance?

The triplet $(\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}], \gamma, L)$ define a Markov Decision Process

- State and input \mathbf{x}, \mathbf{u}
- State transition (model) $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$
- Stage cost L
- Discount γ

Cost-to-go (a.k.a. return)

$$J = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{i+k}, \mathbf{u}_{i+k}) \right]$$



where

- $\mathbf{x}_{i+k}, \mathbf{u}_{i+k}$ is the state trajectory and input sequence
- $\gamma \leq 1$ is the discount factor

MDP, return and discount - How to assess long-term performance?

The triplet $(\mathbb{P}[x_+ | x, u], \gamma, L)$ define a Markov Decision Process

- State and input x, u
- State transition (model) $\mathbb{P}[x_+ | x, u]$
- Stage cost L
- Discount γ

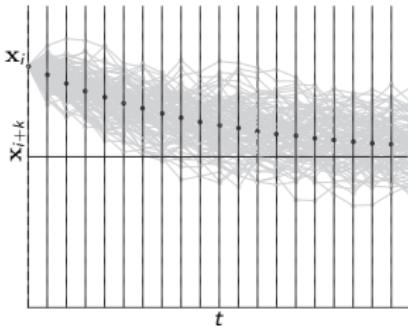
Cost-to-go (a.k.a. return)

$$J = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_{i+k}, u_{i+k}) \right]$$

Why a discount factor?

Policy (decisions) taken such that J is minimized

What is \mathbb{E} expecting over (distribution)?



where

- x_{i+k}, u_{i+k} is the state trajectory and input sequence
- $\gamma \leq 1$ is the discount factor

Policy - How to define decisions?

Deterministic

$$u = \pi(x)$$

Meaning: what input u do we select
provided that we are in state x

Policy - How to define decisions?

Deterministic

$$\mathbf{u} = \pi(\mathbf{x})$$

Meaning: what input \mathbf{u} do we select
provided that we are in state \mathbf{x}

Example LQR:

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$$

$$\mathbf{x}_+ = A\mathbf{x} + B\mathbf{u}$$

yields: $\mathbf{u} = -K\mathbf{x}$

hence: $\pi(\mathbf{x}) = -K\mathbf{x}$

Policy - How to define decisions?

Deterministic

$$\mathbf{u} = \pi(\mathbf{x})$$

Stochastic

$$\pi[\mathbf{u} | \mathbf{x}]$$

Meaning: what input \mathbf{u} do we select
provided that we are in state \mathbf{x}

Example LQR:

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$$

$$\mathbf{x}_+ = A\mathbf{x} + B\mathbf{u}$$

yields: $\mathbf{u} = -K\mathbf{x}$

hence: $\pi(\mathbf{x}) = -K\mathbf{x}$

Policy - How to define decisions?

Deterministic

$$\mathbf{u} = \pi(\mathbf{x})$$

Meaning: what input \mathbf{u} do we select
provided that we are in state \mathbf{x}

Example LQR:

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$$

$$\mathbf{x}_+ = A\mathbf{x} + B\mathbf{u}$$

yields: $\mathbf{u} = -K\mathbf{x}$

hence: $\pi(\mathbf{x}) = -K\mathbf{x}$

Stochastic

$$\pi[\mathbf{u} | \mathbf{x}]$$

Meaning: probability of selecting input \mathbf{u}
provided that we are in state \mathbf{x}

Policy - How to define decisions?

Deterministic

$$\mathbf{u} = \pi(\mathbf{x})$$

Meaning: what input \mathbf{u} do we select
provided that we are in state \mathbf{x}

Example LQR:

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$$

$$\mathbf{x}_+ = A\mathbf{x} + B\mathbf{u}$$

yields: $\mathbf{u} = -K\mathbf{x}$

hence: $\pi(\mathbf{x}) = -K\mathbf{x}$

Stochastic

$$\pi[\mathbf{u} | \mathbf{x}]$$

Meaning: probability of selecting input \mathbf{u}
provided that we are in state \mathbf{x}

Why?

Policy - How to define decisions?

Deterministic

$$u = \pi(x)$$

Meaning: what input u do we select
provided that we are in state x

Example LQR:

$$L(x, u) = \frac{1}{2}x^\top Qx + \frac{1}{2}u^\top Ru$$

$$x_+ = Ax + Bu$$

yields: $u = -Kx$

hence: $\pi(x) = -Kx$

Stochastic

$$\pi[u | x]$$

Meaning: probability of selecting input u
provided that we are in state x

Why?

- unsure of the policy, "try" different decisions

Policy - How to define decisions?

Deterministic

$$\mathbf{u} = \pi(\mathbf{x})$$

Meaning: what input \mathbf{u} do we select
provided that we are in state \mathbf{x}

Example LQR:

$$L(\mathbf{x}, \mathbf{u}) = \frac{1}{2}\mathbf{x}^\top Q\mathbf{x} + \frac{1}{2}\mathbf{u}^\top R\mathbf{u}$$

$$\mathbf{x}_+ = A\mathbf{x} + B\mathbf{u}$$

yields: $\mathbf{u} = -K\mathbf{x}$

hence: $\pi(\mathbf{x}) = -K\mathbf{x}$

Stochastic

$$\pi[\mathbf{u} | \mathbf{x}]$$

Meaning: probability of selecting input \mathbf{u}
provided that we are in state \mathbf{x}

Why?

- unsure of the policy, “try” different decisions
- taking random actions can be better, e.g. Rock-Paper-Scissor. Not an MDP though, but a POMDP

Policy - How to define decisions?

Deterministic

$$u = \pi(x)$$

Meaning: what input u do we select
provided that we are in state x

Example LQR:

$$L(x, u) = \frac{1}{2}x^\top Qx + \frac{1}{2}u^\top Ru$$

$$x_+ = Ax + Bu$$

yields: $u = -Kx$

hence: $\pi(x) = -Kx$

Stochastic

$$\pi(u | x)$$

Meaning: probability of selecting input u
provided that we are in state x

Why?

- unsure of the policy, “try” different decisions
- taking random actions can be better, e.g. Rock-Paper-Scissor. Not an MDP though, but a POMDP
- more on that later...

Policy - How to define decisions?

Deterministic

$$u = \pi(x)$$

Meaning: what input u do we select
provided that we are in state x

Example LQR:

$$L(x, u) = \frac{1}{2}x^\top Qx + \frac{1}{2}u^\top Ru$$

$$x_+ = Ax + Bu$$

yields: $u = -Kx$

hence: $\pi(x) = -Kx$

Stochastic

$$\pi(u | x)$$

Meaning: probability of selecting input u
provided that we are in state x

Why?

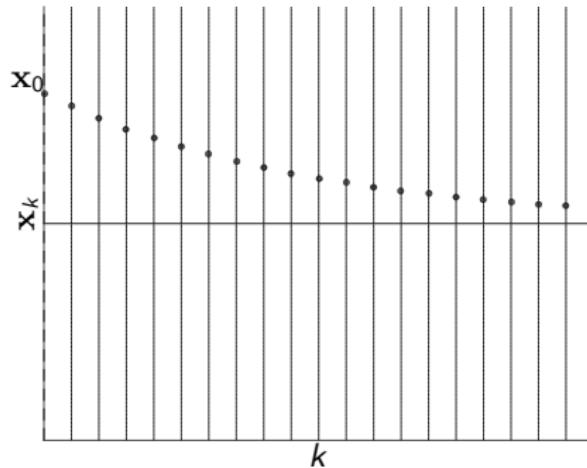
- unsure of the policy, "try" different decisions
- taking random actions can be better, e.g. Rock-Paper-Scissor. Not an MDP though, but a POMDP
- more on that later...

Remarks

- MDPs always have a deterministic optimal policy
- Optimal policy may not be unique
- If several actions are equally good, then a random choice between them is still optimum, hence some MDPs admit stochastic optimal policies

Unpacking the Markov Chain

Unpacking the Markov Chain

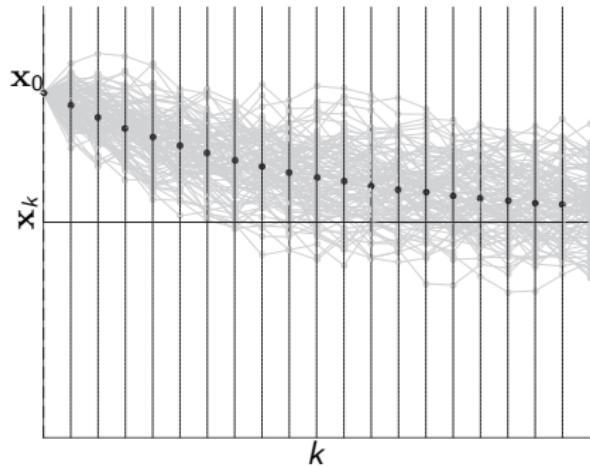


Deterministic state transition

$$\mathbf{x}_+ = \mathbf{f}(\mathbf{x}, \mathbf{u})$$

$$\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}] = \delta(\mathbf{x}_+ - \mathbf{f}(\mathbf{x}, \mathbf{u}))$$

Unpacking the Markov Chain



Stochastic state transition

$\mathbf{x}, \mathbf{u} \rightarrow \mathbf{x}_+ \quad \text{defined by} \quad \mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

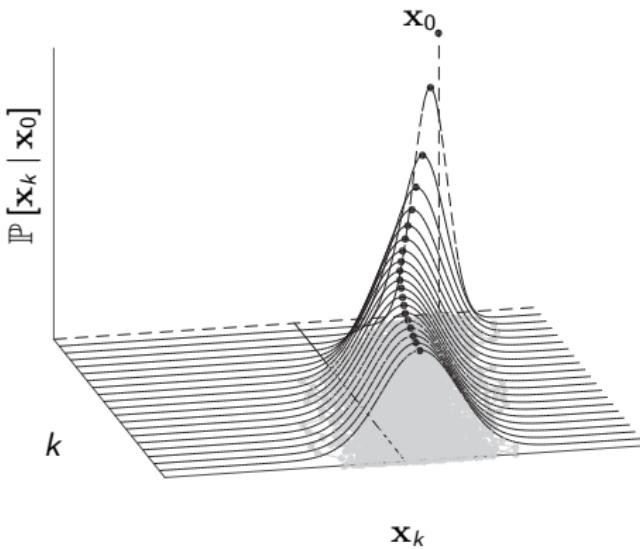
Policy

$\mathbf{u} = \pi(\mathbf{x}), \quad \text{or} \quad \pi[\mathbf{u} | \mathbf{x}]$

Transition $\mathbf{x}_0 \rightarrow \mathbf{x}_k$ is stochastic:

$\mathbb{P}[\mathbf{x}_k | \mathbf{x}_0]$

Unpacking the Markov Chain



Stochastic state transition

$$x, u \rightarrow x_+ \quad \text{defined by} \quad \mathbb{P}[x_+ | x, u]$$

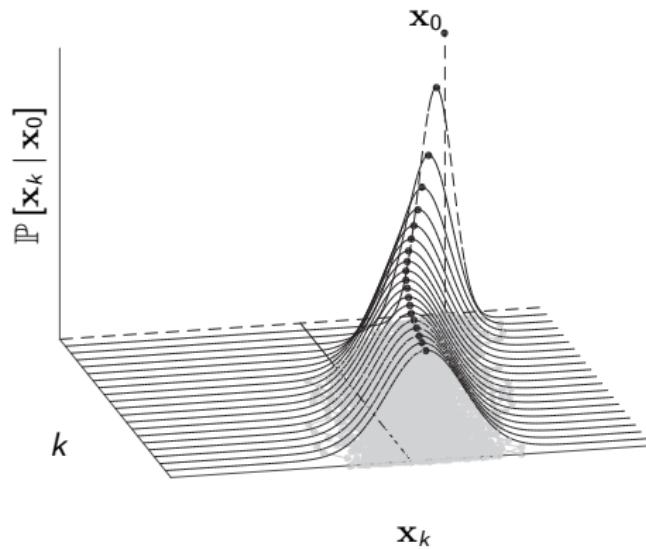
Policy

$$u = \pi(x), \quad \text{or} \quad \pi[u | x]$$

Transition $x_0 \rightarrow x_k$ is stochastic:

$$\mathbb{P}[x_k | x_0]$$

Unpacking the Markov Chain



Stochastic state transition

$$x, u \rightarrow x_+ \quad \text{defined by} \quad P[x_+ | x, u]$$

Policy

$$u = \pi(x), \quad \text{or} \quad \pi[u | x]$$

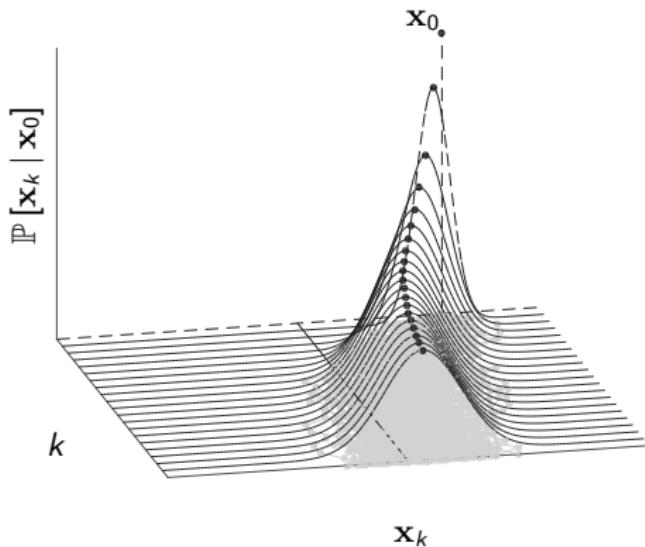
Transition $x_0 \rightarrow x_k$ is stochastic:

$$P[x_k | x_0]$$

1-step distribution:

$$P[x_1 | x_0, u_0 = \pi(x_0)] \equiv P[x_1 | x_0, \pi(x_0)]$$

Unpacking the Markov Chain



Stochastic state transition

$x, u \rightarrow x_+$ defined by $\mathbb{P}[x_+ | x, u]$

Policy

$u = \pi(x)$, or $\pi[u | x]$

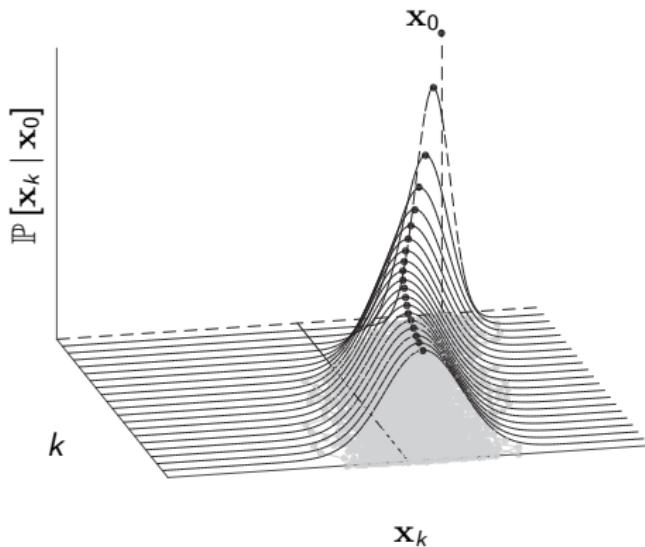
Transition $x_0 \rightarrow x_k$ is stochastic:

$\mathbb{P}[x_k | x_0]$

2-step distribution:

$$\mathbb{P}[x_2 | x_0] = \int \mathbb{P}[x_2 | x_1, \pi(x_1)] \mathbb{P}[x_1 | x_0, u_0] dx_1$$

Unpacking the Markov Chain



Stochastic state transition

$$x, u \rightarrow x_+ \quad \text{defined by} \quad P[x_+ | x, u]$$

Policy

$$u = \pi(x), \quad \text{or} \quad \pi[u | x]$$

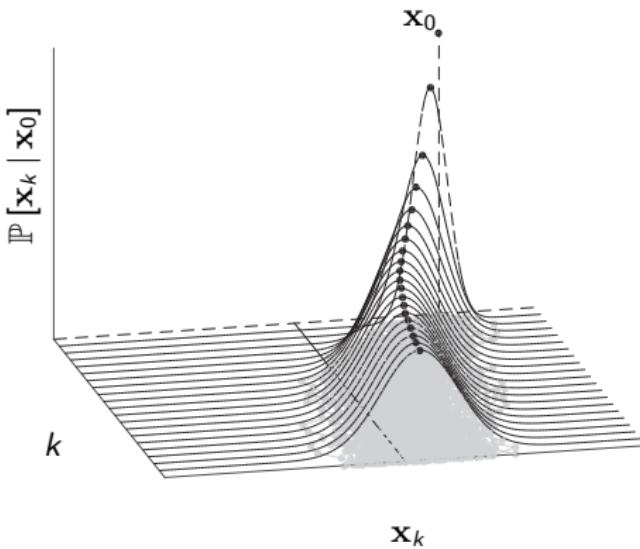
Transition $x_0 \rightarrow x_k$ is stochastic:

$$P[x_k | x_0]$$

3-step distribution:

$$P[x_3 | x_0] = \int P[x_3 | x_2, \pi(x_2)] P[x_2 | x_1, \pi(x_1)] P[x_1 | x_0, \pi(x_0)] dx_1 dx_2$$

Unpacking the Markov Chain



Stochastic state transition

$\mathbf{x}, \mathbf{u} \rightarrow \mathbf{x}_+$ defined by $\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \mathbf{u}]$

Policy

$\mathbf{u} = \pi(\mathbf{x}), \text{ or } \pi[\mathbf{u} | \mathbf{x}]$

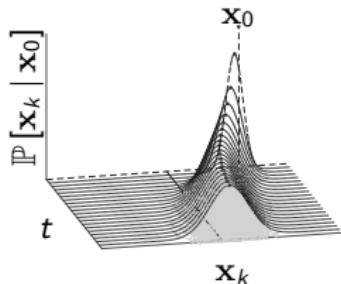
Transition $\mathbf{x}_0 \rightarrow \mathbf{x}_k$ is stochastic:

$\mathbb{P}[\mathbf{x}_k | \mathbf{x}_0]$

n-step distribution:

$$\mathbb{P}[\mathbf{x}_k | \mathbf{x}_0] = \int \prod_{i=0}^{k-1} \mathbb{P}[\mathbf{x}_{i+1} | \mathbf{x}_i, \pi(\mathbf{x}_i)] d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

Unpacking the Markov Chain



Stationary distribution: how often are the different states visited under a given policy, after the Markov chain reaches its stationary dynamics?

$$P[x_\infty | x_0]$$

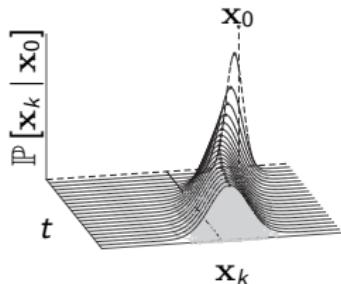
Markov chain distribution: τ_π describes the distribution of trajectories, i.e. the probability distribution that specific trajectories $x_{0,\dots,\infty}$ are observed under policy π

- Distribution $x_0 \rightarrow x_k$

$$P[x_k | x_0] = \int \prod_{i=0}^{k-1} P[x_{i+1} | x_i, \pi(x_i)] \cdot dx_1 \dots dx_{k-1}$$

- Expectation over distribution τ_π , e.g. performance of policy π :

Unpacking the Markov Chain



Stationary distribution: how often are the different states visited under a given policy, after the Markov chain reaches its stationary dynamics?

$$\mathbb{P}[x_\infty | x_0]$$

Markov chain distribution: τ_π describes the distribution of trajectories, i.e. the probability distribution that specific trajectories x_0, \dots, ∞ are observed under policy π

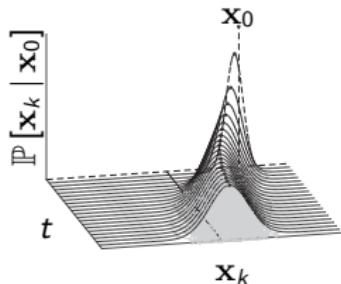
- Distribution $x_0 \rightarrow x_k$

$$\mathbb{P}[x_k | x_0] = \int \prod_{i=0}^{k-1} \mathbb{P}[x_{i+1} | x_i, \pi(x_i)] \cdot dx_1 \dots dx_{k-1}$$

- Expectation over distribution τ_π , e.g. performance of policy π :

$$J(\pi) = \mathbb{E}_{\tau_\pi}[L(x, \pi(x))] = \int \sum_{k=0}^{\infty} \gamma^k L(x_k, \pi(x_k)) \mathbb{P}[x_k | x_0] \mathbb{P}[x_0] dx_0 dx_k$$

Unpacking the Markov Chain



Stationary distribution: how often are the different states visited under a given policy, after the Markov chain reaches its stationary dynamics?

$$\mathbb{P}[x_\infty | x_0]$$

Markov chain distribution: τ_π describes the distribution of trajectories, i.e. the probability distribution that specific trajectories x_0, \dots, ∞ are observed under policy π

- Distribution $x_0 \rightarrow x_k$

$$\mathbb{P}[x_k | x_0] = \int \prod_{i=0}^{k-1} \mathbb{P}[x_{i+1} | x_i, \pi(x_i)] \cdot dx_1 \dots dx_{k-1}$$

- Expectation over distribution τ_π , e.g. performance of policy π :

$$J(\pi) = \sum_{k=0}^{\infty} \int \gamma^k L(x_k, \pi(x_k)) \prod_{i=1}^k \mathbb{P}[x_i | x_{i-1}, \pi(x_{i-1})] \mathbb{P}[x_0] \cdot dx_0 \dots dx_{k-1}$$

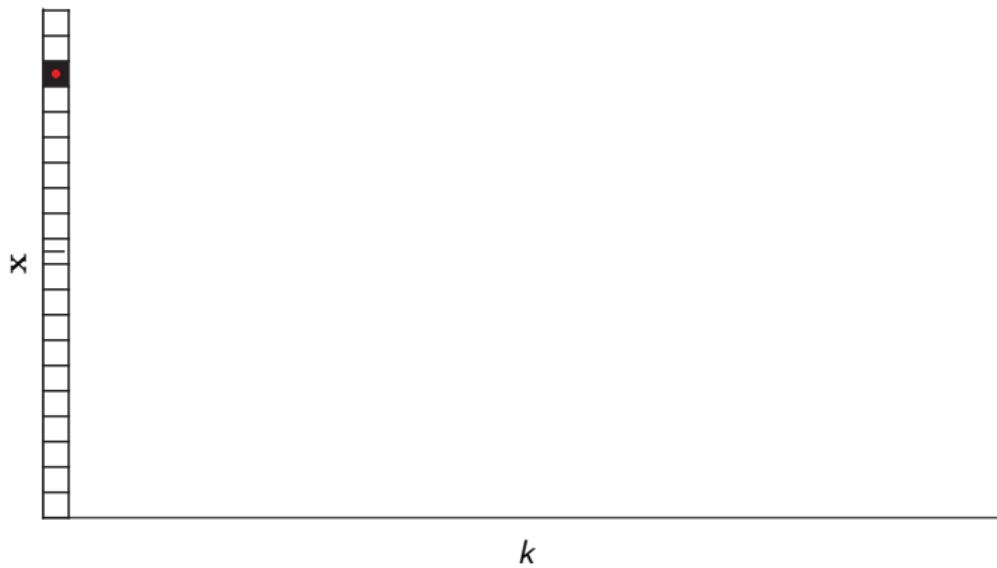
Outline

- 1 Markov Decision Process
- 2 Value functions & Bellman equations
- 3 Function evaluation & Dynamic Programming
- 4 Quick view of Bellman contraction

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

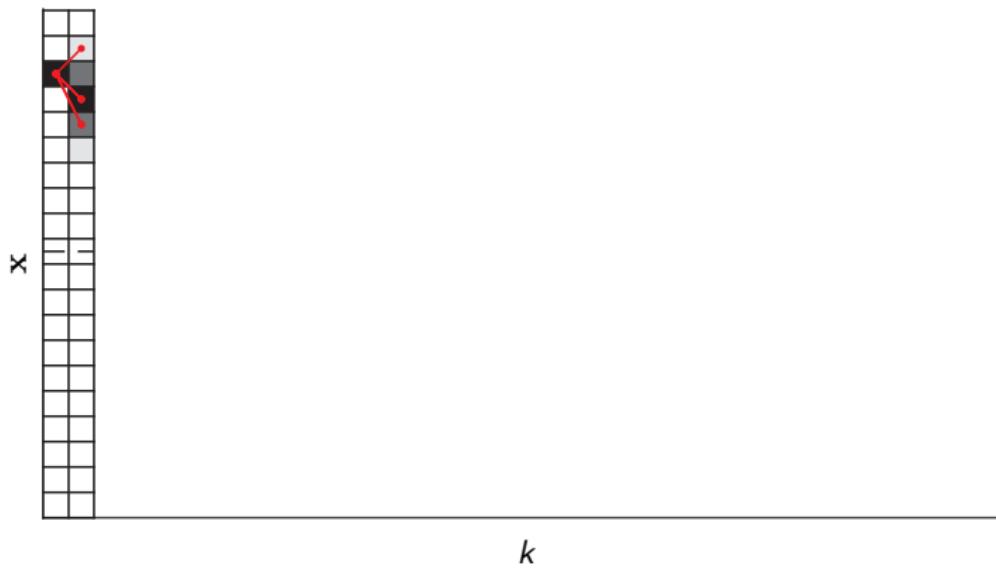
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

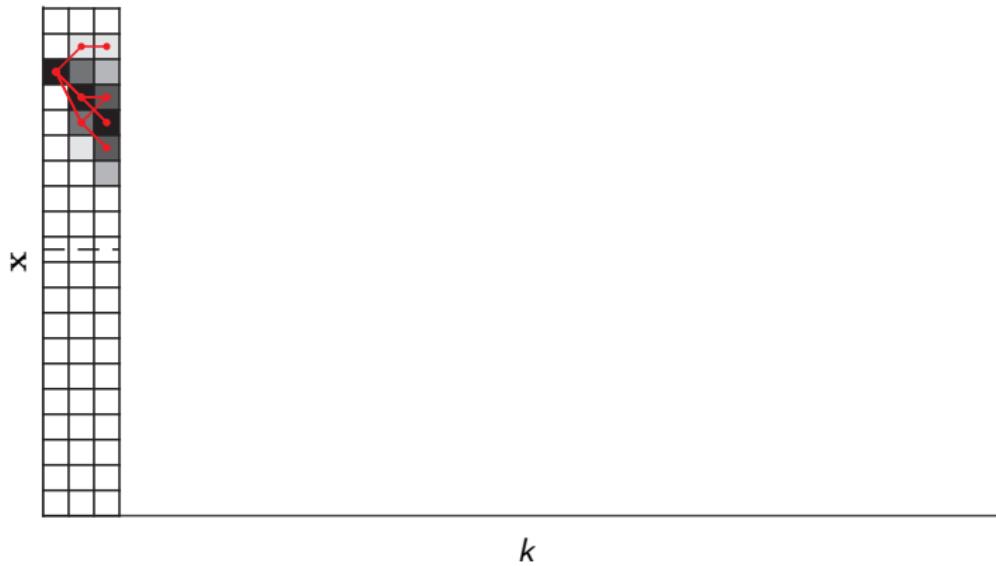


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

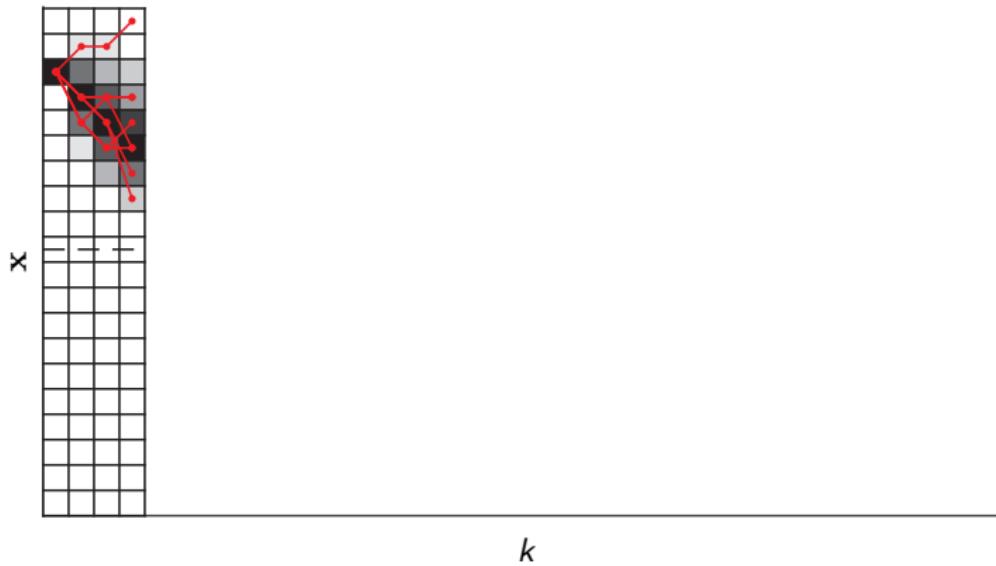
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

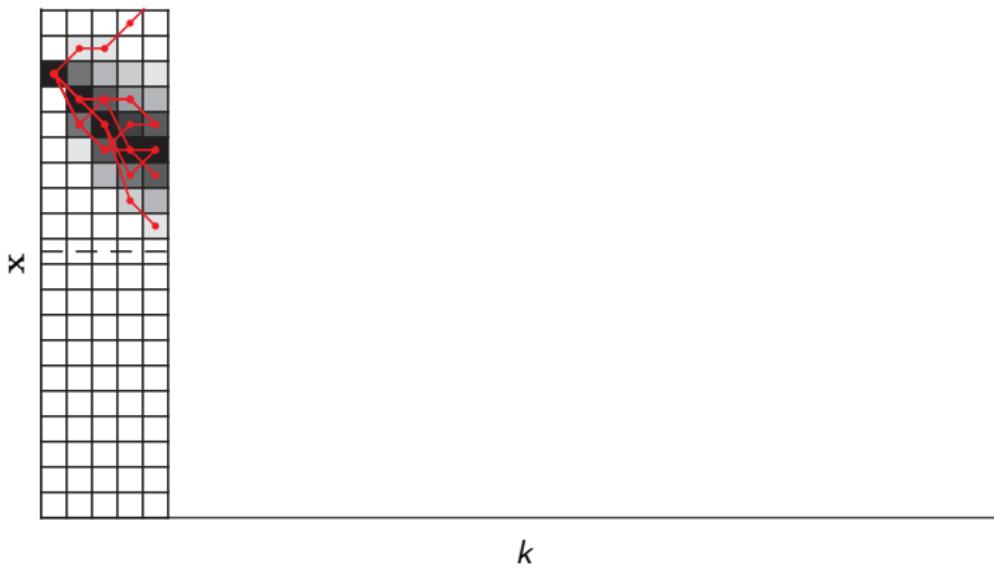


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

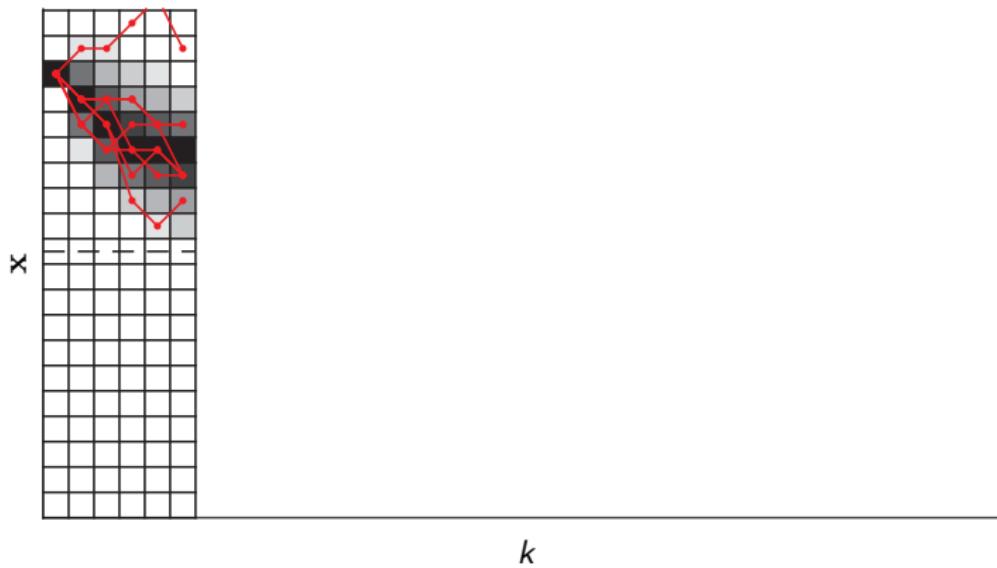


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

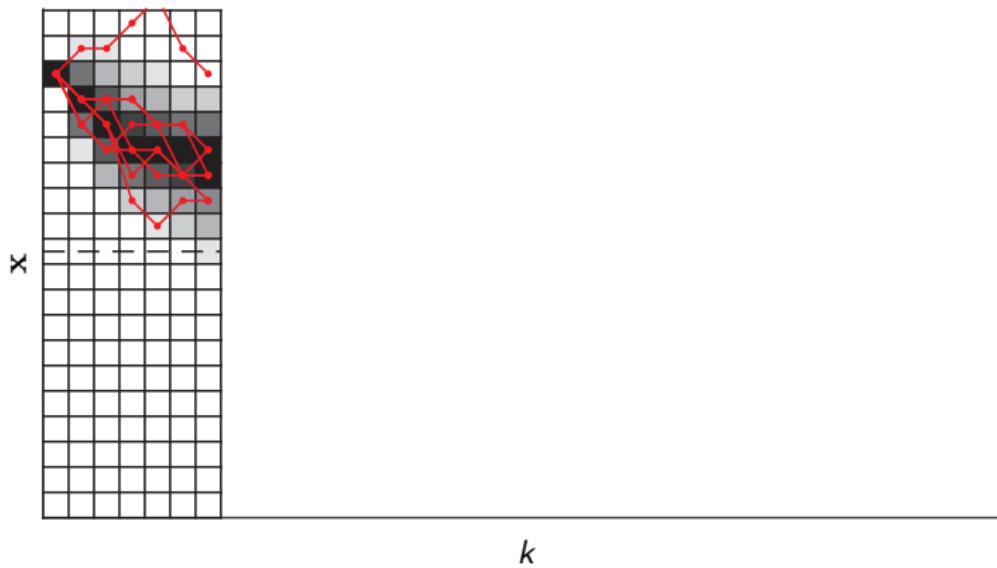


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

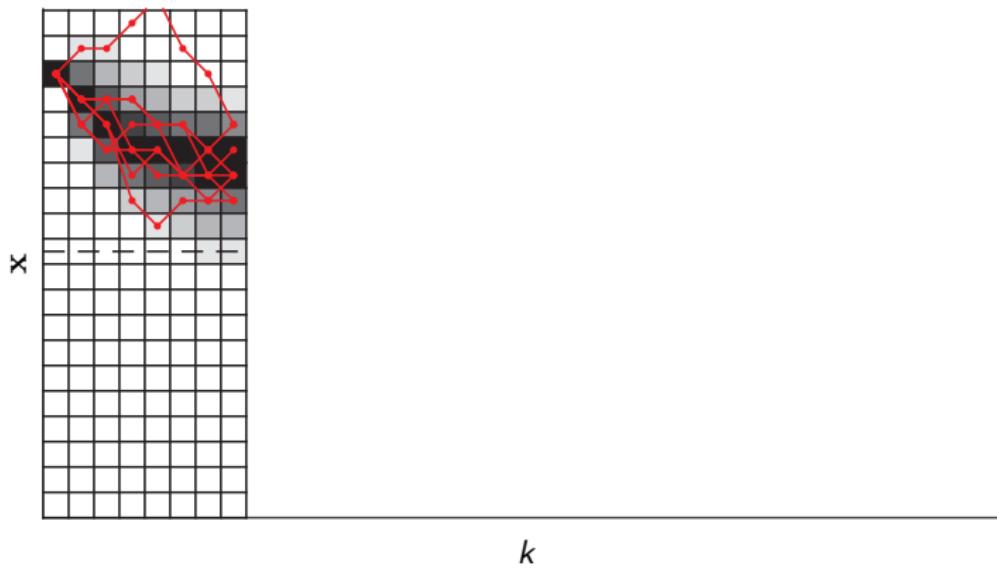


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

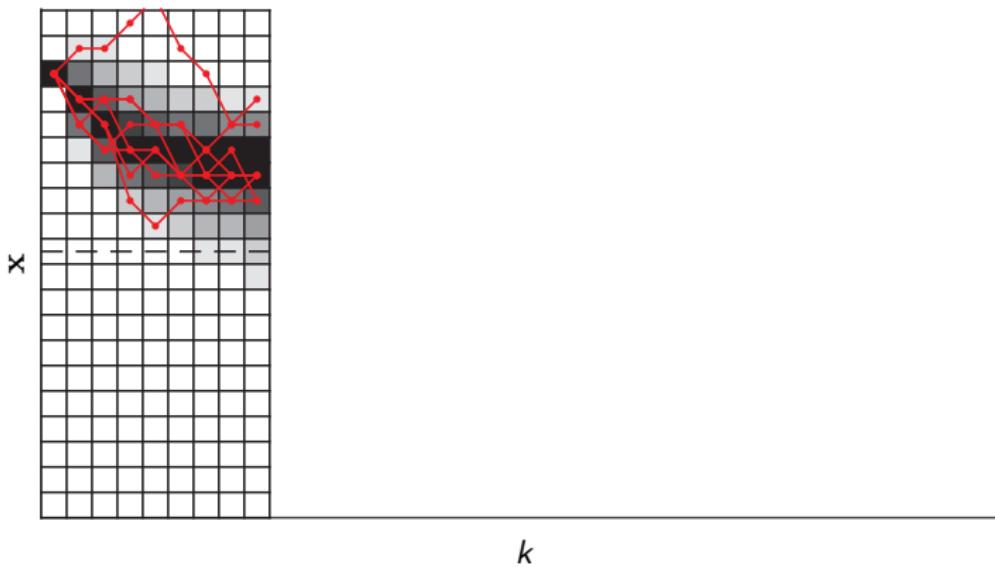


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

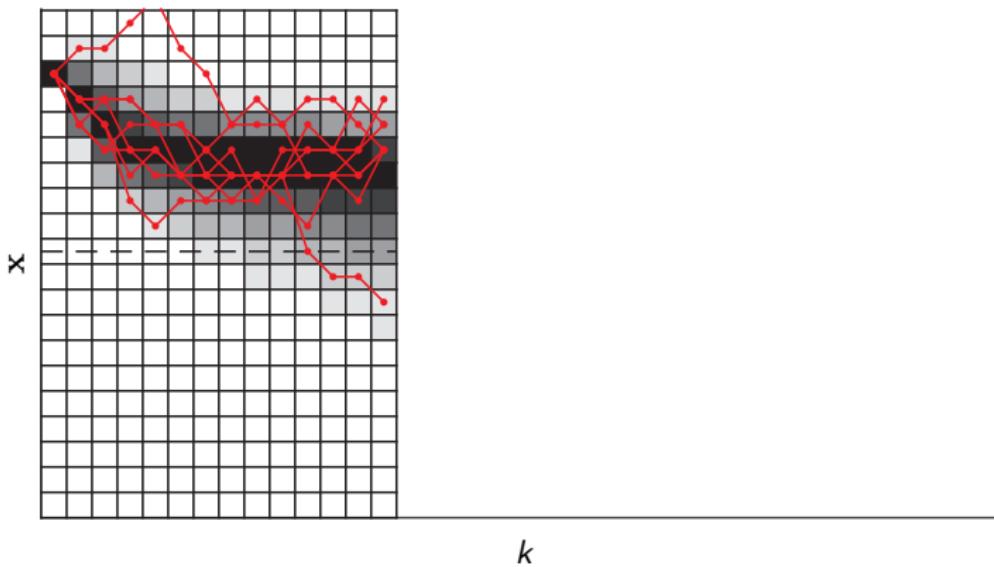


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

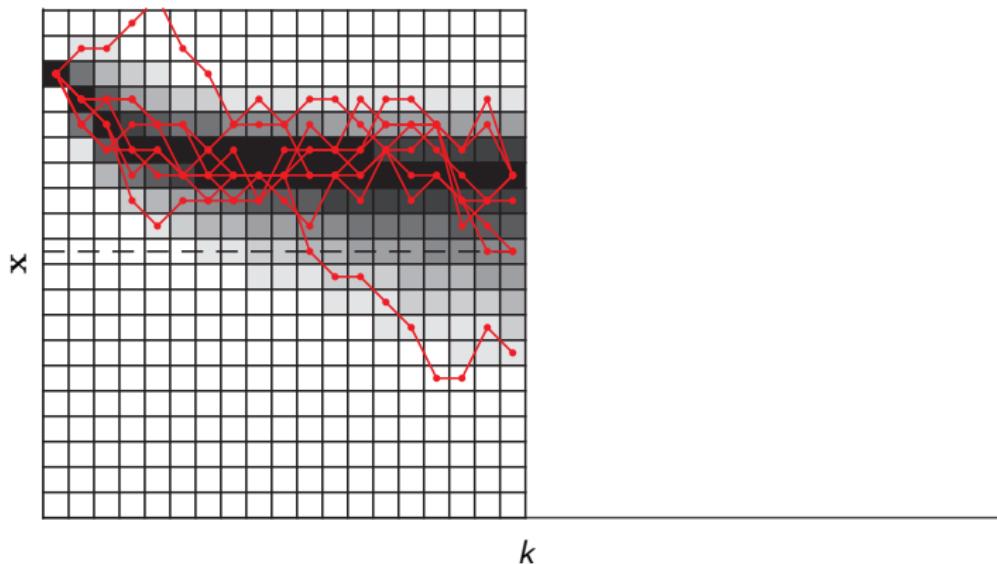


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

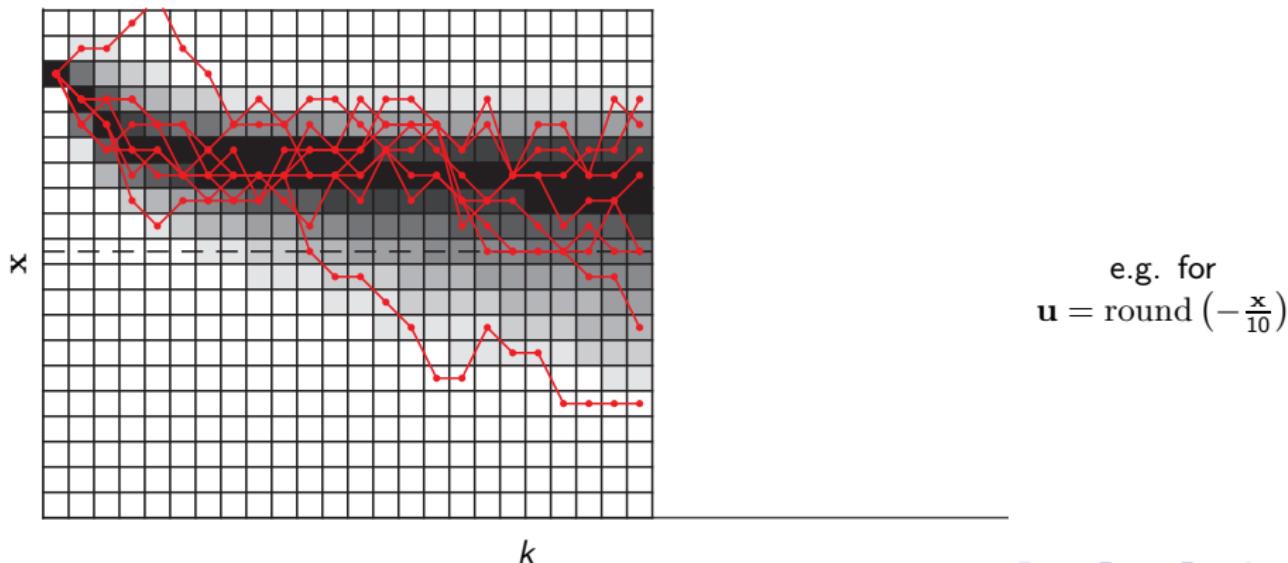


e.g. for
 $u = \text{round}\left(-\frac{x}{10}\right)$

Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

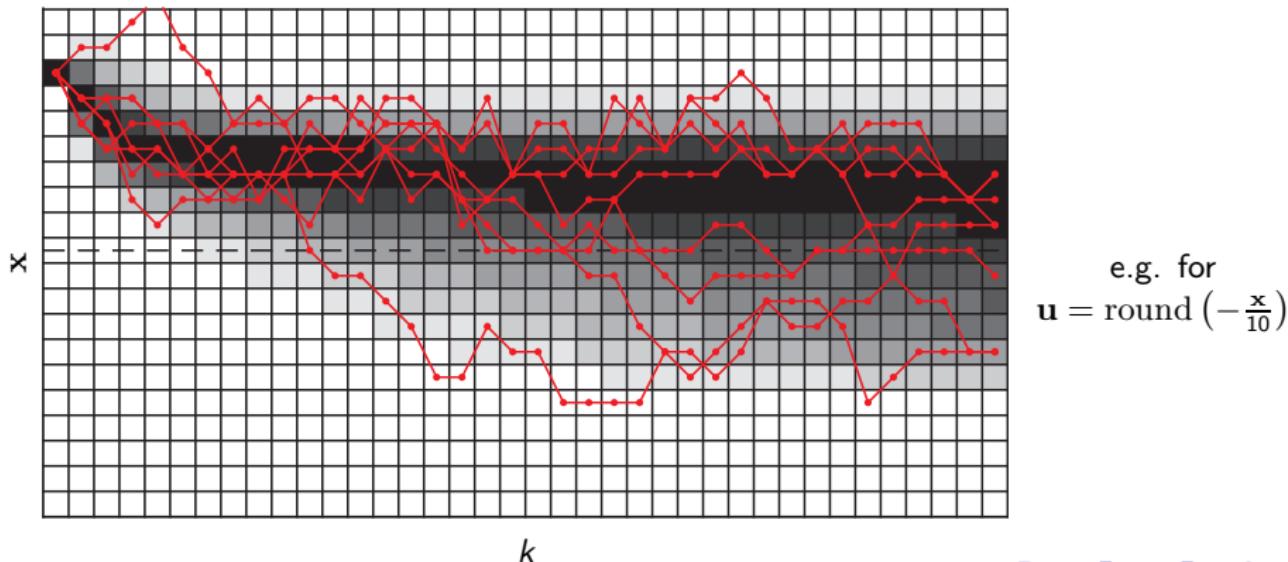
Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$



Let's use the following example (for illustration purposes)

- State $x \in \mathbb{N}$ and input $u \in \mathbb{N}$ locked on a grid
- Input u can only move x on the grid
- Stage cost: $L(x, u) = \frac{1}{2}x^2 + \frac{1}{2}u^2$
- Discount $\gamma = 0.9$

Dynamics: $x_+ = x + u + \text{round}(e)$ where $e \sim \mathcal{N}(0, 1)$

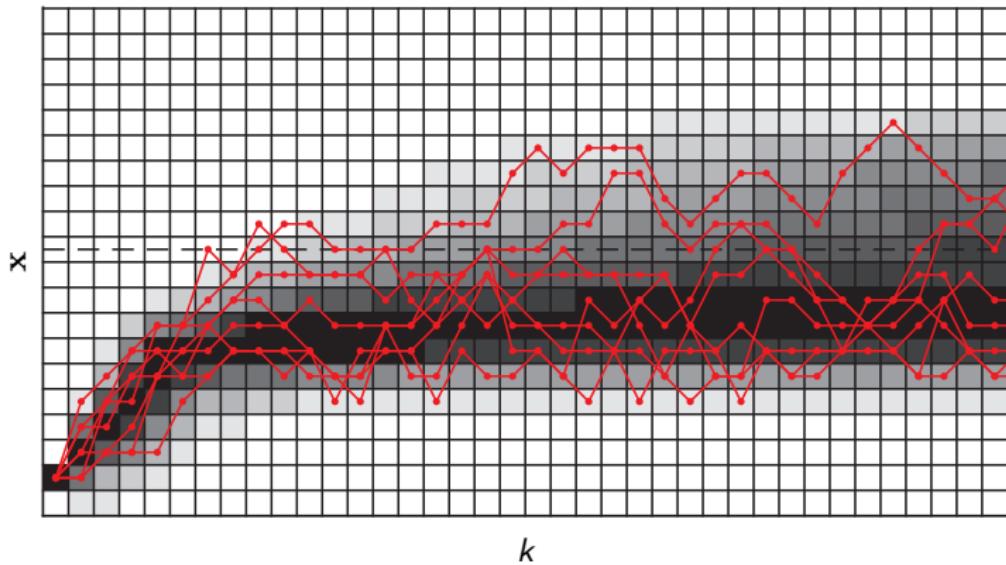


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 155.18$$

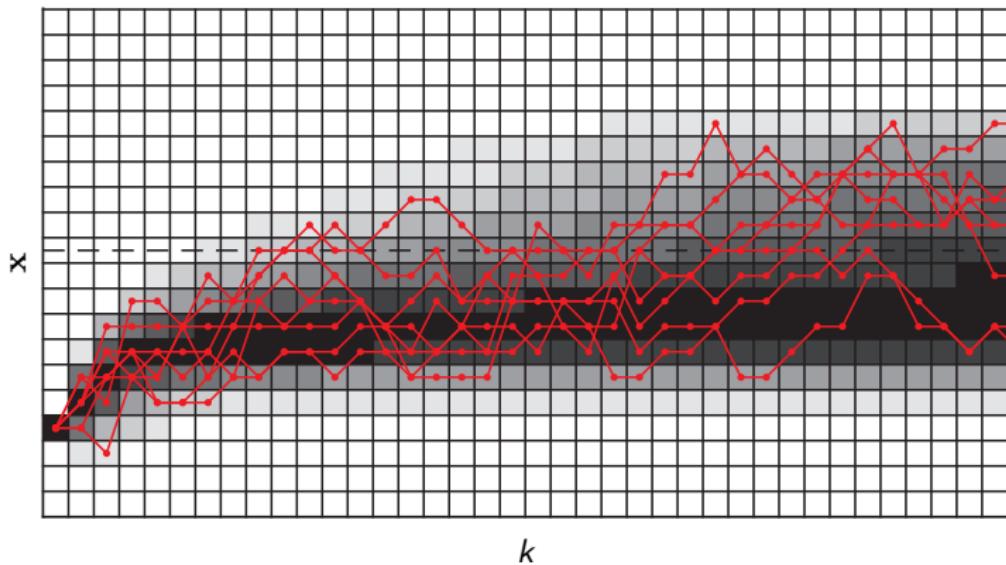


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 99.94$$

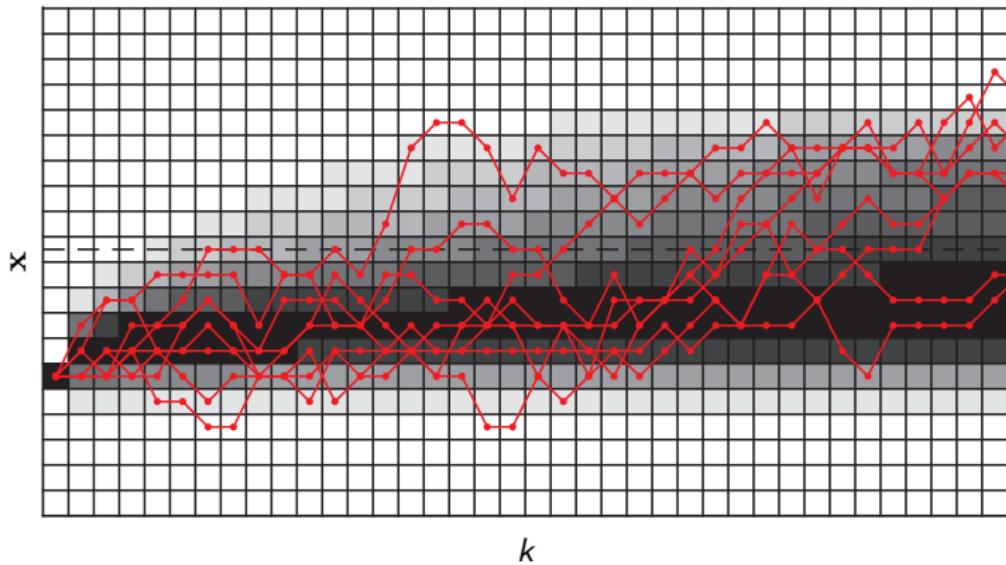


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 65.41$$

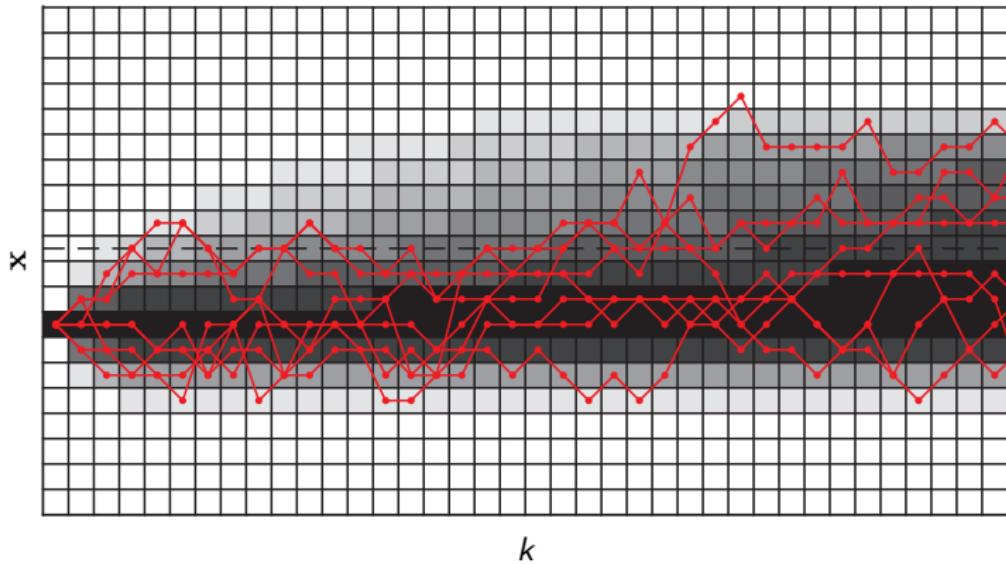


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 47.5$$

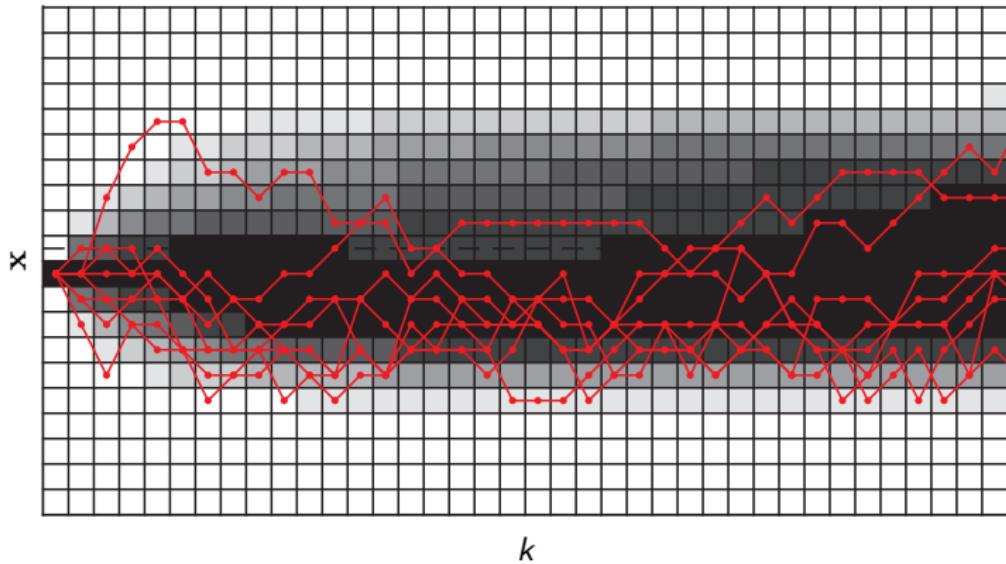


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 28.23$$

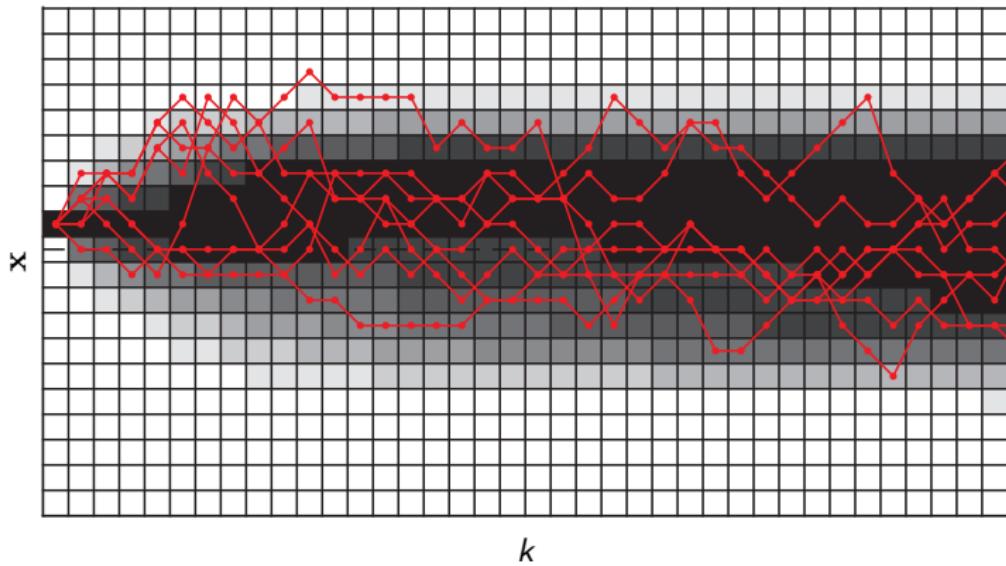


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 28.23$$

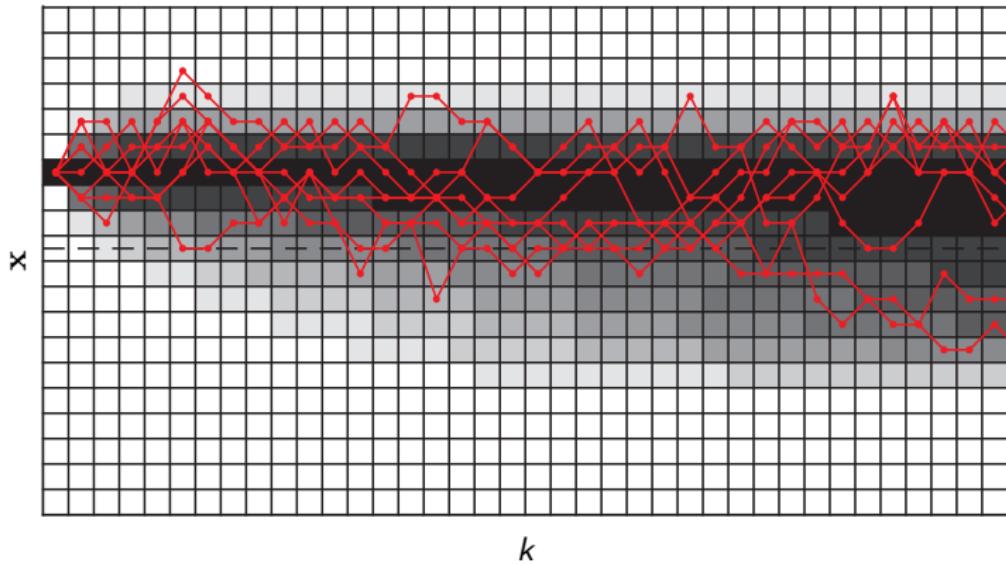


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 47.5$$

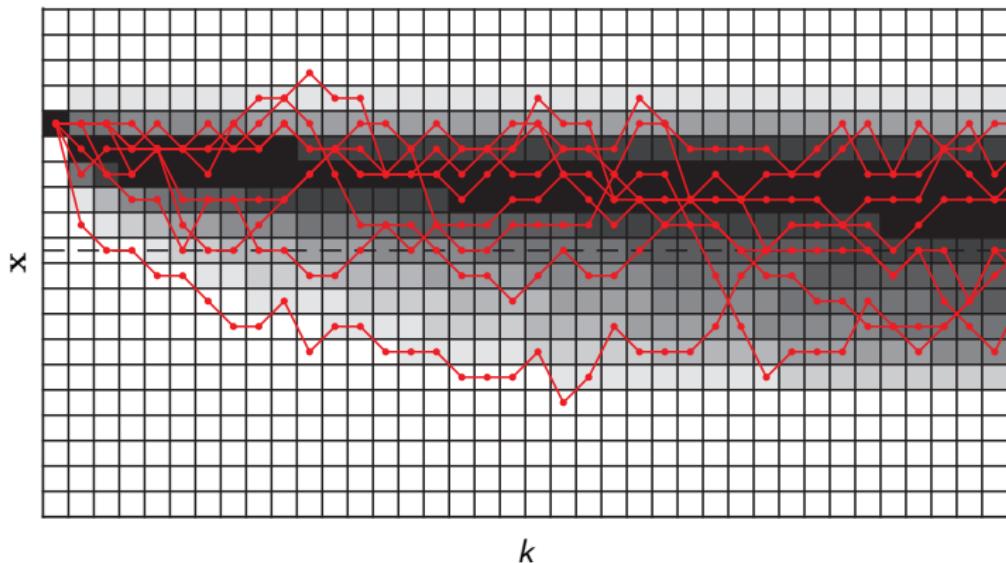


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 65.41$$

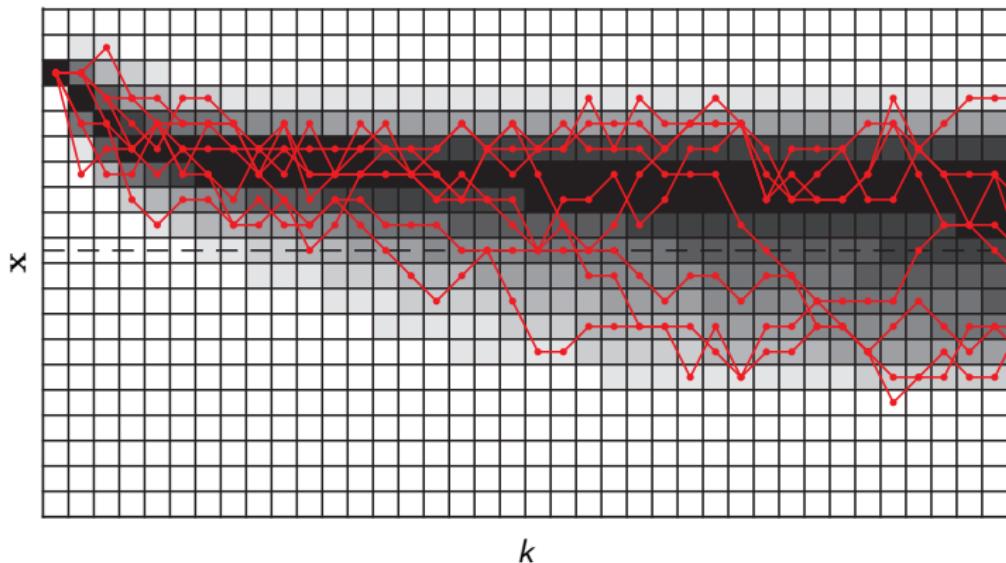


Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

$$V_{\pi}(\mathbf{x}_0) = 99.94$$



Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(x) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(x_k, u_k) \mid u_k = \pi(x_k), x_0 = x \right]$$

Meaning: for a given state x

- We follow policy π for all k (can be stochastic)
- State trajectory $x_{1,\dots,\infty}$ is stochastic as it follows $\mathbb{P}[x_+ | x, u]$
- Initial state $x_0 = x$ is fixed
- Value function at x is the resulting cumulated (discounted) cost
- Expected value $\mathbb{E}[\cdot]$ taken over trajectories, using trajectories distributions...

Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Formally:

$$V_{\pi}(\mathbf{x}) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0 = \mathbf{x}] d\mathbf{x}_k$$

where we use the distributions $\mathbf{x}_0 \rightarrow \mathbf{x}_k$, recall

$$\mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0] = \int \prod_{i=1}^k \mathbb{P}[\mathbf{x}_i \mid \mathbf{x}_{i-1}, \pi(\mathbf{x}_{i-1})] \cdot d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Formally:

$$V_{\pi}(\mathbf{x}) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0 = \mathbf{x}] d\mathbf{x}_k$$

where we use the distributions $\mathbf{x}_0 \rightarrow \mathbf{x}_k$, recall

$$\mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0] = \int \prod_{i=1}^k \mathbb{P}[\mathbf{x}_i \mid \mathbf{x}_{i-1}, \pi(\mathbf{x}_{i-1})] \cdot d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

i.e.

$$V_{\pi}(\mathbf{x}_0) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \prod_{i=1}^k \mathbb{P}[\mathbf{x}_i \mid \mathbf{x}_{i-1}, \pi(\mathbf{x}_{i-1})] \cdot d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

Note that this construction has all the properties of an expected value (e.g. linearity)

Value function - How good is it to be in a given state?

Value function for policy π :

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{u}_k = \pi(\mathbf{x}_k), \mathbf{x}_0 = \mathbf{x} \right]$$

Formally:

$$V_{\pi}(\mathbf{x}) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0 = \mathbf{x}] d\mathbf{x}_k$$

Relation $J(\pi)$ to $V_{\pi}(\mathbf{x})$?

$$\begin{aligned} J(\pi) &= \mathbb{E}[V_{\pi}(\mathbf{x})] \\ &= \int V_{\pi}(\mathbf{x}_0) \mathbb{P}[\mathbf{x}_0] d\mathbf{x}_0 \end{aligned}$$

where we use the distributions $\mathbf{x}_0 \rightarrow \mathbf{x}_k$, recall

$$\mathbb{P}[\mathbf{x}_k \mid \mathbf{x}_0] = \int \prod_{i=1}^k \mathbb{P}[\mathbf{x}_i \mid \mathbf{x}_{i-1}, \pi(\mathbf{x}_{i-1})] \cdot d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

i.e.

$$V_{\pi}(\mathbf{x}_0) = \sum_{k=0}^{\infty} \int \gamma^k L(\mathbf{x}_k, \pi(\mathbf{x}_k)) \prod_{i=1}^k \mathbb{P}[\mathbf{x}_i \mid \mathbf{x}_{i-1}, \pi(\mathbf{x}_{i-1})] \cdot d\mathbf{x}_1 \dots d\mathbf{x}_{k-1}$$

Note that this construction has all the properties of an expected value (e.g. linearity)

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Value function backup

$$V_{\pi}(\mathbf{x}_0) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right]$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Value function backup

$$V_{\pi}(\mathbf{x}_0) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] = L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \mathbb{E} \left[\sum_{k=1}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right]$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Value function backup

$$\begin{aligned} V_{\pi}(\mathbf{x}_0) &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] = L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \mathbb{E} \left[\sum_{k=1}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] \\ &= L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \gamma \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \mid \mathbf{x}_0, \pi \right] \end{aligned}$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Value function backup

$$\begin{aligned} V_\pi(\mathbf{x}_0) &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] = L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \mathbb{E} \left[\sum_{k=1}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] \\ &= L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \gamma \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \mid \mathbf{x}_0, \pi \right] \\ &= L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \gamma \mathbb{E} [V_\pi(\mathbf{x}_1) \mid \mathbf{x}_0, \pi] \end{aligned}$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_\pi(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

Value function backup

$$\begin{aligned} V_\pi(\mathbf{x}_0) &= \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] = L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \mathbb{E} \left[\sum_{k=1}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \mathbf{x}_0, \pi \right] \\ &= L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \gamma \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_{k+1}, \mathbf{u}_{k+1}) \mid \mathbf{x}_0, \pi \right] \\ &= L(\mathbf{x}_0, \pi(\mathbf{x}_0)) + \gamma \mathbb{E} [V_\pi(\mathbf{x}_1) \mid \mathbf{x}_0, \pi] \end{aligned}$$

where

$$\mathbb{E} [V_\pi(\mathbf{x}_1) \mid \mathbf{x}_0] = \int V_\pi(\mathbf{x}_1) \overbrace{\mathbb{P}[\mathbf{x}_1 \mid \mathbf{x}_0, \pi(\mathbf{x}_0)]}^{Model} d\mathbf{x}_1$$

Bellman equation - Computing the value function

Value function for policy π (either $\mathbf{u} = \pi(\mathbf{x})$ or $\pi[\mathbf{u} | \mathbf{x}]$):

$$V_{\pi}(\mathbf{x}) = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k L(\mathbf{x}_k, \mathbf{u}_k) \mid \pi, \mathbf{x}_0 = \mathbf{x} \right]$$

We get:

Bellman equation

$$V_{\pi}(\mathbf{x}) = L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V_{\pi}(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

sometimes called optimality backup

Action-value function

Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

Action-value function

Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

Interpretation:

- From state x , use given input u , yields (stochastic) transition $x, u \rightarrow x_+$
- Then from x_+ follow policy π
- Action-value function $Q_{\pi}(x, u)$ is the expected return of the resulting trajectory

Action-value function

Bellman equation for the value function

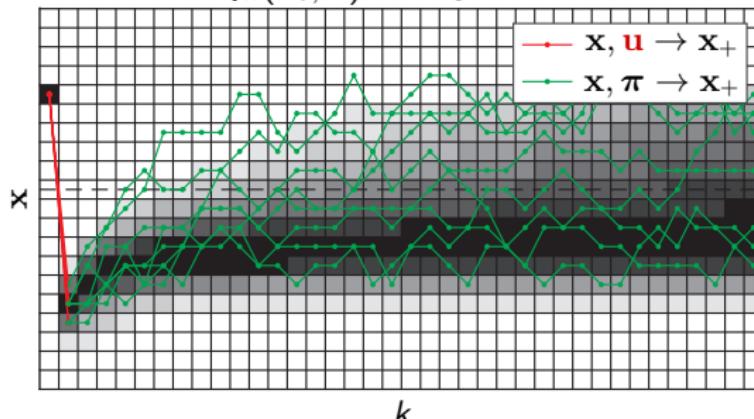
$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

$$Q_{\pi}(x_0, u) = 147.94$$



Action-value function

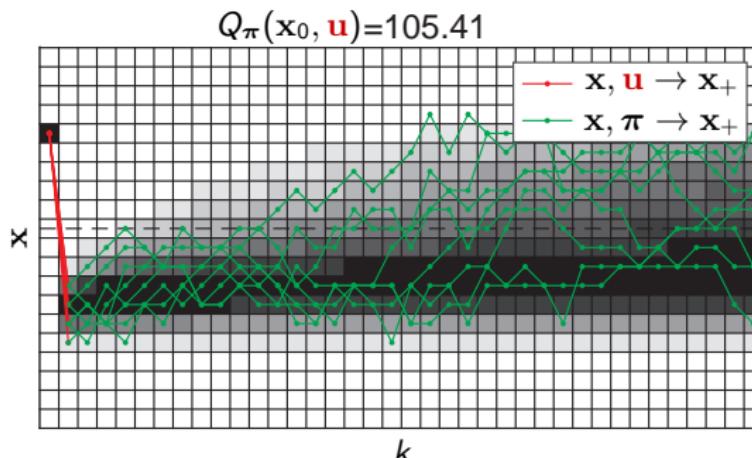
Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$



Action-value function

Bellman equation for the value function

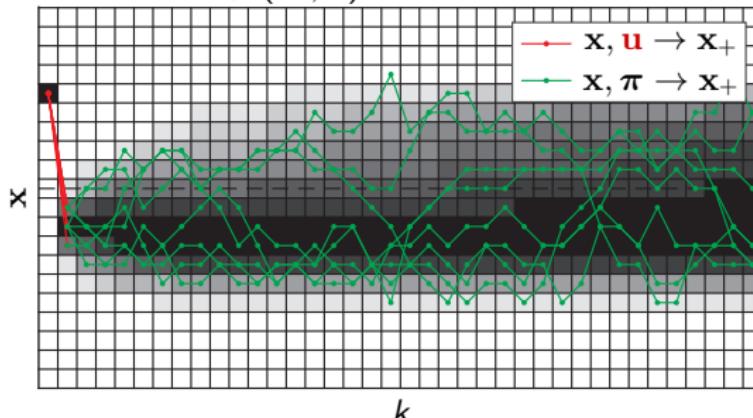
$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

$$Q_{\pi}(x_0, u) = 71.05$$



Action-value function

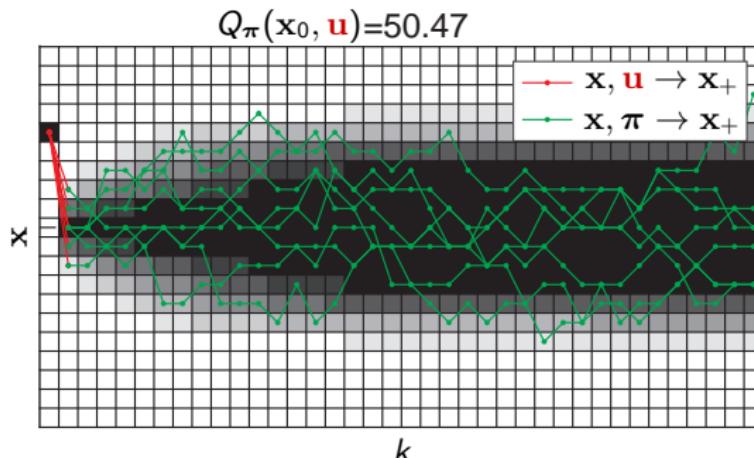
Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$



Action-value function

Bellman equation for the value function

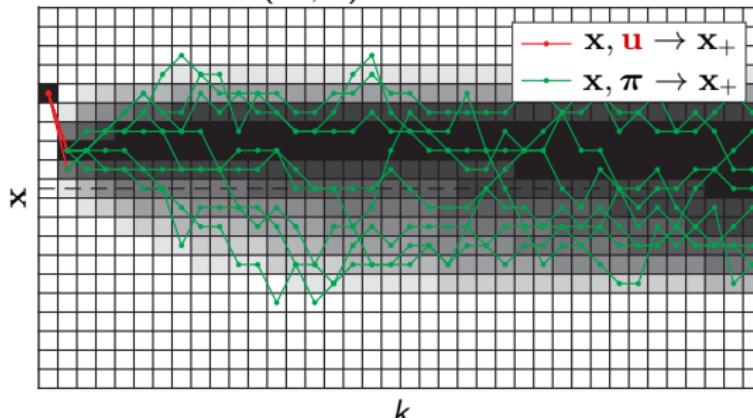
$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

$$Q_{\pi}(x_0, u) = 51.05$$



Action-value function

Bellman equation for the value function

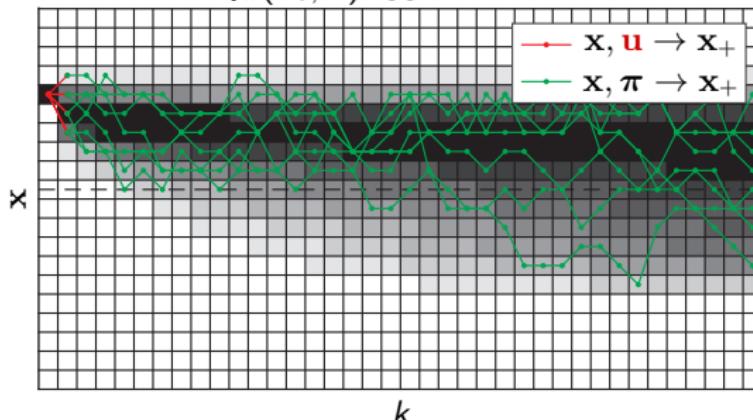
$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

$$Q_{\pi}(x_0, u) = 65.41$$



Action-value function

Bellman equation for the value function

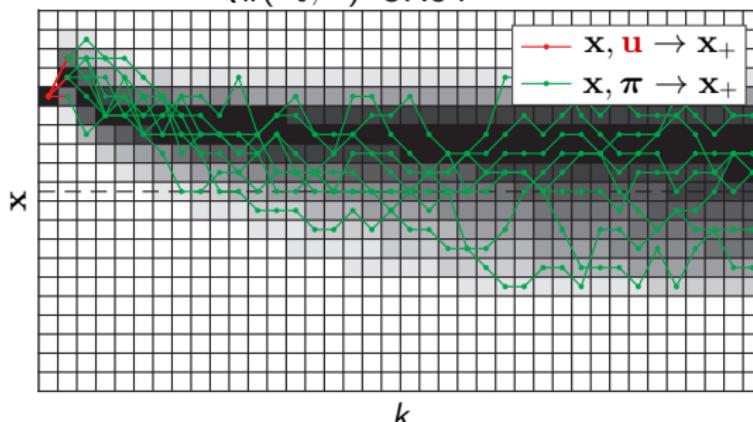
$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

$$Q_{\pi}(x_0, u) = 87.94$$



Action-value function

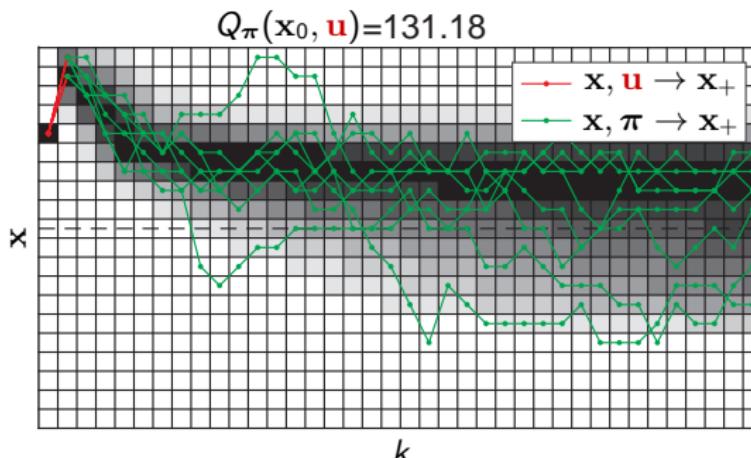
Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$



Action-value function

Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

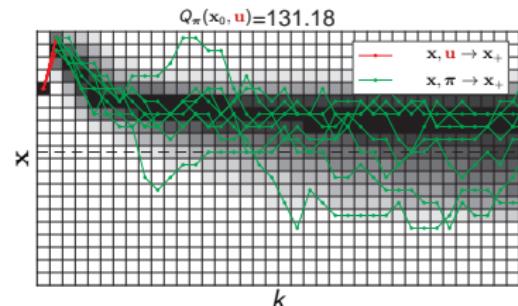
Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

Interpretation:

- From state x , use given input u , yields (stochastic) transition $x, u \rightarrow x_+$
- Then from x_+ onward follow policy π
- Action-value function $Q_{\pi}(x, u)$ is the expected return of the resulting trajectory



Action-value function

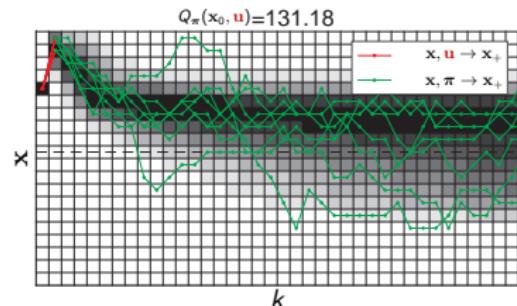
Bellman equation for the value function

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

Define

Action-value function:

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$



Interpretation:

- From state x , use given input u , yields (stochastic) transition $x, u \rightarrow x_+$
- Then from x_+ onward follow policy π
- Action-value function $Q_{\pi}(x, u)$ is the expected return of the resulting trajectory

Relationship Q_{π} and V_{π}

$$Q_{\pi}(x, \pi(x)) = V_{\pi}(x)$$

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Note: this is optimization over a policy, i.e. on the functions $\pi(x)$, rather than on inputs!

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Note: this is optimization over a policy, i.e. on the functions $\pi(x)$, rather than on inputs!

Note:

$$Q_{\pi_*}(x, u) = L(x, u) + \gamma \mathbb{E}[V_*(x_+) | x, u]$$

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Note: this is optimization over a policy, i.e. on the functions $\pi(x)$, rather than on inputs!

Note:

$$Q_{\pi_*}(x, u) = L(x, u) + \gamma \mathbb{E}[V_*(x_+) | x, u] \leq L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u] = Q_\pi(x, u)$$

for all x, u and π

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Note: this is optimization over a policy, i.e. on the functions $\pi(x)$, rather than on inputs!

Note:

$$Q_{\pi_*}(x, u) = L(x, u) + \gamma \mathbb{E}[V_*(x_+) | x, u] \leq L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u] = Q_\pi(x, u)$$

for all x, u and π

$Q_{\pi_*}(x, u) \leq Q_\pi(x, u)$
Hence minimizer π_* of V_π is
minimizer of Q_π

Optimal policy & value function

- Ordering over policies: π is better or equal to π' if $V_\pi(x) \leq V_{\pi'}(x)$ for all x
- Optimal policy π_* defined by $V_{\pi_*}(x) \leq V_\pi(x)$ for all x and π (maybe not unique)
- Define $V_*(x) = V_{\pi_*}(x)$ the optimal value function (unique)

Optimal value function

$$V_*(x) = \min_{\pi} V_{\pi}(x)$$

Note: this is optimization over a policy, i.e. on the functions $\pi(x)$, rather than on inputs!

Note:

$$Q_{\pi_*}(x, u) = L(x, u) + \gamma \mathbb{E}[V_*(x_+) | x, u] \leq L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u] = Q_\pi(x, u)$$

for all x, u and π

Optimal action-value function

$$Q_*(x, u) = \min_{\pi} Q_{\pi}(x, u) = Q_{\pi_*}(x, u)$$

$Q_{\pi_*}(x, u) \leq Q_\pi(x, u)$
Hence minimizer π_* of V_π is
minimizer of Q_π

Bellman optimality equation

Relationships

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Bellman optimality equation

Relationships

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Yield Bellman optimality equation for value function:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Bellman optimality equation

Relationships

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Yield Bellman optimality equation for value function:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

Note:

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

Bellman optimality equation

Relationships

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Yield Bellman optimality equation for value function:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Note:

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Yields Bellman optimality equation for action-value function:

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Bellman optimality equation

Relationships

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Yield Bellman optimality equation for value function:

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Note:

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

Yields Bellman optimality equation for action-value function:

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

Observe that $V_* \rightarrow Q_*$ exchanges \mathbb{E} and $\min_{\mathbf{u}}$

Advantage function

Associated to a policy:

$$A_\pi(x, u) = Q_\pi(x, u) - V_\pi(x)$$

Properties:

$$A_\pi(x, \pi(x)) = 0$$

$$A_\pi(x, u) < 0 \Rightarrow u \text{ better than } \pi(x)$$

- Assesses the value of an input u compared to a base policy (for given state)
- When $A_\pi(x, \pi(x)) \geq 0$ for all x , we have an optimal policy
- Will be useful for policy gradient methods

Note:

- $A_\pi(x, u) < 0 \Rightarrow u \text{ better than } \pi(x) \text{ for } x$
- $A_\pi(x, u) > 0 \Rightarrow u \text{ worse than } \pi(x) \text{ for } x$

allows for “comparing” arbitrary input u to baseline policy

Advantage function

Associated to a policy:

$$A_\pi(\mathbf{x}, \mathbf{u}) = Q_\pi(\mathbf{x}, \mathbf{u}) - V_\pi(\mathbf{x})$$

Properties:

$$A_\pi(\mathbf{x}, \pi(\mathbf{x})) = 0$$

$$A_\pi(\mathbf{x}, \mathbf{u}) < 0 \Rightarrow \mathbf{u} \text{ better than } \pi(\mathbf{x})$$

Optimal:

$$A_*(\mathbf{x}, \mathbf{u}) = Q_*(\mathbf{x}, \mathbf{u}) - V_*(\mathbf{x})$$

Properties:

$$A_*(\mathbf{x}, \pi_*(\mathbf{x})) = 0$$

$$A_*(\mathbf{x}, \mathbf{u}) \geq 0$$

- Assesses the value of an input \mathbf{u} compared to a base policy (for given state)
- When $A_\pi(\mathbf{x}, \pi(\mathbf{x})) \geq 0$ for all \mathbf{x} , we have an optimal policy
- Will be useful for policy gradient methods

Note:

- $A_\pi(\mathbf{x}, \mathbf{u}) < 0 \Rightarrow \mathbf{u} \text{ better than } \pi(\mathbf{x}) \text{ for } \mathbf{x}$
- $A_\pi(\mathbf{x}, \mathbf{u}) > 0 \Rightarrow \mathbf{u} \text{ worse than } \pi(\mathbf{x}) \text{ for } \mathbf{x}$

allows for “comparing” arbitrary input \mathbf{u} to baseline policy

Outline

- 1 Markov Decision Process
- 2 Value functions & Bellman equations
- 3 Function evaluation & Dynamic Programming
- 4 Quick view of Bellman contraction

How to compute value functions?

How to solve the Bellman equations?

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

$$Q_\pi(x, u) = L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u]$$

How to compute value functions?

How to solve the Bellman equations?

$$V_{\pi}(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, \pi]$$

$$V_{\pi}(x) = Q_{\pi}(x, \pi(x))$$

$$Q_{\pi}(x, u) = L(x, u) + \gamma \mathbb{E}[V_{\pi}(x_+) | x, u]$$

Policy evaluation: for a given policy π start with arbitrary value function $V_0(x)$, iterate:

$$V_{k+1}(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_k(x_+) | x, \pi] \quad (\text{sweep over } x)$$

How to compute value functions?

How to solve the Bellman equations?

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

$$Q_\pi(x, u) = L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u]$$

Policy evaluation: for a given policy π start with arbitrary value function $V_0(x)$, iterate:

$$V_{k+1}(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_k(x_+) | x, \pi] \quad (\text{sweep over } x)$$

Theorem

For $\gamma < 1$

$$\lim_{k \rightarrow \infty} V_k(x) = V_\pi(x)$$

holds for any $V_0(x)$ (finite)

How to compute value functions?

How to solve the Bellman equations?

$$V_\pi(x) = L(x, \pi(x)) + \gamma \mathbb{E}[V_\pi(x_+) | x, \pi]$$

$$V_\pi(x) = Q_\pi(x, \pi(x))$$

$$Q_\pi(x, u) = L(x, u) + \gamma \mathbb{E}[V_\pi(x_+) | x, u]$$

Policy evaluation: for a given policy π start with arbitrary value function $V_0(x)$, iterate:

$$V_{k+1}(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V_k(x_+) | x, \pi] \quad (\text{sweep over } x)$$

Theorem

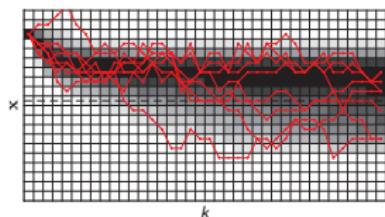
For $\gamma < 1$

$$\lim_{k \rightarrow \infty} V_k(x) = V_\pi(x)$$

holds for any $V_0(x)$ (finite)

Note: convergence for $\gamma = 1$ needs extra assumptions!

Computing V_π - Illustration

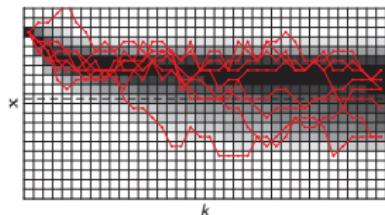


Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}

Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}

Algorithm: π -evaluation

Input: V (e.g. $= 0$) and $\text{tol} > 0$

Iter = True

while Iter **do**

for All \mathbf{x} **do**

$$V_+(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

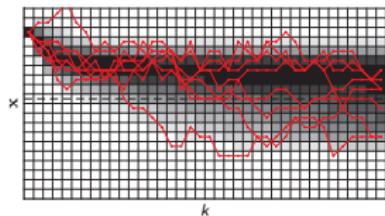
if $\|V_+ - V\|_\infty \leq \text{tol}$ **then**

 Iter = False

$$V \leftarrow V_+$$

return $V_\pi \leftarrow V_+$

Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}

Algorithm: π -evaluation

Input: V (e.g. = 0) and $\text{tol} > 0$

Iter = True

while Iter **do**

for All \mathbf{x} **do**

$$V_+(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \pi]$$

if $\|V_+ - V\|_\infty \leq \text{tol}$ **then**

 Iter = False

$$V \leftarrow V_+$$

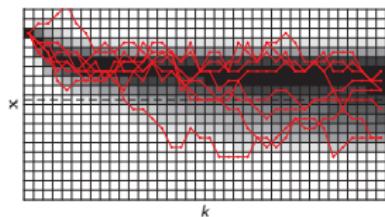
return $V_\pi \leftarrow V_+$

Remarks:

- Coding π -evaluation requires a discretized state-space $\mathbf{x} \in \mathbb{X}$
- V , V_+ are stored as tables (say of size n)
- Each iteration requires n^2 accesses
- Expected value for discrete state

$$\begin{aligned} \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \pi] \\ = \sum_{\mathbf{x}_+ \in \mathbb{X}} V(\mathbf{x}_+) \underbrace{\mathbb{P}[\mathbf{x}_+ | \mathbf{x}, \pi(\mathbf{x})]}_{\text{Model}} \end{aligned}$$

Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

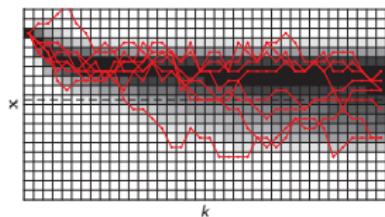
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

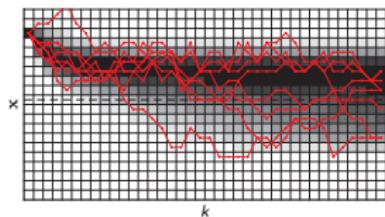
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

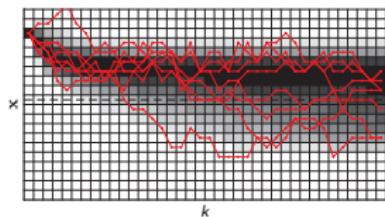
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

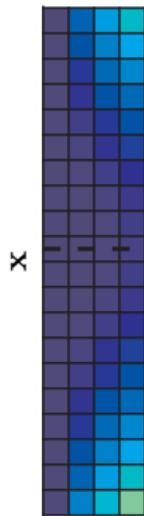
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

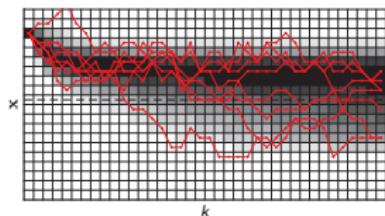
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

Computing V_π - Illustration



Start with e.g. $V(x) = 0$ and iterate

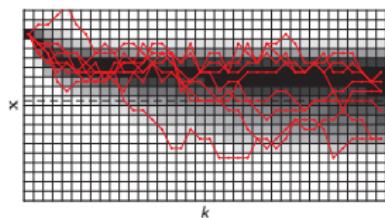
$$V(x) \leftarrow L(x, \pi(x)) + \gamma \mathbb{E}[V(x_+) | x, \pi]$$

as a sweep over all x



Iteration of π evaluation

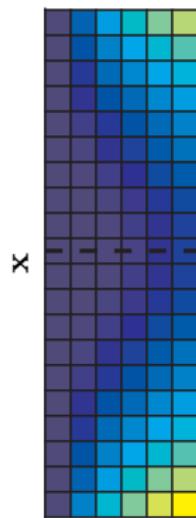
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

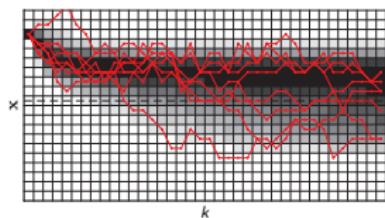
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

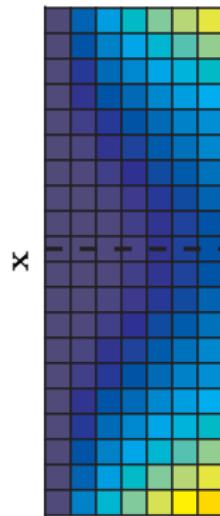
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

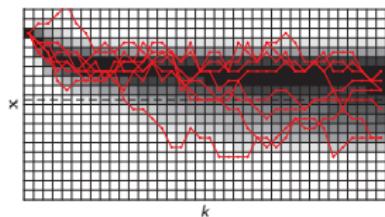
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

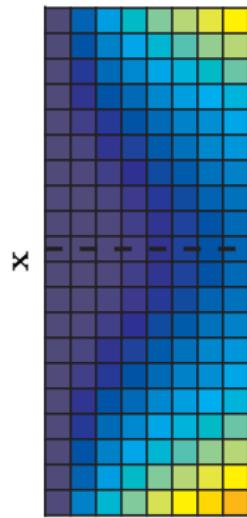
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

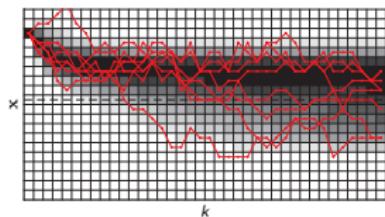
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

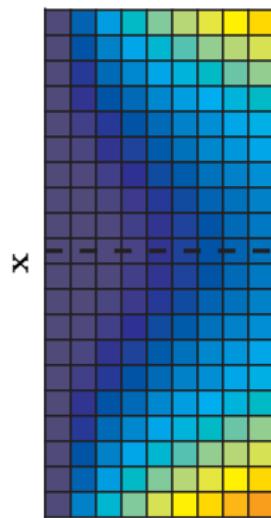
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

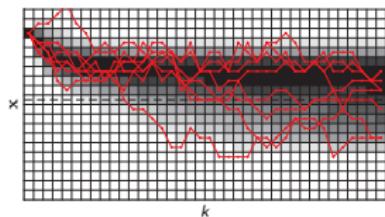
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

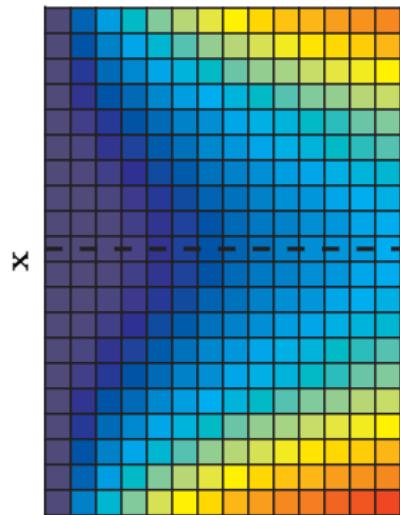
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

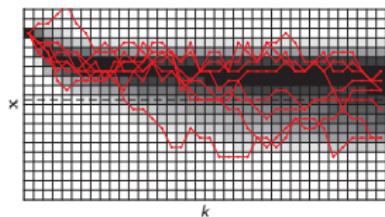
$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

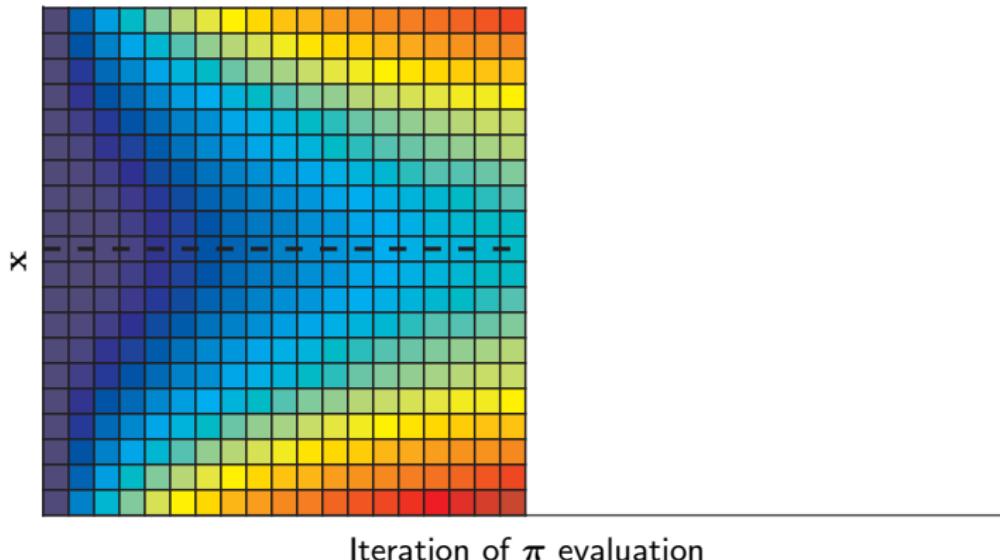
Computing V_π - Illustration



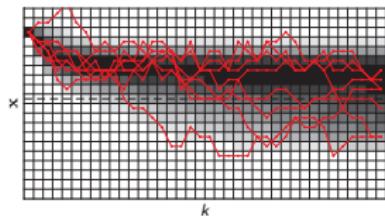
Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



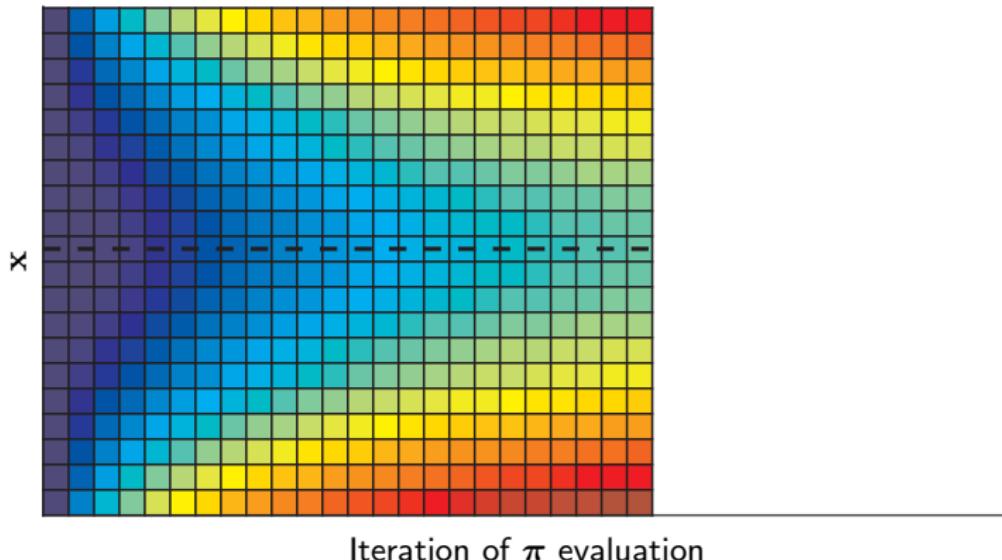
Computing V_π - Illustration



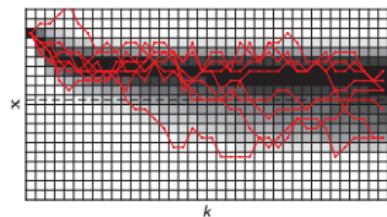
Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



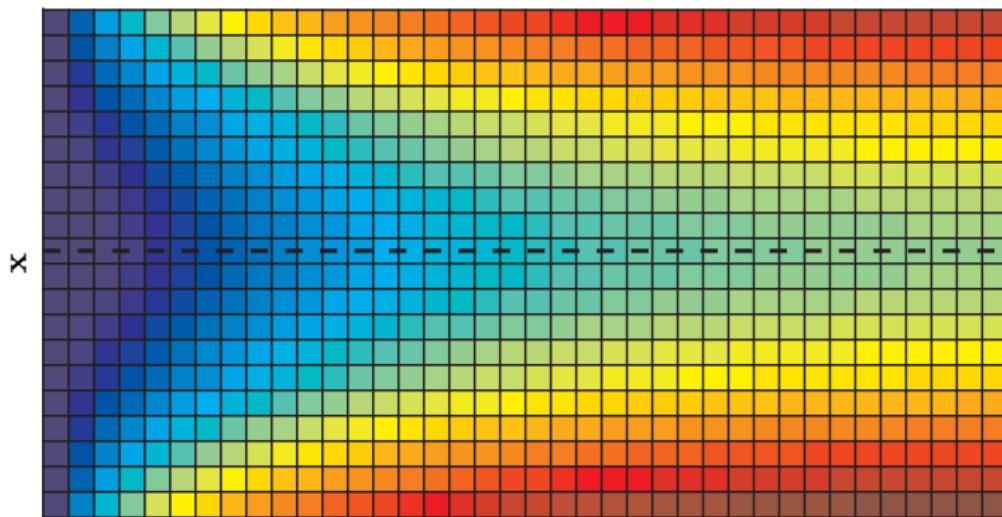
Computing V_π - Illustration



Start with e.g. $V(\mathbf{x}) = 0$ and iterate

$$V(\mathbf{x}) \leftarrow L(\mathbf{x}, \pi(\mathbf{x})) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \pi]$$

as a sweep over all \mathbf{x}



Iteration of π evaluation

Computing Q_π - Illustration

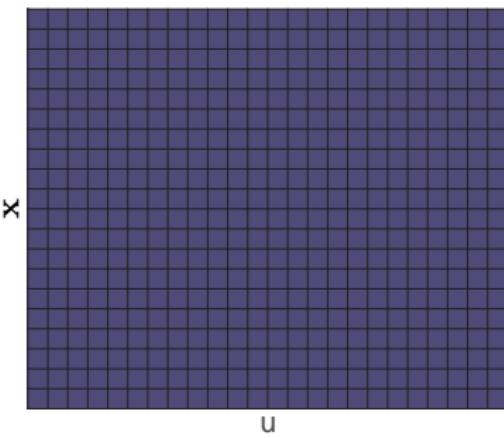
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

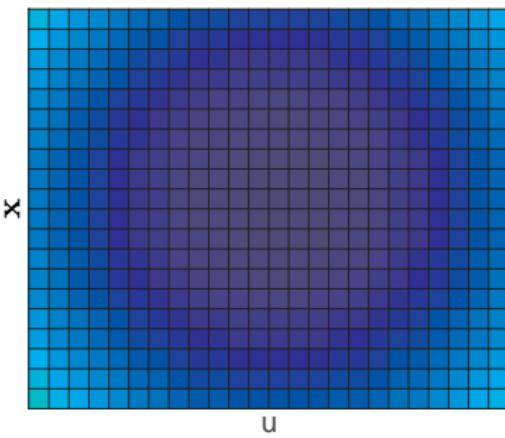
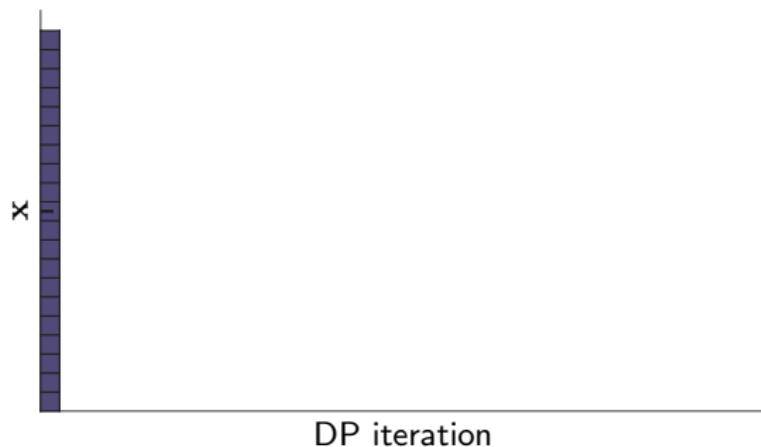
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

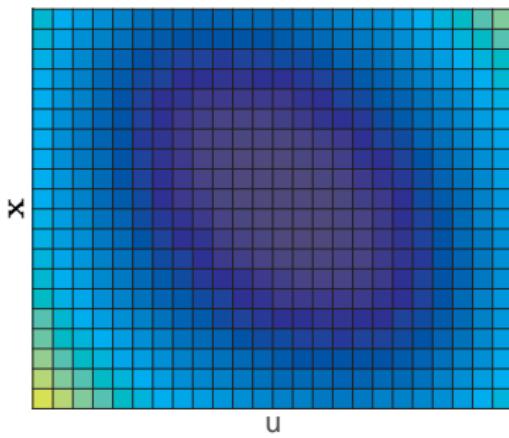
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

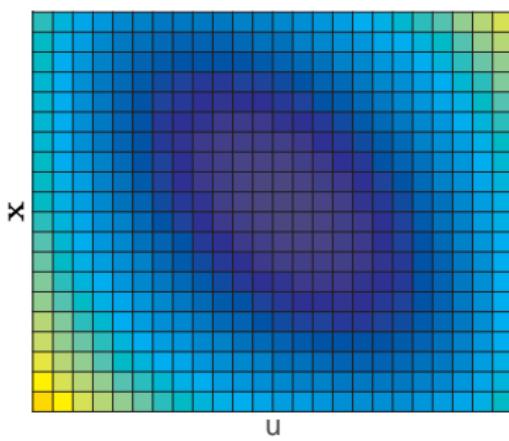
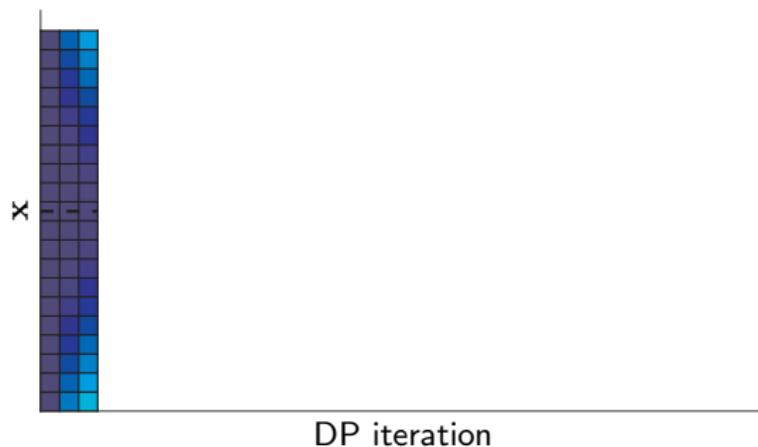
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

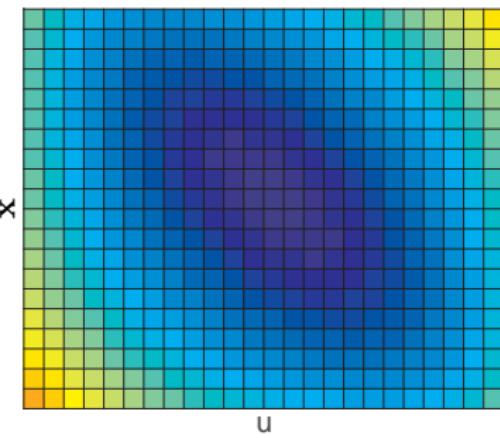
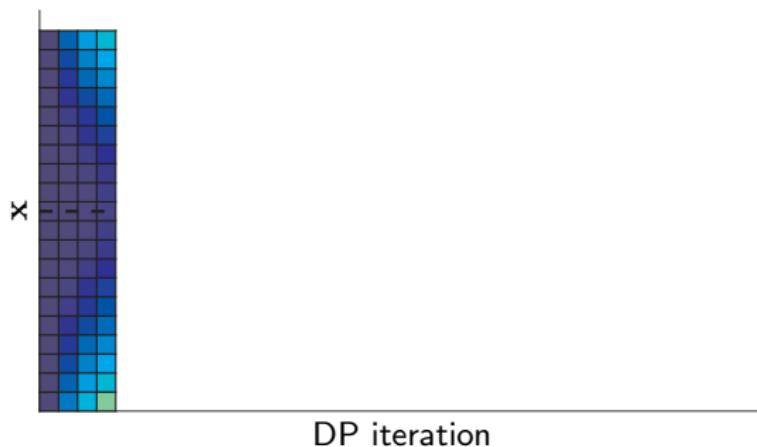
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

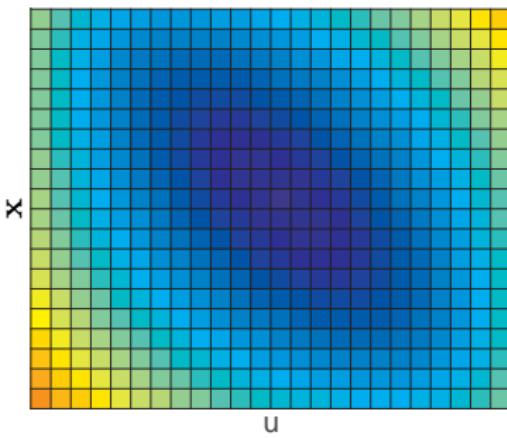
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

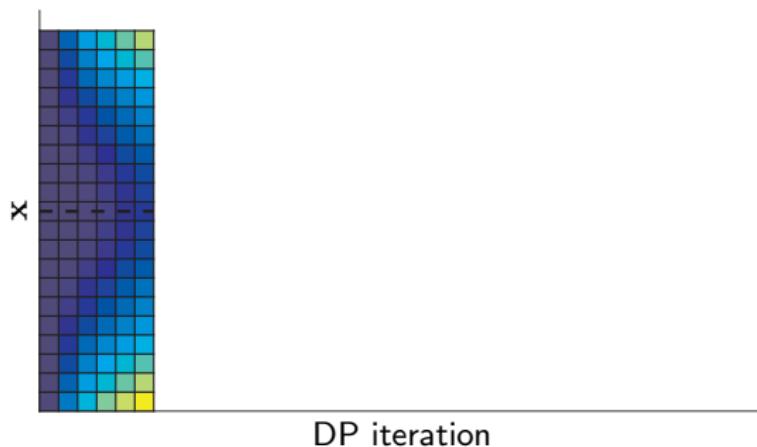
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

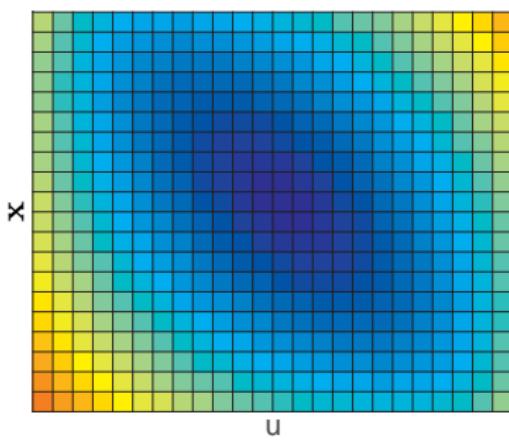
$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



DP iteration



Computing Q_π - Illustration

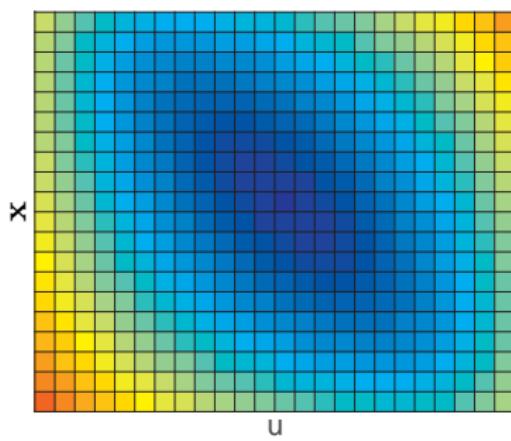
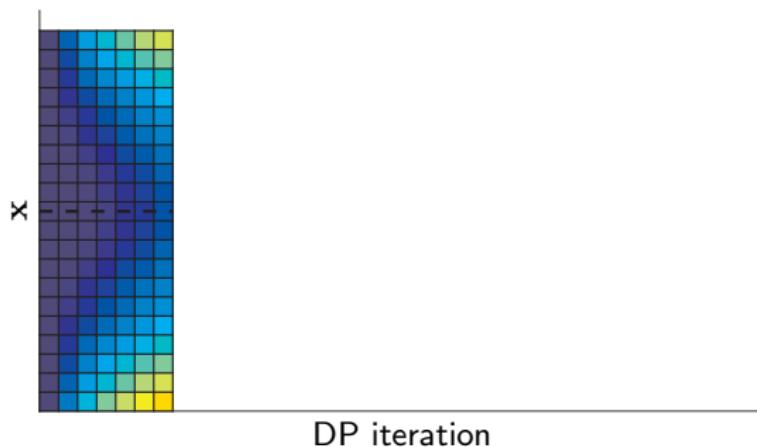
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

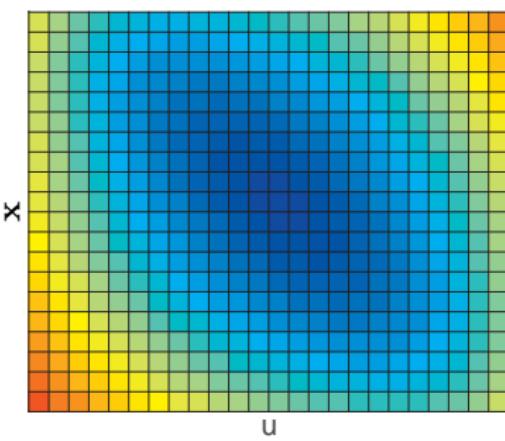
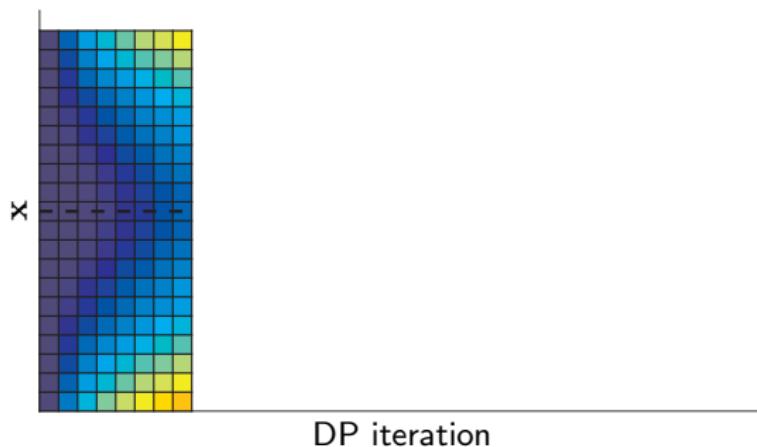
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

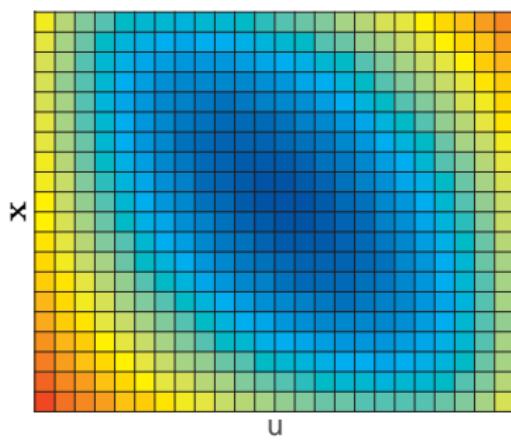
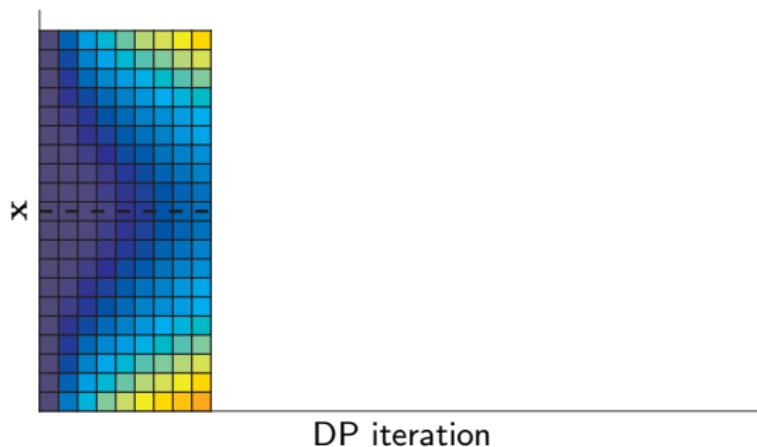
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

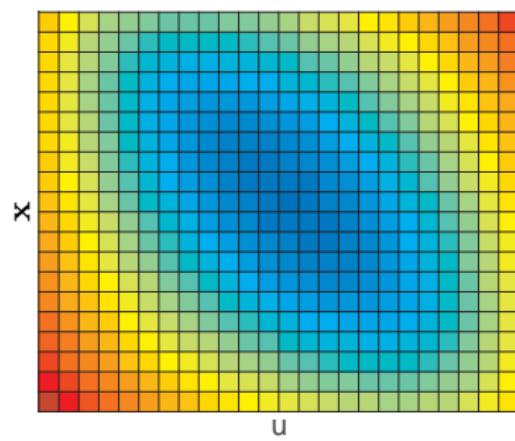
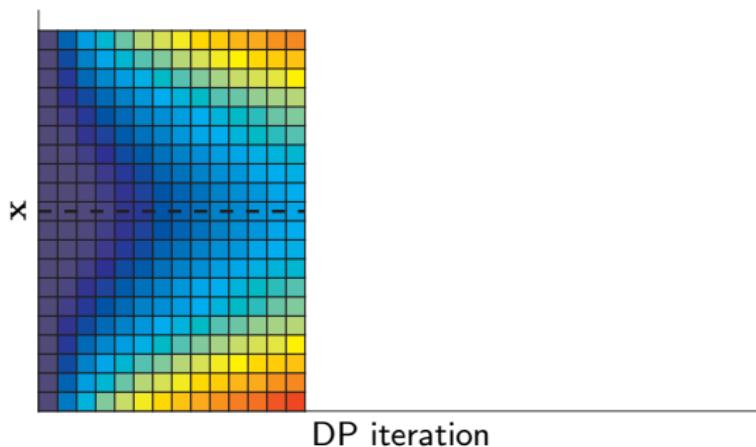
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

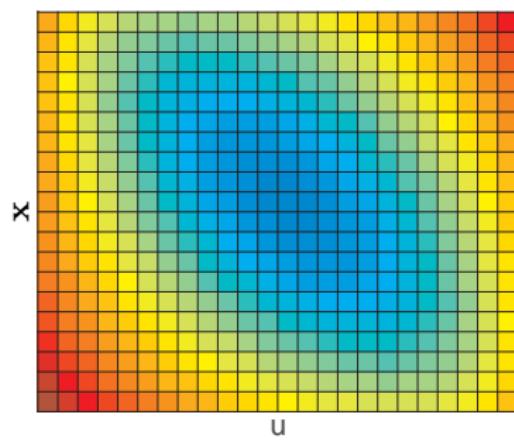
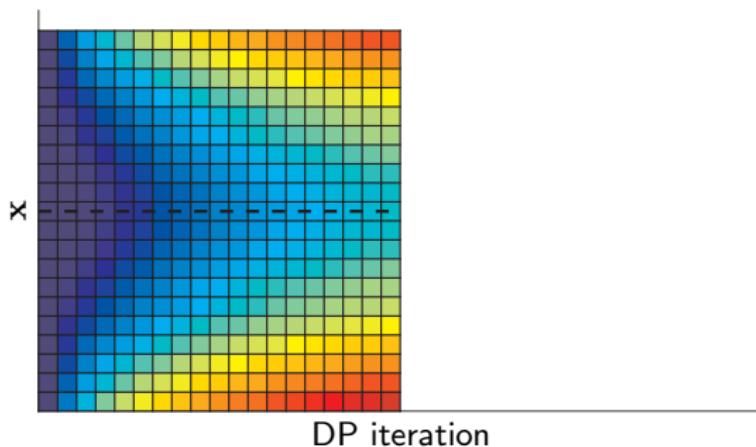
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

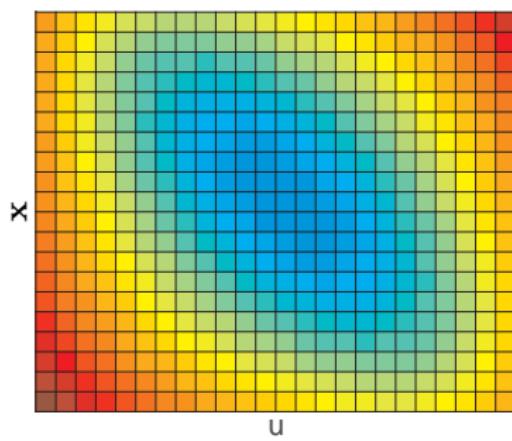
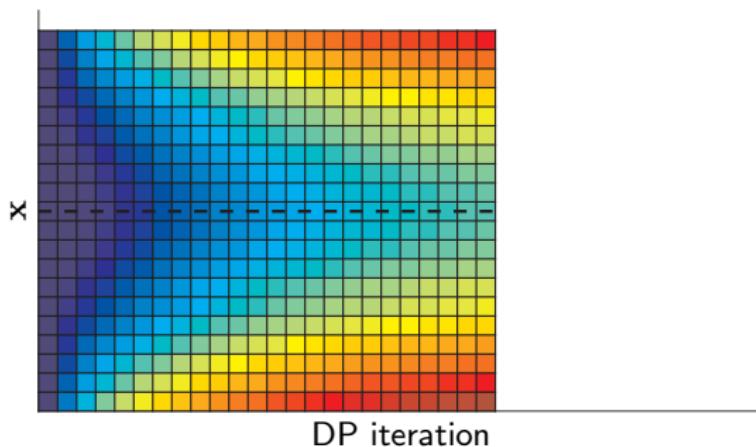
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



Computing Q_π - Illustration

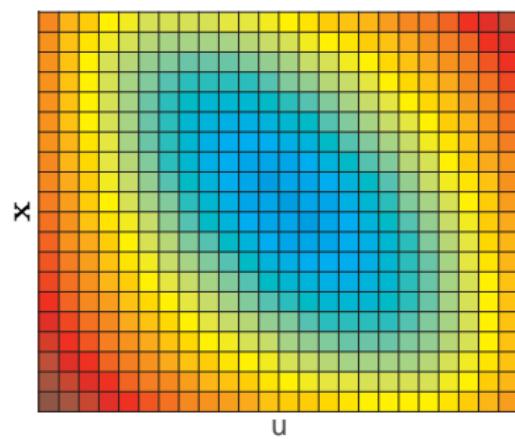
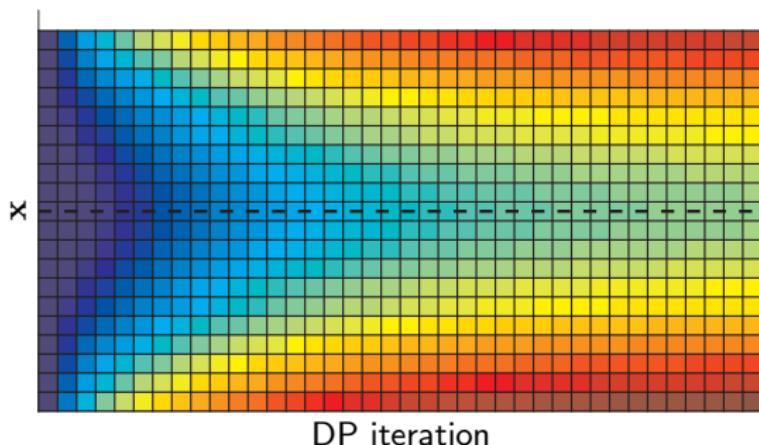
Can build V_π , Q_π jointly using:

$$Q(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow Q(\mathbf{x}, \pi(\mathbf{x}))$$

V , Q converge to V_π , Q_π

Need to sweep over \mathbf{x} and \mathbf{u} to build Q



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') \mid \mathbf{x}, \mathbf{u} \right]$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Policy improvement: take any given policy π with associated Q_π and define:

$$\pi'(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$$

How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Policy improvement: take any given policy π with associated Q_π and define:

$$\pi'(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$$

Then $V_{\pi'}(\mathbf{x}) \underbrace{\leq}_{\text{not obvious}} Q_\pi(\mathbf{x}, \pi'(\mathbf{x})) \leq Q_\pi(\mathbf{x}, \pi(\mathbf{x})) = V_\pi(\mathbf{x}), \quad \forall \mathbf{x}$

How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Policy improvement: take any given policy π with associated Q_π and define:

$$\pi'(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$$

Then $V_{\pi'}(\mathbf{x}) \underbrace{\leq}_{\text{not obvious}} Q_\pi(\mathbf{x}, \pi'(\mathbf{x})) \leq Q_\pi(\mathbf{x}, \pi(\mathbf{x})) = V_\pi(\mathbf{x}), \quad \forall \mathbf{x}$

π -improvement: select a $\pi(\mathbf{x})$, alternate

① Compute $Q_\pi(\mathbf{x})$

② Update policy

$$\pi(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$$

How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

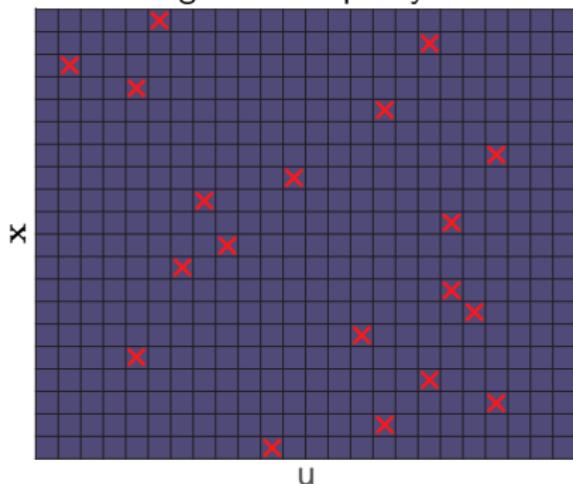
if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$

E.g. random policy π



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

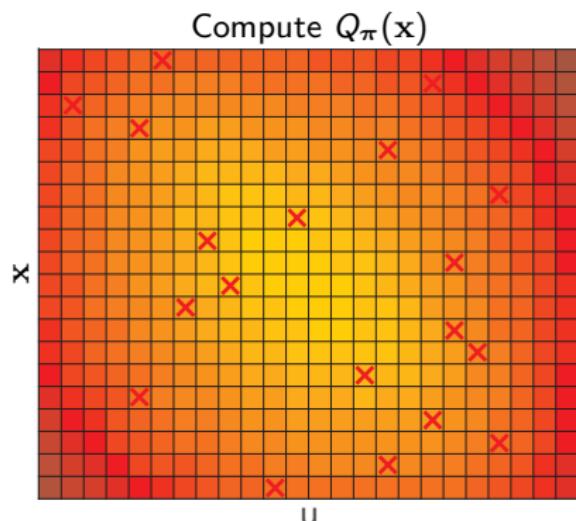
 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

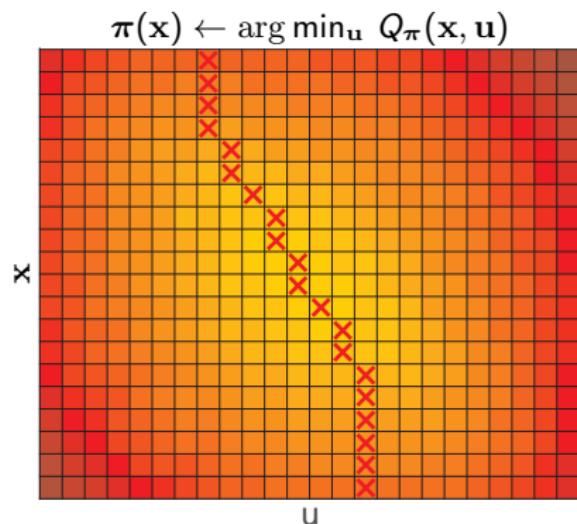
 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

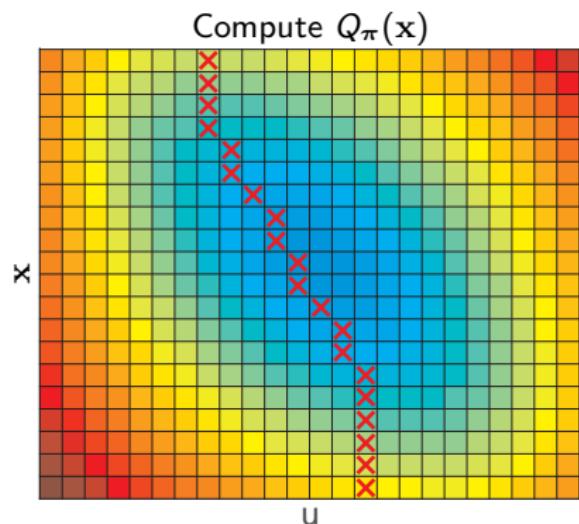
 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_*(\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

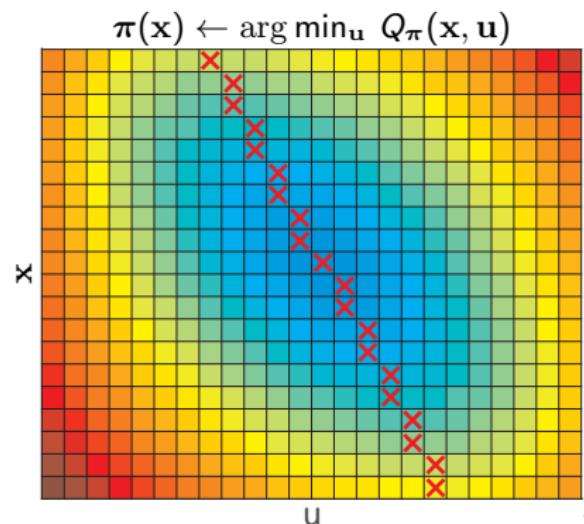
 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} [V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

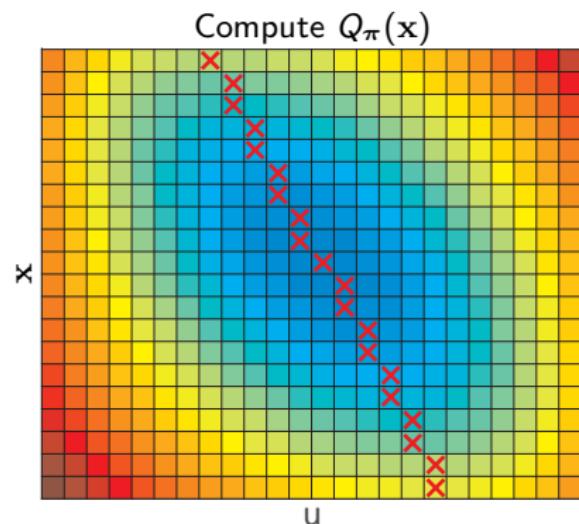
 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$



How to compute optimal policies?

How to solve the Bellman optimality equations?

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V_*(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$Q_*(\mathbf{x}, \mathbf{u}) = L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E} \left[\min_{\mathbf{u}'} Q_*(\mathbf{x}_+, \mathbf{u}') | \mathbf{x}, \mathbf{u} \right]$$

$$V_*(\mathbf{x}) = \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u}),$$

$$\pi_* (\mathbf{x}) = \arg \min_{\mathbf{u}} Q_*(\mathbf{x}, \mathbf{u})$$

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

 Update $\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

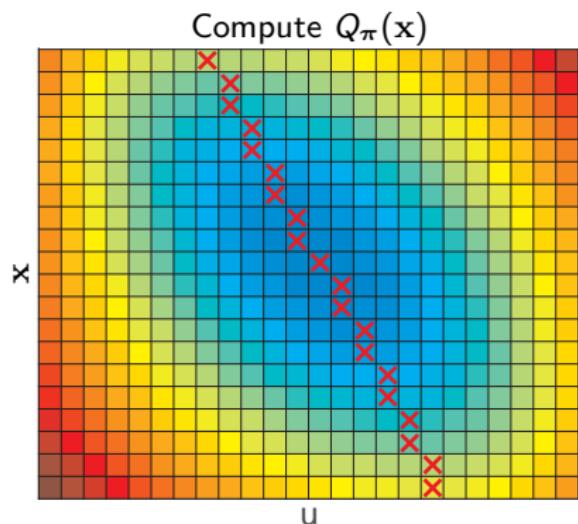
if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$

Converges to π_* if initial policy yields
 $Q_\pi < \infty$ for some \mathbf{x}, \mathbf{u}



Dynamic Programming

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

 Update

$\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$

Dynamic Programming

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

 Update

$\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$

$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

if $\|Q_+ - Q\| < \text{tol}$ **then**

 Iter = False

$Q \leftarrow Q_+$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

Dynamic Programming

Algorithm: π -improvement

Input: π and $\text{tol} > 0$

Iter = True

while Iter **do**

 Compute Q_π

 Update

$\pi_+(\mathbf{x}) \leftarrow \arg \min_{\mathbf{u}} Q_\pi(\mathbf{x}, \mathbf{u})$

if $\|\pi_+ - \pi\| < \text{tol}$ **then**

 Iter = False

$\pi_+ \leftarrow \pi$

return $\pi_* \leftarrow \pi_+$

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$

$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

if $\|Q_+ - Q\| < \text{tol}$ **then**

 Iter = False

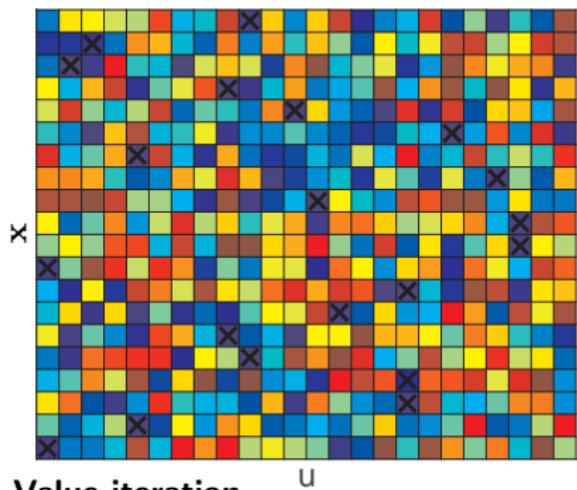
$Q \leftarrow Q_+$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

Value iteration

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

$\text{Iter} = \text{True}$, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

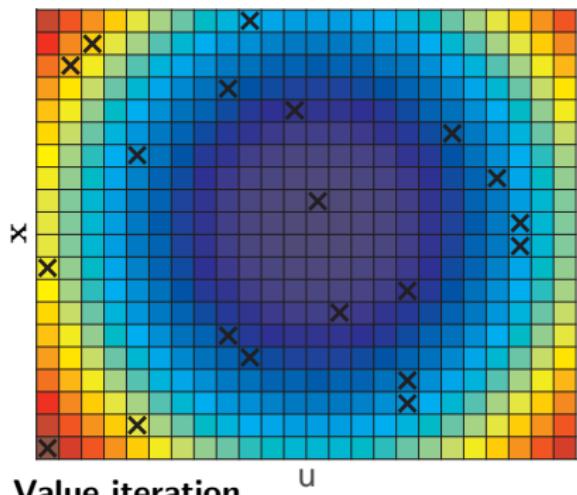
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

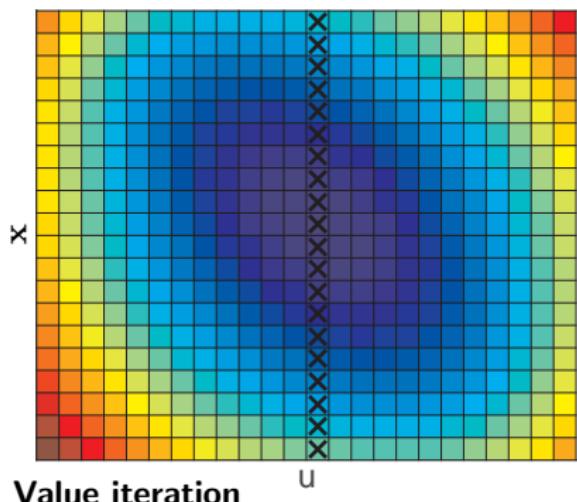
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

$\text{Iter} = \text{True}$, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

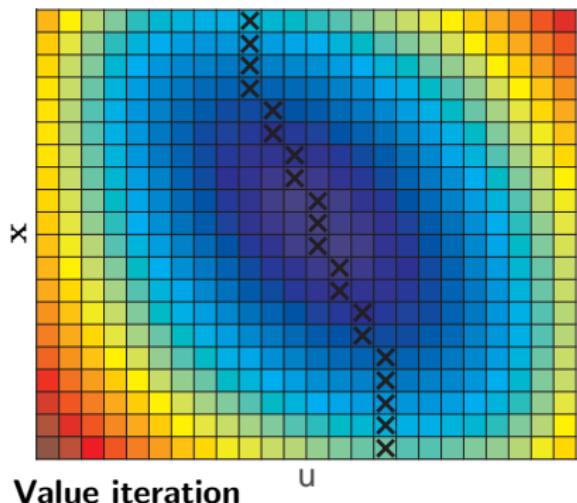
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

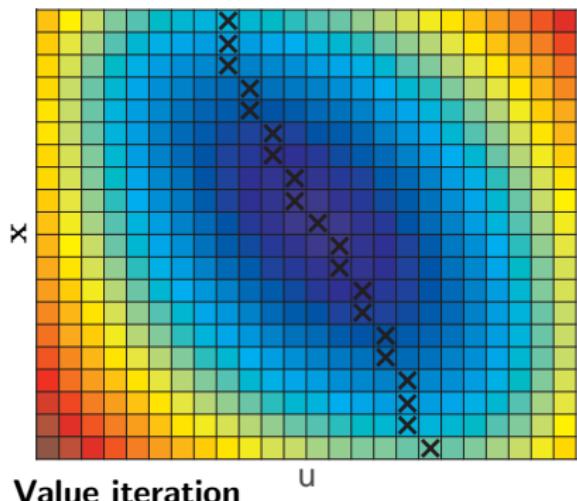
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

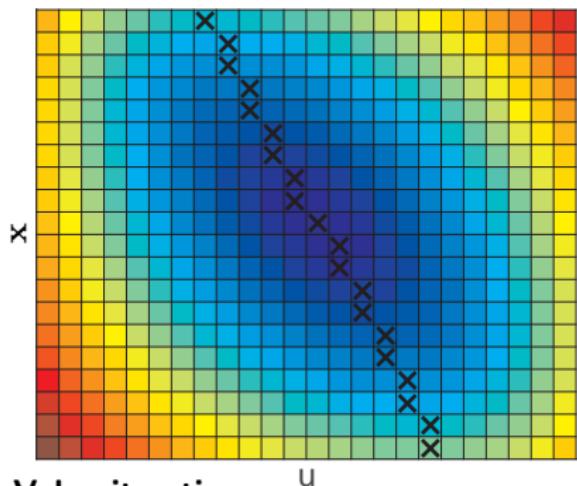
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

$\text{Iter} = \text{True}$, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

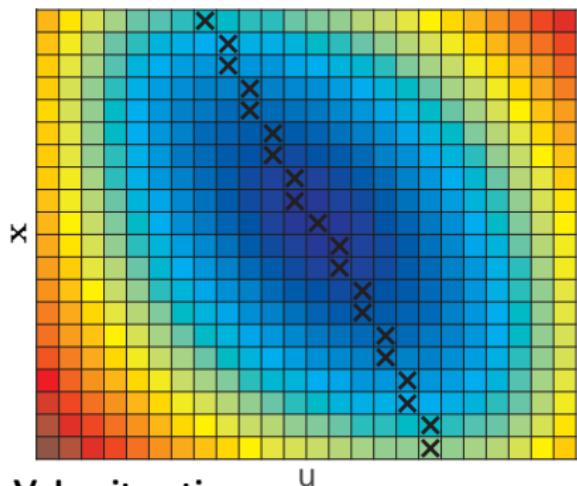
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

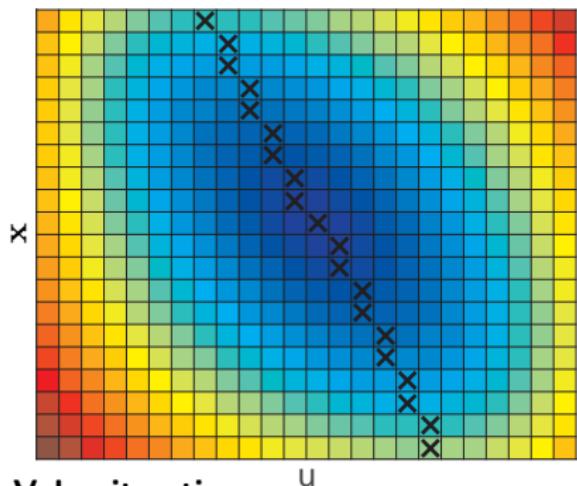
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

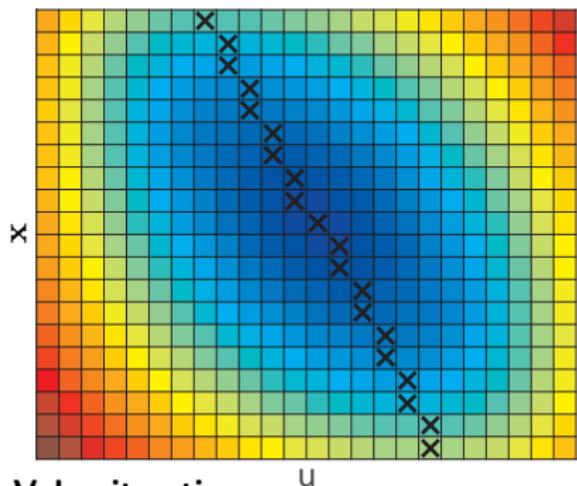
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

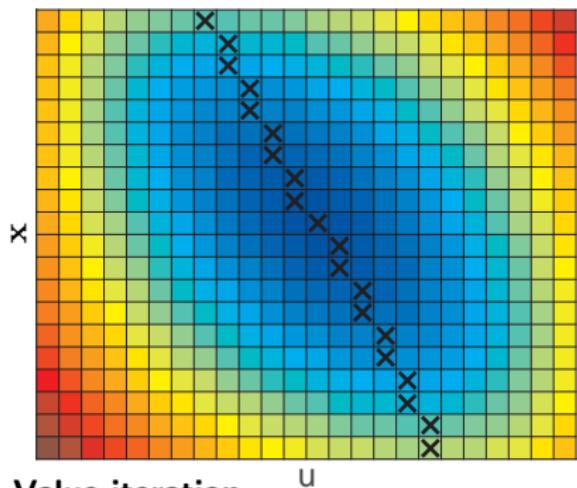
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

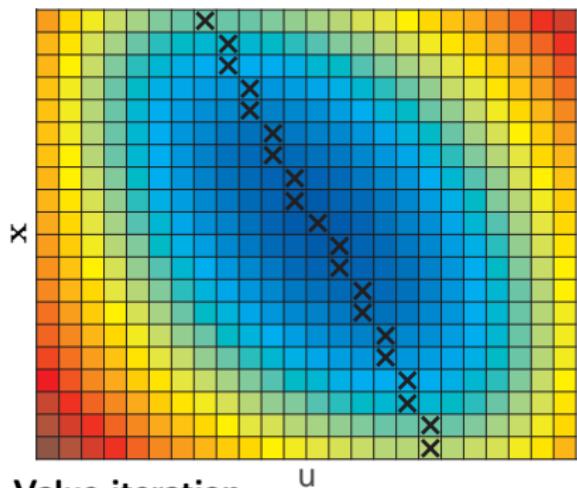
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

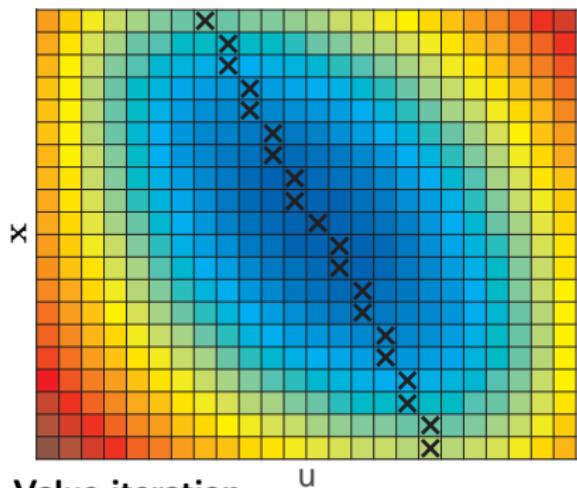
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

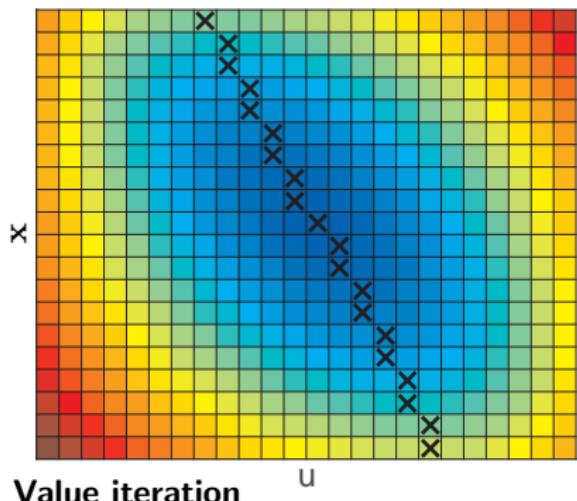
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

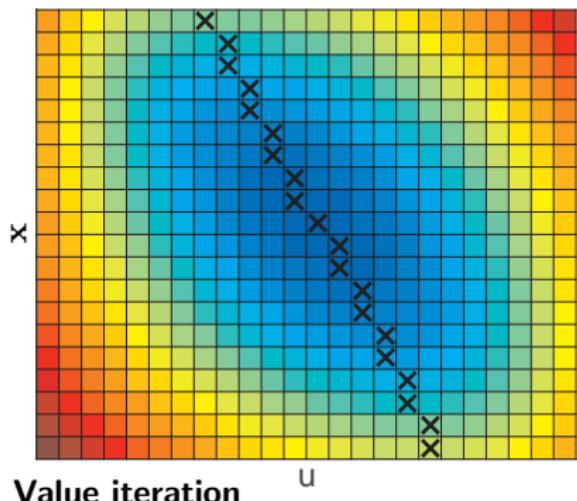
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

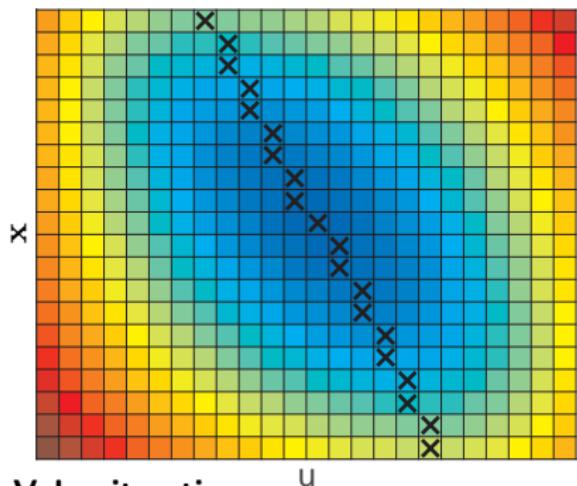
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

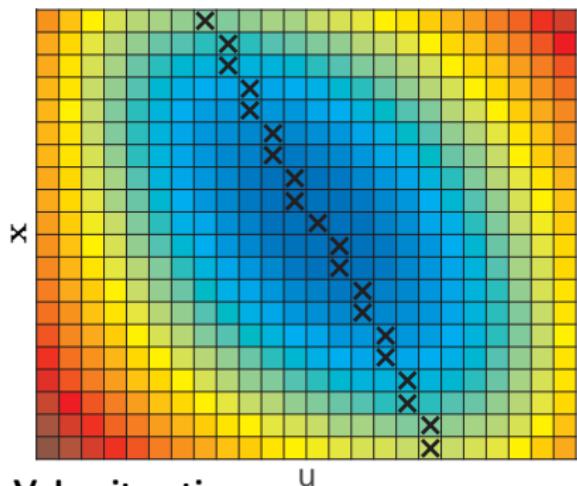
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

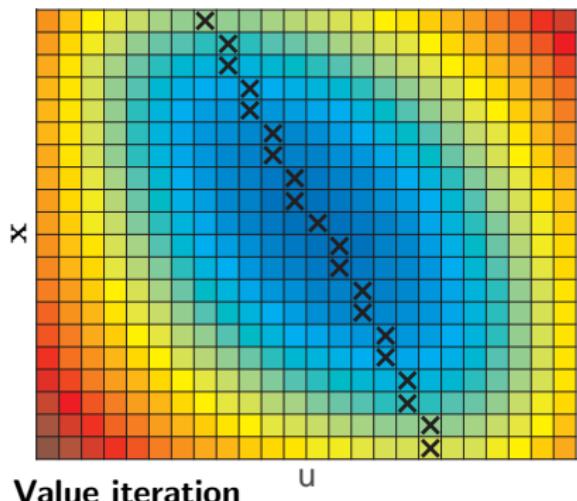
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

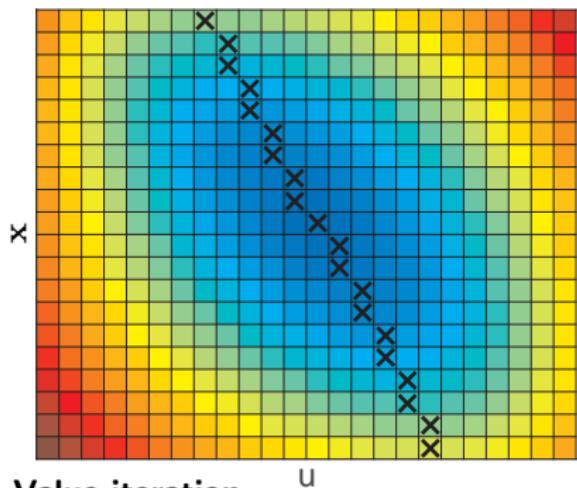
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

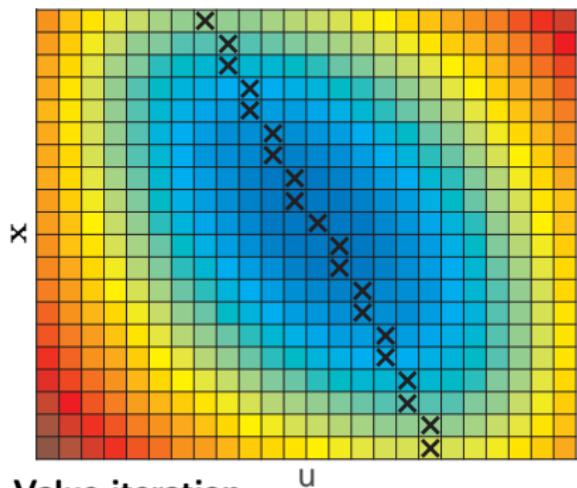
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

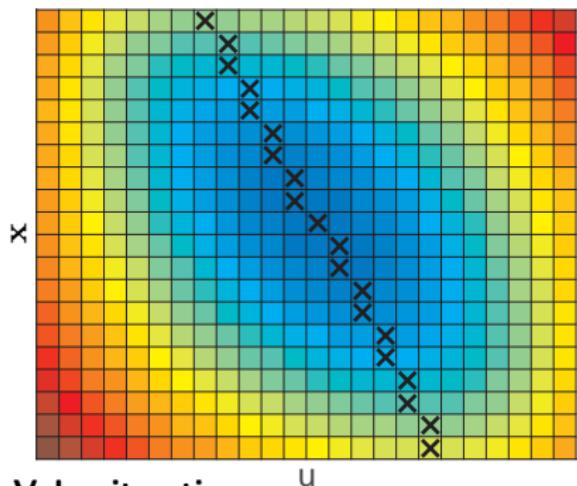
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Value iteration

Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

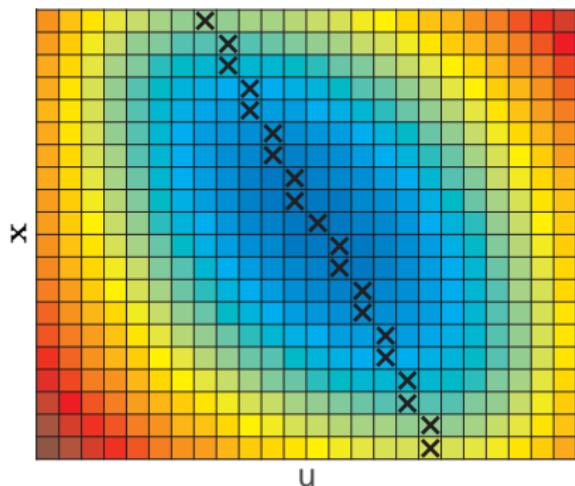
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Iteration often written as $V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$
- No need to compute Q_π (saves a lot of computation)
- Before convergence $V(\mathbf{x})$ may not correspond to any policy
- For $\gamma < 1$, converges to V_* , Q_* , π_* for any starting V (finite)
- For $\gamma = 1$, convergence requires extra assumptions (may depend on starting V)

Dynamic Programming



Algorithm: Value iteration

Input: V and $\text{tol} > 0$

Iter = True, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) \mid \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

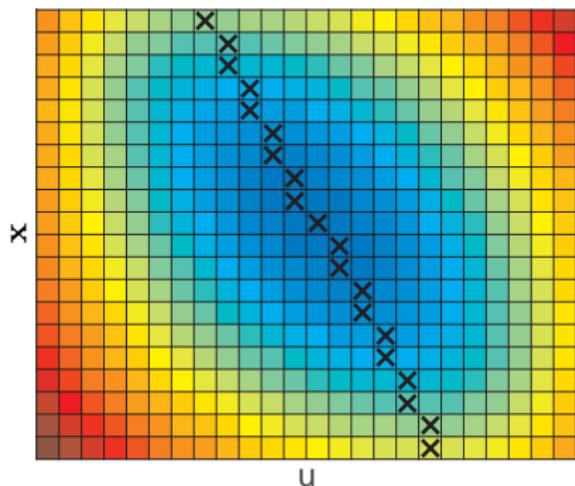
 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Policy π converges to π_* after a few iterations
- Action-value function Q converges to Q_* after (many) more iterations
- Hence π becomes correct “long” before Q_π is fully converged

Dynamic Programming



Algorithm: Value iteration

Input: V and $\text{tol} > 0$

$\text{Iter} = \text{True}$, $Q = 0$

while Iter **do**

 Compute

$$Q_+(\mathbf{x}, \mathbf{u}) \leftarrow L(\mathbf{x}, \mathbf{u}) + \gamma \mathbb{E}[V(\mathbf{x}_+) | \mathbf{x}, \mathbf{u}]$$

$$V(\mathbf{x}) \leftarrow \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$$

if $\|Q_+ - Q\| < \text{tol}$ **then**

 Iter = False

$$Q \leftarrow Q_+$$

return $\pi_* \leftarrow \arg \min_{\mathbf{u}} Q_+(\mathbf{x}, \mathbf{u})$

- Policy π converges to π_* after a few iterations
- Action-value function Q converges to Q_* after (many) more iterations
- Hence π becomes correct “long” before Q_π is fully converged

Value iteration often labelled Dynamic Programming, for $\gamma < 1$ convergence to optimal policy is guaranteed for any (finite) starting V

Curse of dimensionality

The Bellman equations & DP are very powerful mathematical tools

- One can treat *any* (possibly) stochastic optimal control problem
- The global optimum is achieved
- Policy $u = \pi(x)$ provides a feedback law, no extra online work needed

why don't we just use that?!

Curse of dimensionality

The Bellman equations & DP are very powerful mathematical tools

- One can treat *any* (possibly) stochastic optimal control problem
- The global optimum is achieved
- Policy $u = \pi(x)$ provides a feedback law, no extra online work needed

why don't we just use that?!

Dynamic Programming is “cursed”

Curse of dimensionality

The Bellman equations & DP are very powerful mathematical tools

- One can treat *any* (possibly) stochastic optimal control problem
- The global optimum is achieved
- Policy $u = \pi(x)$ provides a feedback law, no extra online work needed

why don't we just use that?!

Dynamic Programming is “cursed”

E.g.

- Imagine a system with **6 states, 2 inputs**
- Create **grids of 100 discrete values** for each

Curse of dimensionality

The Bellman equations & DP are very powerful mathematical tools

- One can treat *any* (possibly) stochastic optimal control problem
- The global optimum is achieved
- Policy $\mathbf{u} = \pi(\mathbf{x})$ provides a feedback law, no extra online work needed

why don't we just use that?!

Dynamic Programming is “cursed”

E.g.

- Imagine a system with **6 states, 2 inputs**
- Create **grids of 100 discrete values** for each

then e.g. storing $Q(\mathbf{x}, \mathbf{u})$ requires:

$$100^8 = 10^{16} \text{ numbers} \equiv 10000 \text{ terabytes of memory}$$

and each DP iteration requires going through 10^{32} numbers

Outline

- 1 Markov Decision Process
- 2 Value functions & Bellman equations
- 3 Function evaluation & Dynamic Programming
- 4 Quick view of Bellman contraction

Contraction of Bellman backup

- We can define the Bellman backup operator \mathcal{T}_π :

$$\mathcal{T}_\pi(V) = L(\cdot, \pi(\cdot)) + \mathbb{E}[V(x_+) | \cdot, \pi(\cdot)]$$

that transforms a function V into a new V . E.g. π -evaluation iterates $V \leftarrow \mathcal{T}_\pi(V)$

¹otherwise we need a bit of measure theory

Contraction of Bellman backup

- We can define the Bellman backup operator \mathcal{T}_π :

$$\mathcal{T}_\pi(V) = L(\cdot, \pi(\cdot)) + \mathbb{E}[V(x_+) | \cdot, \pi(\cdot)]$$

that transforms a function V into a new V . E.g. π -evaluation iterates $V \leftarrow \mathcal{T}_\pi(V)$

- We can define a norm over functions $\|V\|_\infty = \max_x |V(x)|$

¹otherwise we need a bit of measure theory

Contraction of Bellman backup

- We can define the Bellman backup operator \mathcal{T}_π :

$$\mathcal{T}_\pi(V) = L(\cdot, \pi(\cdot)) + \mathbb{E}[V(x_+) | \cdot, \pi(\cdot)]$$

that transforms a function V into a new V . E.g. π -evaluation iterates $V \leftarrow \mathcal{T}_\pi(V)$

- We can define a norm over functions $\|V\|_\infty = \max_x |V(x)|$

Then for V bounded¹...

The operator \mathcal{T}_π is a γ -**contraction**, i.e.

$$\|\mathcal{T}_\pi(V) - \mathcal{T}_\pi(V')\|_\infty \leq \gamma \|V - V'\|_\infty$$

This follows from observing that for any function V and any distribution, $\mathbb{E}[V] \leq \|V\|_\infty$

¹otherwise we need a bit of measure theory

Contraction of Bellman backup

- We can define the Bellman backup operator \mathcal{T}_π :

$$\mathcal{T}_\pi(V) = L(\cdot, \pi(\cdot)) + \mathbb{E}[V(x_+) | \cdot, \pi(\cdot)]$$

that transforms a function V into a new V . E.g. π -evaluation iterates $V \leftarrow \mathcal{T}_\pi(V)$

- We can define a norm over functions $\|V\|_\infty = \max_x |V(x)|$

Then for V bounded¹...

The operator \mathcal{T}_π is a γ -**contraction**, i.e.

$$\|\mathcal{T}_\pi(V) - \mathcal{T}_\pi(V')\|_\infty \leq \gamma \|V - V'\|_\infty$$

This follows from observing that for any function V and any distribution, $\mathbb{E}[V] \leq \|V\|_\infty$

Consequences

- For $\gamma < 1$ π -evaluation techniques converge to V_π regardless of the starting V
- Convergence is linear

¹otherwise we need a bit of measure theory

Contraction of Bellman Optimality backup

We can define the Bellman optimality backup operator \mathcal{T}_* :

$$\mathcal{T}_*(V) = \min_{\mathbf{u}} L(\cdot, \mathbf{u}) + \mathbb{E}[V(\mathbf{x}_+) | \cdot, \mathbf{u}]$$

that transforms a function V into a new V . E.g. DP iterates $V \leftarrow \mathcal{T}_*(V)$

²otherwise we need a bit of measure theory

Contraction of Bellman Optimality backup

We can define the Bellman optimality backup operator \mathcal{T}_* :

$$\mathcal{T}_*(V) = \min_{\mathbf{u}} L(\cdot, \mathbf{u}) + \mathbb{E}[V(\mathbf{x}_+) | \cdot, \mathbf{u}]$$

that transforms a function V into a new V . E.g. DP iterates $V \leftarrow \mathcal{T}_*(V)$

Then for V bounded²...

The operator \mathcal{T}_* is a γ -contraction, i.e.

$$\|\mathcal{T}_*(V) - \mathcal{T}_*(V')\|_\infty \leq \gamma \|V - V'\|_\infty$$

²otherwise we need a bit of measure theory

Contraction of Bellman Optimality backup

We can define the Bellman optimality backup operator \mathcal{T}_* :

$$\mathcal{T}_*(V) = \min_{\mathbf{u}} L(\cdot, \mathbf{u}) + \mathbb{E}[V(\mathbf{x}_+) | \cdot, \mathbf{u}]$$

that transforms a function V into a new V . E.g. DP iterates $V \leftarrow \mathcal{T}_*(V)$

Then for V bounded²...

The operator \mathcal{T}_* is a **γ -contraction**, i.e.

$$\|\mathcal{T}_*(V) - \mathcal{T}_*(V')\|_\infty \leq \gamma \|V - V'\|_\infty$$

Consequences

- For $\gamma < 1$ DP techniques converge to V_* regardless of the starting V
- Convergence is linear

²otherwise we need a bit of measure theory