

AIL 722: Reinforcement Learning

Lecture 11: Fitted Value Iteration

Raunak Bhattacharyya



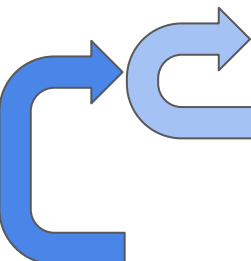
ScAI

YARDI SCHOOL OF ARTIFICIAL INTELLIGENCE
INDIAN INSTITUTE OF TECHNOLOGY DELHI

Outline

- Approximating the value function
- Fitted value iteration
- Fitted Q iteration

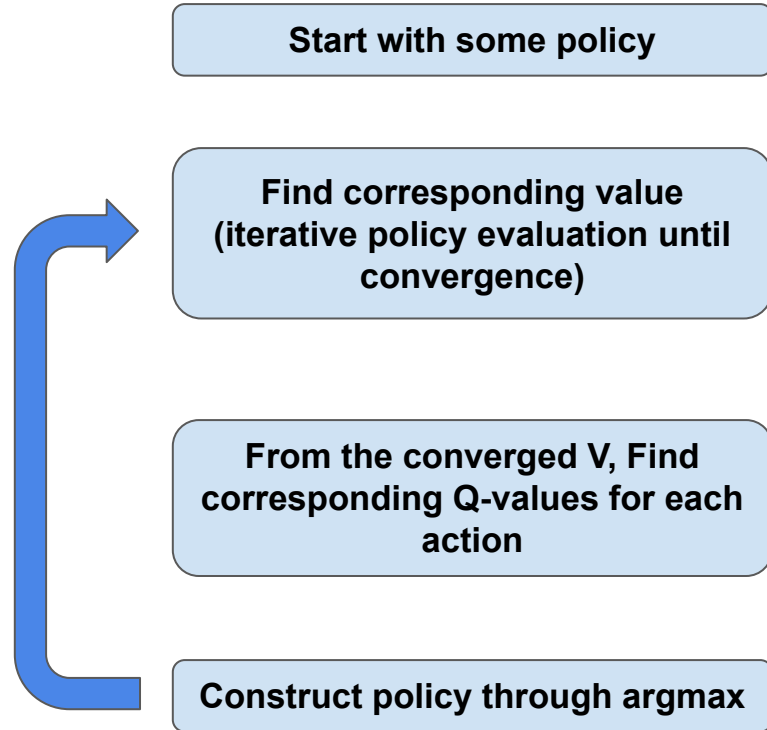
Policy Iteration: Using Q-values



1. $V^\pi(s) = r(s, \pi(s)) + \gamma \cdot \mathbb{E}_{p(s'|s, \pi(s))} [V^\pi(s')]$
2. Set $\pi \leftarrow \pi_{\text{new}}$

$$\pi_{\text{new}} = \begin{cases} 1 & \text{if } a = \arg \max_a Q^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Workflow



Value Iteration

$a_?$
$a_?$
$a_?$
$a_?$
$a_?$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

a_2
a_1
a_3
a_3
a_1

$a_?$
$a_?$
$a_?$
$a_?$
$a_?$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

a_2
a_1
a_3
a_3
a_1


$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

Value Iteration

Start with a random value function $V(s)$

- 
1. Set $Q(s, a) \leftarrow r(s, a) + \gamma \cdot \mathbb{E}_{p(s'|s,a)} \left[V^\pi(s') \right]$
 2. Set $V(s) \leftarrow \max_a Q(s, a)$

Value Iteration Demo

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

0.00 ↖	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↗
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕					0.00 ↕				0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕ R 1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕		0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕ R -1.0		0.00 ↕ R -1.0	0.00 ↕ R -1.0	0.00 ↕	0.00 ↕	0.00 ↕
0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕	0.00 ↕

Fitted Value Iteration

Toy Domains to Reality

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

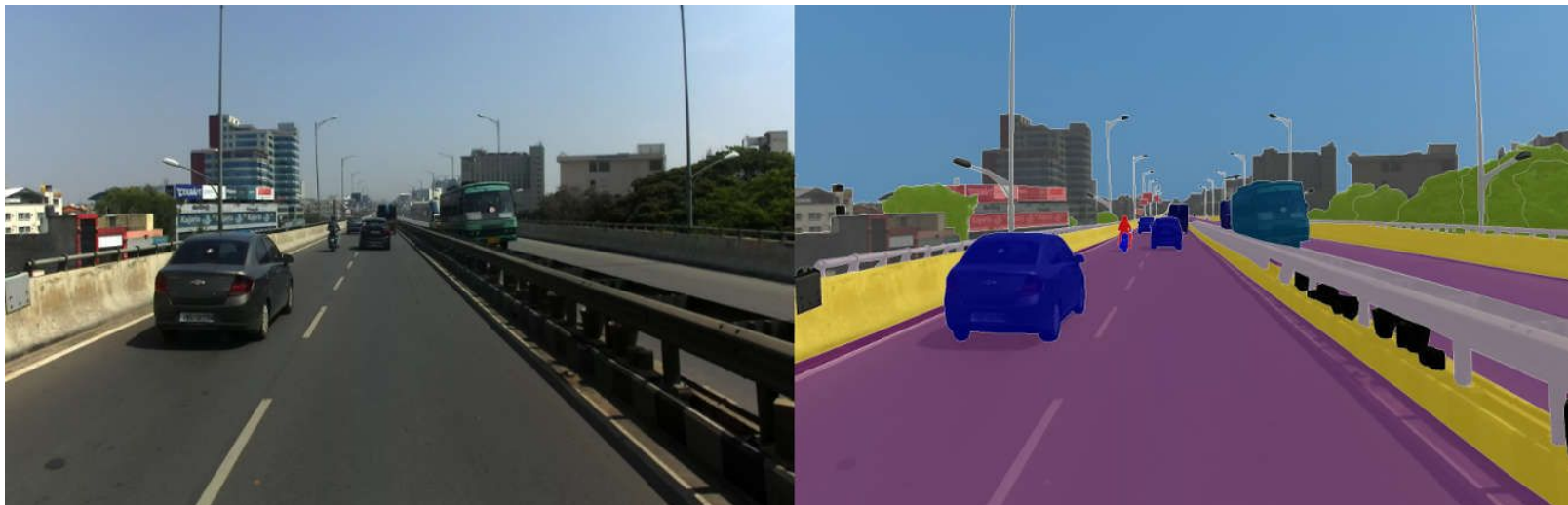
0.00 ↖	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ↘
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	█	█	█	█	0.00 ⬇	█	█	█	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇ R -1.0	█	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	█	0.00 ⬇ R -1.0	0.00 ⬇ R -1.0	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	█	0.00 ⬇ R 1.0	0.00 ⬇ R -1.0	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	█	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇ R -1.0	█	0.00 ⬇ R -1.0	0.00 ⬇ R -1.0	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇
0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇	0.00 ⬇



© Authors of ICRA 2018 Paper 1799 Thu AM Post Q.2

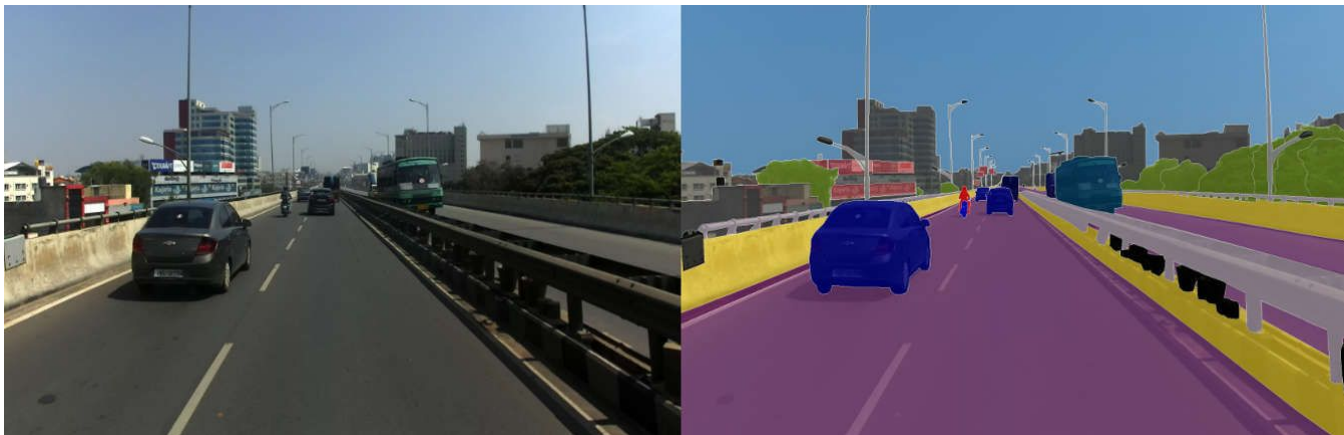


Toy Domains to Reality



How do we represent V ?

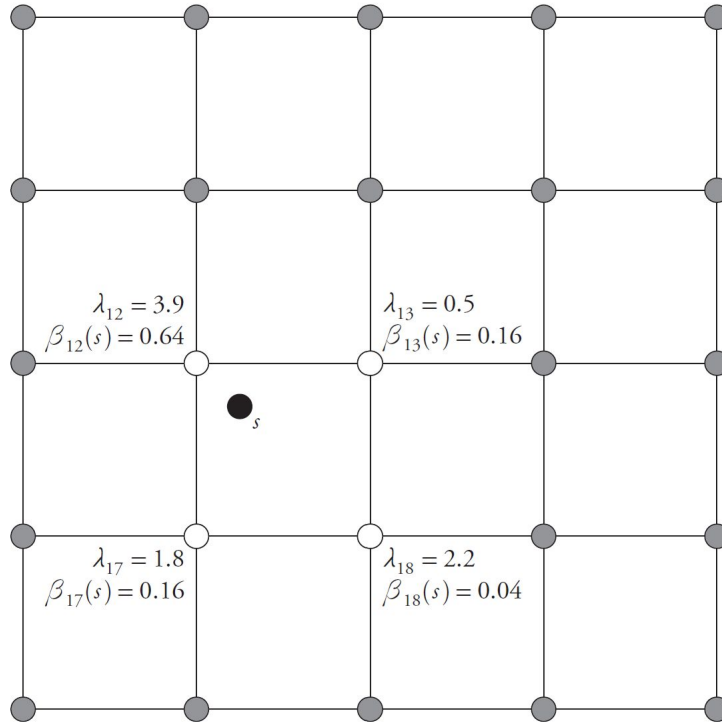
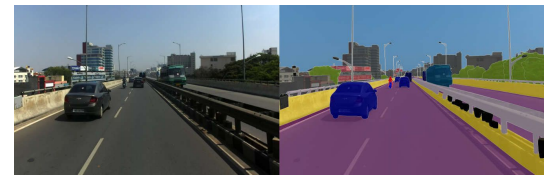
Toy Domains to Reality



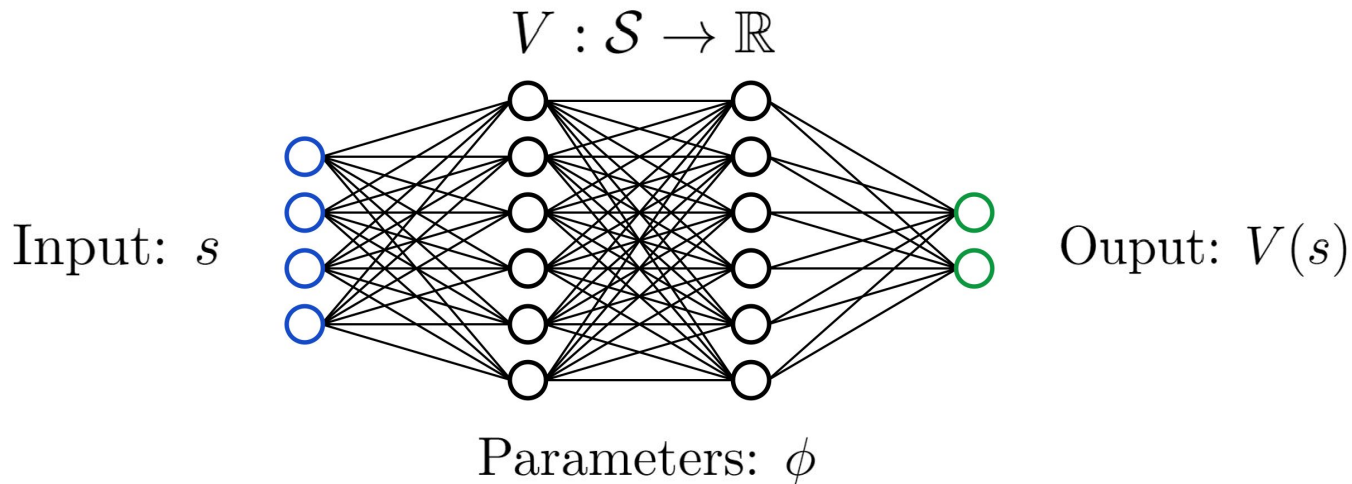
Curse of dimensionality

$$|\mathcal{S}| = (255^3)^{600 \times 600}$$

Approximating the Value Function




Approximating V




What should we train it on?

Value Iteration

Start with a random value function $V(s)$

- 
1. Set $Q(s, a) \leftarrow r(s, a) + \gamma \cdot \mathbb{E}_{p(s'|s,a)} \left[V^\pi(s') \right]$
 2. Set $V(s) \leftarrow \max_a Q(s, a)$

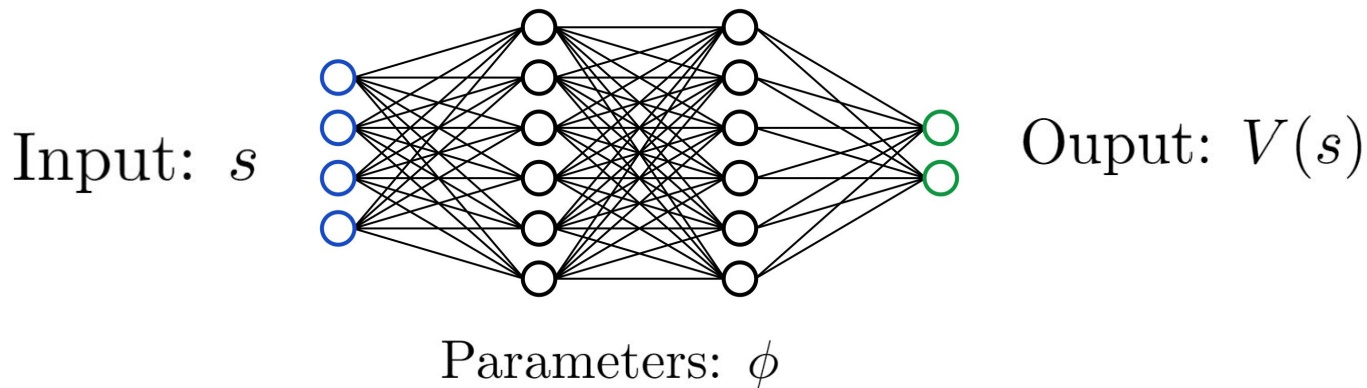
Value Iteration

- 
1. Set $Q(s, a) \leftarrow r(s, a) + \gamma \cdot \mathbb{E}_{p(s'|s,a)} \left[V^\pi(s') \right]$
 2. Set $V(s) \leftarrow \max_a Q(s, a)$

$Q(s_1, a_1)$	$Q(s_1, a_2)$	$Q(s_1, a_3)$
$Q(s_2, a_1)$	$Q(s_2, a_2)$	$Q(s_2, a_3)$
$Q(s_3, a_1)$	$Q(s_3, a_2)$	$Q(s_3, a_3)$
$Q(s_4, a_1)$	$Q(s_4, a_2)$	$Q(s_4, a_3)$
$Q(s_5, a_1)$	$Q(s_5, a_2)$	$Q(s_5, a_3)$

$V(s_1)$
$V(s_2)$
$V(s_3)$
$V(s_4)$
$V(s_5)$

Loss Function



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

How do we instantiate this?

Fitted Value Iteration

1. Set $y_i \leftarrow \max_a \left(r(s_i, a_i) + \gamma \cdot \mathbb{E} \left[V_\phi(s'_i) \right] \right)$

2. Set $\phi \leftarrow \arg \min_\phi \sum_i \frac{1}{2} \|V_\phi(s_i) - y_i\|^2$

Fitted Value Iteration

1. Set $y_i \leftarrow \max_a \left(r(s_i, a_i) + \gamma \cdot \mathbb{E} \left[V_\phi(s'_i) \right] \right)$

2. Set $\phi \leftarrow \arg \min_\phi \sum_i \frac{1}{2} \|V_\phi(s_i) - y_i\|^2$

- We will work with samples
- We have a (finite) sampled set of states
- At each state, we compute the Q values corresp to each action, then take the max over those to create our target y_i
- Compute NN parameters through linear regression to make V close to $\max Q$

Why did we not need samples in PI and VI?

Fitted Value Iteration

$$\text{Dataset: } \left\{ (s_i, a_i, s'_i, r_i) \right\}$$

1. Set $y_i \leftarrow \max_a \left(r(s_i, a_i) + \gamma \cdot \mathbb{E} \left[V_\phi(s'_i) \right] \right)$
2. Set $\phi \leftarrow \arg \min_\phi \sum_i \frac{1}{2} \|V_\phi(s_i) - y_i\|^2$

Announcements

- Pytorch tutorial tomorrow



[Course webpage](#)