Value Function Methods

CS 285: Deep Reinforcement Learning, Decision Making, and Control Sergey Levine

Class Notes

- 1. Homework 2 is due in one week (next Monday)
- 2. Remember to start forming final project groups and writing your proposal!
 - Proposal due 9/25, this Wednesday!

Today's Lecture

- 1. What if we just use a critic, without an actor?
- 2. Extracting a policy from a value function
- 3. The Q-learning algorithm
- 4. Extensions: continuous actions, improvements
- Goals:
 - Understand how value functions give rise to policies
 - Understand the Q-learning algorithm
 - Understand practical considerations for Q-learning

Recap: actor-critic

batch actor-critic algorithm:

1. sample {s_i, a_i} from π_θ(a|s) (run it on the robot)
2. fit V^π_φ(s) to sampled reward sums
3. evaluate Â^π(s_i, a_i) = r(s_i, a_i) + V^π_φ(s'_i) - V^π_φ(s_i)
4. ∇_θJ(θ) ≈ ∑_i ∇_θ log π_θ(a_i|s_i)Â^π(s_i, a_i)
5. θ ← θ + α∇_θJ(θ)



Can we omit policy gradient completely?

 $A^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$: how much better is \mathbf{a}_t than the average action according to π

 $\arg \max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$: best action from \mathbf{s}_t , if we then follow π

at *least* as good as any $\mathbf{a}_t \sim \pi(\mathbf{a}_t | \mathbf{s}_t)$

regardless of what $\pi(\mathbf{a}_t|\mathbf{s}_t)$ is!





fit A^{π} (or Q^{π} or V^{π})



Policy iteration

High level idea:

policy iteration algorithm:

1. evaluate $A^{\pi}(\mathbf{s}, \mathbf{a}) \longleftarrow$ how to do this? 2. set $\pi \leftarrow \pi'$

$$\pi'(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

as before: $A^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')] - V^{\pi}(\mathbf{s})$ let's evaluate $V^{\pi}(\mathbf{s})!$



Dynamic programming

Let's assume we know $p(\mathbf{s}'|\mathbf{s}, \mathbf{a})$, and \mathbf{s} and \mathbf{a} are both discrete (and small)

0.2	<mark>0.3</mark>	0.4	0.3
0.3	0.3	<mark>0.5</mark>	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5

16 states, 4 actions per state can store full $V^{\pi}(\mathbf{s})$ in a table! \mathcal{T} is $16 \times 16 \times 4$ tensor

bootstrapped update: $V^{\pi}(\mathbf{s}) \leftarrow E_{\mathbf{a} \sim \pi(\mathbf{a}|\mathbf{s})}[r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}'|\mathbf{s}, \mathbf{a})}[V^{\pi}(\mathbf{s}')]]$ just use the current estimate here

$$\pi'(\mathbf{a}_t|\mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg\max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases} \longrightarrow \text{ deterministic policy } \pi(\mathbf{s}) = \mathbf{a} \end{cases}$$

simplified: $V^{\pi}(\mathbf{s}) \leftarrow r(\mathbf{s}, \pi(\mathbf{s})) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s}))}[V^{\pi}(\mathbf{s}')]$

Policy iteration with dynamic programming



0.2	<mark>0.</mark> 3	0.4	0.3
0.3	0.3	<mark>0.</mark> 5	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5

16 states, 4 actions per state can store full $V^{\pi}(\mathbf{s})$ in a table! \mathcal{T} is $16 \times 16 \times 4$ tensor

Even simpler dynamic programming

$$\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

$$A^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')] - V^{\pi}(\mathbf{s})$$

 $\arg\max_{\mathbf{a}_t} A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \arg\max_{\mathbf{a}_t} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$

 $Q^{\pi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\pi}(\mathbf{s}')]$ (a bit simpler)

skip the policy and compute values directly!

value iteration algorithm:

1. set
$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$$

2. set $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$





$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \mathbf{a})}[V^{\pi}(\mathbf{s}')]$$



Fitted value iteration

how do we represent $V(\mathbf{s})$?

big table, one entry for each discrete **s** neural net function $V : S \to \mathbb{R}$



$$\mathcal{L}(\phi) = \frac{1}{2} \left\| V_{\phi}(\mathbf{s}) - \max_{\mathbf{a}} Q^{\pi}(\mathbf{s}, \mathbf{a}) \right\|^{2}$$

fitted value iteration algorithm:

1. set
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$

2. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$

$$s = 0: V(s) = 0.2$$

 $s = 1: V(s) = 0.3$
 $s = 2: V(s) = 0.5$



curse of dimensionality

 $|\mathcal{S}| = (255^3)^{200 \times 200}$ (more than atoms in the universe)

$$Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \mathbf{a})}[V^{\pi}(\mathbf{s}')]$$



What if we don't know the transition dynamics?

fitted value iteration algorithm:

1. set
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$

2. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$

need to know outcomes for different actions!

Back to policy iteration...

policy iteration: 1. evaluate $Q^{\pi}(\mathbf{s}, \mathbf{a})$ 2. set $\pi \leftarrow \pi'$ $\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$ policy evaluation: $V^{\pi}(\mathbf{s}) \leftarrow r(\mathbf{s}, \pi(\mathbf{s})) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \pi(\mathbf{s}))}[V^{\pi}(\mathbf{s}')]$ $Q^{\pi}(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E_{\mathbf{s}' \sim p(\mathbf{s}' | \mathbf{s}, \mathbf{a})}[Q^{\pi}(\mathbf{s}', \pi(\mathbf{s}'))]$ can fit this using samples

Can we do the "max" trick again?

policy iteration:

 $\begin{array}{c} \bullet \\ \bullet \\ 1. \text{ evaluate } V^{\pi}(\mathbf{s}) \\ 2. \text{ set } \pi \leftarrow \pi' \end{array}$

fitted value iteration algorithm:

1. set
$$\mathbf{y}_i \leftarrow \max_{\mathbf{a}_i} (r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)])$$

2. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|V_{\phi}(\mathbf{s}_i) - \mathbf{y}_i\|^2$

forget policy, compute value directly

can we do this with Q-values **also**, without knowing the transitions?

fitted Q iteration algorithm:

1. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)] \leftarrow q = approximate E[V(\mathbf{s}'_i)] \approx \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$ 2. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$ doesn't require simulation of actions!

+ works even for off-policy samples (unlike actor-critic)

+ only one network, no high-variance policy gradient

- no convergence guarantees for non-linear function approximation (more on this later)

Fitted Q-iteration

full fitted Q-iteration algorithm:

1. collect dataset
$$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$$
 using some policy
2. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$
3. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$

parameters

dataset size N, collection policy iterations Kgradient steps S



Review

- Value-based methods
 - Don't learn a policy explicitly
 - Just learn value or Q-function
- If we have value function, we have a policy
- Fitted Q-iteration



Break

Why is this algorithm off-policy?

full fitted Q-iteration algorithm:

1. collect dataset $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ using some policy 2. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma(\max_{\mathbf{a}'_i} Q_\phi(\mathbf{s}'_i, \mathbf{a}'_i))$ 3. set $\phi \leftarrow \arg\min_{\phi} \frac{1}{2} \sum_i ||Q_\phi(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i||^2$ given \mathbf{s} and \mathbf{a} , transition is independent of π this approximates the value of π' at \mathbf{s}'_i $\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg\max_{\mathbf{a}_t} Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$



What is fitted Q-iteration optimizing?

full fitted Q-iteration algorithm:

1. collect dataset $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$ using some policy 2. set $\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'_i} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i) \leftarrow \mathbf{this} \max$ improves the policy (tabular case) 3. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i ||Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i||^2$ error \mathcal{E} $\mathcal{E} = \frac{1}{2} E_{(\mathbf{s}, \mathbf{a}) \sim \beta} \left[\left(Q_{\phi}(\mathbf{s}, \mathbf{a}) - [r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}', \mathbf{a}')] \right)^2 \right]$

if $\mathcal{E} = 0$, then $Q_{\phi}(\mathbf{s}, \mathbf{a}) = r(\mathbf{s}, \mathbf{a}) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}', \mathbf{a}')$

this is an *optimal* Q-function, corresponding to optimal policy π' :

$$\pi'(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) & \text{maximizes reward} \\ 0 \text{ otherwise} & \text{sometimes written } Q^* \text{ and } \pi^* \end{cases}$$

most guarantees are lost when we leave the tabular case (e.g., when we use neural network function approximation)

Online Q-learning algorithms

full fitted Q-iteration algorithm:





off policy, so many choices here!

online Q iteration algorithm:

1. take some action
$$\mathbf{a}_i$$
 and observe $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$

2.
$$\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$$

3.
$$\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$$

Exploration with Q-learning

online Q iteration algorithm:

1. take some action \mathbf{a}_i and observe $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ 2. $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$ 3. $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$

final policy:

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_\phi(\mathbf{s}_t, \mathbf{a}_t) \\ 0 \text{ otherwise} \end{cases}$$

why is this a bad idea for step 1?

$$\pi(\mathbf{a}_t | \mathbf{s}_t) = \begin{cases} 1 - \epsilon \text{ if } \mathbf{a}_t = \arg \max_{\mathbf{a}_t} Q_{\phi}(\mathbf{s}_t, \mathbf{a}_t) \\ \epsilon/(|\mathcal{A}| - 1) \text{ otherwise} \end{cases}$$

"epsilon-greedy"

 $\pi(\mathbf{a}_t|\mathbf{s}_t) \propto \exp(Q_\phi(\mathbf{s}_t,\mathbf{a}_t))$

"Boltzmann exploration"

We'll discuss exploration in more detail in a later lecture!

Review

- Value-based methods
 - Don't learn a policy explicitly
 - Just learn value or Q-function
- If we have value function, we have a policy
- Fitted Q-iteration
 - Batch mode, off-policy method
- Q-learning
 - Online analogue of fitted Qiteration



Value function learning theory

value iteration algorithm:

1. set
$$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$$

2. set $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$

0.20.30.40.30.30.30.50.30.40.40.60.40.50.50.70.5

does it converge?

and if so, to what?

stacked vector of rewards at all states for action **a** define an operator \mathcal{B} : $\mathcal{B}V = \max_{\mathbf{a}} r_{\mathbf{a}} + \gamma \mathcal{T}_{\mathbf{a}} V$ matrix of transitions for action **a** such that $\mathcal{T}_{\mathbf{a},i,j} = p(\mathbf{s}' = i | \mathbf{s} = j, \mathbf{a})$

 V^{\star} is a fixed point of \mathcal{B} $V^{\star}(\mathbf{s}) = \max_{\mathbf{a}} r(\mathbf{s}, \mathbf{a}) + \gamma E[V^{\star}(\mathbf{s}')], \text{ so } V^{\star} = \mathcal{B}V^{\star}$

always exists, is always unique, always corresponds to the optimal policy

...but will we reach it?

Value function learning theory

value iteration algorithm:

1. set $Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \gamma E[V(\mathbf{s}')]$ 2. set $V(\mathbf{s}) \leftarrow \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a})$

 V^* is a fixed point of \mathcal{B} $V^*(\mathbf{s}) = \max_{\mathbf{a}} r(\mathbf{s}, \mathbf{a}) + \gamma E[V^*(\mathbf{s}')], \text{ so } V^* = \mathcal{B}V^*$

we can prove that value iteration reaches V^* because \mathcal{B} is a *contraction*

contraction: for any V and \overline{V} , we have $\|\mathcal{B}V - \mathcal{B}\overline{V}\|_{\infty} \leq \gamma \|V - \overline{V}\|_{\infty}$ gap always gets smaller by γ ! (with respect to ∞ -norm)

what if we choose V^* as \overline{V} ? $\mathcal{B}V^* = V^*$!

 $\|\mathcal{B}V - V^{\star}\|_{\infty} \leq \gamma \|V - V^{\star}\|_{\infty}$

0.2	<mark>0.</mark> 3	0.4	0.3
0.3	0.3	<mark>0.</mark> 5	0.3
0.4	0.4	0.6	0.4
0.5	0.5	0.7	0.5



Non-tabular value function learning

value iteration algorithm (using \mathcal{B}): 1. $V \leftarrow \mathcal{B}V$

fitted value iteration algorithm (using \mathcal{B} and Π): 1. $V \leftarrow \Pi \mathcal{B} V$ fitted value iteration algorithm:

define new operator Π : $\Pi V = \arg \min_{V' \in \Omega} \frac{1}{2} \sum ||V'(\mathbf{s}) - V(\mathbf{s})||^2$ Π is a *projection* onto Ω (in terms of ℓ_2 norm)

updated value function

$$\mathcal{B}V \qquad V' \leftarrow \arg\min_{V' \in \Omega} \frac{1}{2} \sum \|V'(\mathbf{s}) - (\mathcal{B}V)(\mathbf{s})\|^2$$
set Ω (e.g., neural nets)
all value functions represented by, e.g., neural nets

Non-tabular value function learning

fitted value iteration algorithm (using \mathcal{B} and Π): 1. $V \leftarrow \Pi \mathcal{B} V$

 \mathcal{B} is a contraction w.r.t. ∞ -norm ("max" norm)

 Π is a contraction w.r.t. ℓ_2 -norm (Euclidean distance)

but... $\Pi \mathcal{B}$ is not a contraction of any kind





 $\|\mathcal{B}V - \mathcal{B}\bar{V}\|_{\infty} \le \gamma \|V - \bar{V}\|_{\infty}$ $\|\Pi V - \Pi\bar{V}\|^{2} \le \|V - \bar{V}\|^{2}$

Conclusions:

value iteration converges (tabular case) fitted value iteration does **not** converge not in general often not in practice

What about fitted Q-iteration?

fitted Q iteration algorithm:

1. set
$$\mathbf{y}_i \leftarrow r(\mathbf{s}_i, \mathbf{a}_i) + \gamma E[V_{\phi}(\mathbf{s}'_i)]$$

2. set $\phi \leftarrow \arg \min_{\phi} \frac{1}{2} \sum_i \|Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i\|^2$

define an operator $\mathcal{B}: \mathcal{B}Q = r + \gamma \mathcal{T} \max_{\mathbf{a}} Q$

max now after the transition operator

define an operator Π : $\Pi Q = \arg \min_{Q' \in \Omega} \frac{1}{2} \sum \|Q'(\mathbf{s}, \mathbf{a}) - Q(\mathbf{s}, \mathbf{a})\|^2$

fitted Q-iteration algorithm (using \mathcal{B} and Π): 1. $Q \leftarrow \Pi \mathcal{B} Q$

Applies also to online Q-learning

 \mathcal{B} is a contraction w.r.t. ∞ -norm ("max" norm) Π is a contraction w.r.t. ℓ_2 -norm (Euclidean distance) $\Pi \mathcal{B}$ is not a contraction of any kind

But... it's just regression!

online Q iteration algorithm:

1. take some action \mathbf{a}_i and observe $(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)$ 2. $\mathbf{y}_i = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i)$ 3. $\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i)(Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - \mathbf{y}_i)$ isn't this just gradient descent? that converges, right?

ue

Q-learning is *not* gradient descent!

$$\phi \leftarrow \phi - \alpha \frac{dQ_{\phi}}{d\phi}(\mathbf{s}_i, \mathbf{a}_i) (Q_{\phi}(\mathbf{s}_i, \mathbf{a}_i) - r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \max_{\mathbf{a}'} Q_{\phi}(\mathbf{s}'_i, \mathbf{a}'_i))$$

no gradient through target val

A sad corollary

batch actor-critic algorithm:

1. sample {s_i, a_i} from π_θ(a|s) (run it on the robot)
2. fit V^π_φ(s) to sampled reward sums
3. evaluate Â^π(s_i, a_i) = r(s_i, a_i) + V^π_φ(s'_i) - V^π_φ(s_i)
4. ∇_θJ(θ) ≈ ∑_i ∇_θ log π_θ(a_i|s_i)Â^π(s_i, a_i)
5. θ ← θ + α∇_θJ(θ)



An aside regarding terminology

 V^{π} : value function for policy π this is what the critic does

 V^* : value function for optimal policy π^* this is what value iteration does

fitted bootstrapped policy evaluation doesn't converge!

Review

- Value iteration theory
 - Linear operator for backup
 - Linear operator for projection
 - Backup is contraction
 - Value iteration converges
- Convergence with function approximation
 - Projection is also a contraction
 - Projection + backup is **not** a contraction
 - Fitted value iteration does not in general converge
- Implications for Q-learning
 - Q-learning, fitted Q-iteration, etc. does not converge with function approximation
- But we can make it work in practice!
 - Sometimes tune in next time

