

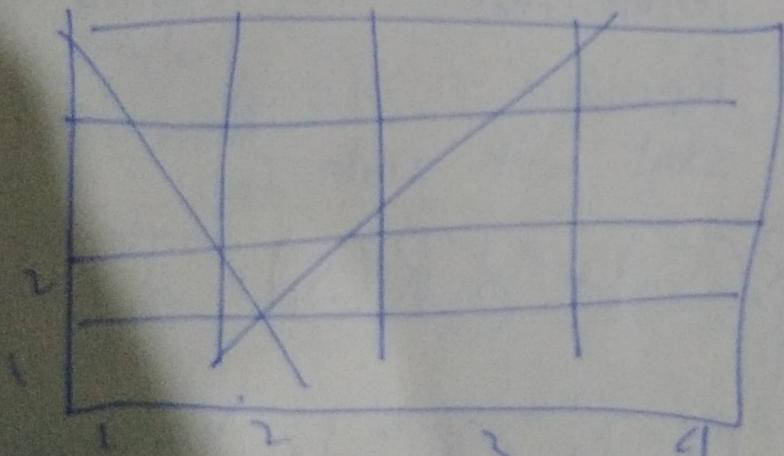
Define V^π, V^*, π^*

✓ For a policy $\pi, V^\pi: S \rightarrow R$

s.t $V^\pi(s)$ is the expected total payoff for starting in state s and executing π .

$$V^\pi = E[R(S_0) + \gamma R(S_1) + \dots | \pi, S_0 = s]$$

(\Rightarrow Value function of a policy π .)



π (A specific policy for our MDP)

3	→	→	→	+1
2	↓	B	→	-1
1	→	→	↑	↑

↳ Not a great policy.

V^π (for the above policy π):

.52	.73	.77	.11
-.90	.22	-.82	-.1
-.088	-.087	-.085	-.106

Bellman's equation:

$$V^\pi(s) = R(s_0) + \gamma \sum_s P_{s\pi(s)}(s') V^\pi(s')$$

\downarrow
 Immediate
 Reward \downarrow
 Expected
 Future
 Reward

Let say you start with state s_0 :

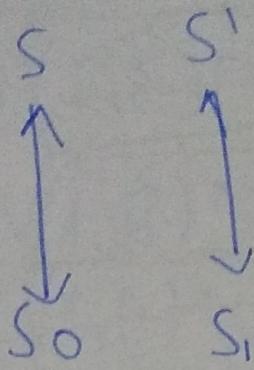
$$V^\pi(s) = E[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

$$V^\pi(s) = E[R(s_0) + \gamma (R(s_1) + \gamma R(s_2) + \dots)] | \bar{\pi}, s_0 = s_1$$

$$V^\pi(s) = R(s_0) + \gamma \sum_s P_{s\pi(s)} V^\pi(s') \quad (\text{Bellman's equation})$$

(*)

Above equation mapping to our example.



$$\textcircled{1} \quad V^\pi(s) = E[R(s) + \gamma V^\pi(s')]$$

$s' \sim P_{s\pi(s)}$ (s' distribution).

In state s , take action
 $a = \pi(s)$.

→ Now again taking $\textcircled{1}$; and taking into s' distribution we have Bellman's equation

$$V^\pi(s) = E(R(s)) + E(\gamma V^\pi(s'))$$

$$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{s\pi(s)} V^\pi(s')$$

→ Given π , get a linear system of equations in terms of $V^\pi(s)$ (By Bellman's equation).

Example:



Let say you are in the position (3, 1).

191

$$\circ \pi(3,1) = \uparrow$$

$$\boxed{V^\pi((3,1))} = R(3,1) + \gamma \left(0.8 \boxed{V^\pi(3,2)} + 0.1 \boxed{V^\pi(2,1)} + 0.1 \boxed{V^\pi(4,1)} \right)$$

→ Are unknown variables and if you have 11 states then you will have 11 equations (Value function for each state) and 11 unknowns which can be solved through system of 11 linear equations.

→ Given any policy π , you can have 11 unknown variable vector and solve it.

$$\begin{bmatrix} V^\pi((1,2)) \\ V^\pi((1,2)) \\ \vdots \\ V^\pi((4,3)) \end{bmatrix} \rightarrow 11 \text{ unknowns} - \mathbf{y} \in \mathbb{R}^{11}$$

→ Now lets define V^*

→

(92)

V^* is the optimal value function:

$$V^*(s) = \max_{\pi} V^\pi(s)$$

↳ Over all the combinatorially large number of possible policies for this MDP, let's take the value of the next possible value for that state.

→ Bellman's equations for $V^*(s)$ (optimal value):

$$V^*(s) = R(s) + \max_a \sum_{s'} P_{s'a}^{(s)} V^*(s')$$

→ Written in neat form:

$$V^*(s) = R(s) + \max_a \sum_{s'} P_{s'a}(s') V(s')$$

↓
Immediate
Reward

Action ↓ which
maximizes
the future expected
reward.

→ Now based on $V^*(s)$ we will
compute:

$$\boxed{\pi^*(s)} = \operatorname{argmax}_{\pi} \sum_{s'} P_{sa} V(s')$$

↓
Optimal
Policy

↳ What is the policy that maximizes the future expected value/reward.

→ Practice with confusion notation:

$$V^*(s) = V^{\pi^*}(s) \geq V^{\pi}(s)$$

for every $\boxed{\pi, s}$

→ Strategy to find optimal Policy

- 1) Find V^* argmax equation to
- 2) Use find π^*

① First Strategy: Value Iteration

→ Algorithm:

Initialize $V(s) := 0$ for every state s .

For every s update:

$$V(s) := R(s) + \max_a \sum_{s'} P_{sa}(s') V(s')$$

→ For example in our map, you will
create a 11-dimensional vector ea for in
order to represent value function for
each state:

$$\begin{bmatrix} V((1,1)) \\ V((1,2)) \\ \vdots \\ V((4,3)) \end{bmatrix} \rightarrow \in \mathbb{R}^{11}$$

(Initialize to zero)

→ Two ways:

a) Synchronous update:

$$\rightarrow \text{compute } \max_a \sum_{s'} P_{sa}(s') V(s')$$

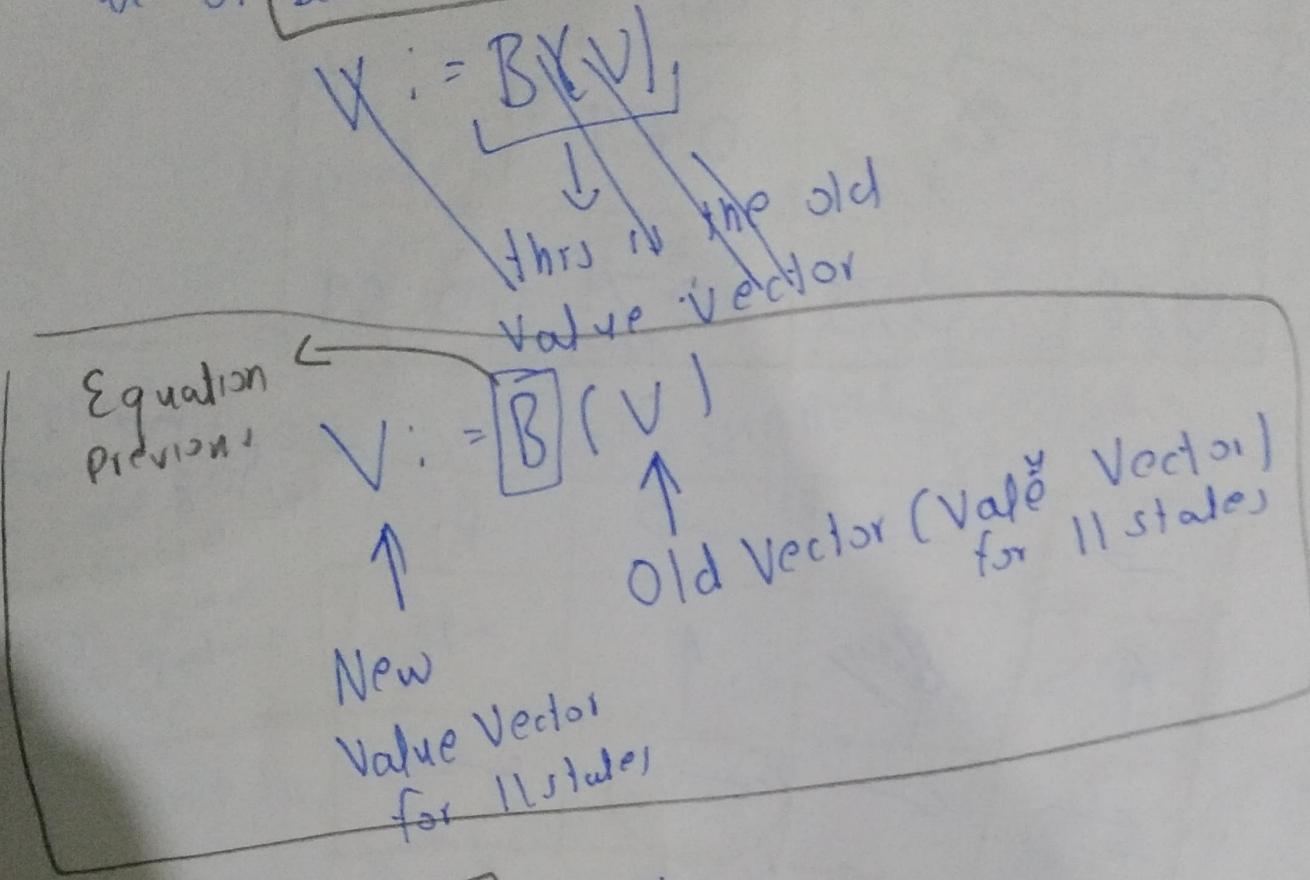
Simultaneously for 11 states and
update all 11 states

b) Asynchronous update

→ Overwriting each value
one at a time.

(195)

- * In Synchronous update we make use of Bellman's backup operator:-



Value Iteration (cause)

V to converge to V^*
 (there is a formal proof).
 after a specific number of
 iteration. (100 or few iteration).

- * If you run Value iteration on our MDP you end up with:-

3	0.86	0.90	0.93	+1
2	0.81	0.81	0.69	-1
1	0.78	0.75	0.71	0.99

↗ Optimal policy.

3	→	→	→	+1
2	↑	0.81	↑	-1
1	←	←	←	9

Let say you are in state (3,1) now
 do you compute optimal policy (using π^*)
 equation) (Future Reward for each action a in state (3,1))

$$\textcircled{1} \quad W = \sum_{s'} P_{sa} V(s') = 0.8 \times 0.75 + 0.1 \times 0.69 + 0.1 \times \cancel{0.49} = 0.740$$

$$\textcircled{2} \quad N = \frac{0.8 \times 0.69}{0.676} + 0.1 \times \cancel{0.71} + 0.1 \times \cancel{0.49} = 0.53$$

$$\textcircled{3} \quad E\$ = 0.8 \times \frac{0.75}{0.69} + 0.1 \times \cancel{0.75} + 0.1 \times \cancel{0.49} = 0.53$$

↳ Based on this optimal policy it's better to go left (higher value)

④ Second Strategy:- Policy Iteration:

→ Algorithm:

Initialize π randomly

Repeat {

Set $V_i = V\pi$ (i.e solve Bellman equation to get V^π)

Set $\pi(s) = \operatorname{argmax}_a \sum_{s'} P_{sa}(s')V(s')$

}

As previously shown that for fixed π solve for V^π through system of linear equations

→ In value iteration our focus is on V^* (coming up with policy at end)

→ In policy iteration our focus is on policy π (coming up with new policy at each iteration)

→ Pros & Cons: (VI vs PI)

→ Policy iteration requires solving system of linear equations which become quite expensive in large state space so in those cases Value iteration is more feasible.

198

→ Value iteration only come closer to V^* (but is never exact) whereas policy iteration solves for exact V^*

→ What if don't know P_{sa} ?

→ In practical cases you don't know the state transition probabilities

→ In many MDP you need to estimate $P_{sa}(s')$ from data

$\rightarrow P_{sa}(s') = \frac{H}{H} \text{ of times agent lands in state } s' \text{ given state } s \text{ & action } a$

H of time agent is in state s and execute action a .

$= \frac{H}{H} \text{ of time took action } a \text{ in state } s \text{ and got to } s'$

$\frac{H}{H} \text{ of times at took action } a \text{ in states }$

→ or $\frac{1}{|S|}$ if above ab is 0
 ↪ total number of states -

→ Putting it together

Repeat {

ε-greedy

0.9 chance w.r.t π
0.1 chance randomly

(See on
next page)
!!!

- Take action w.r.t π to get experience in MDP

• Update estimates of P_{sa} based on observation of data

• Solve Bellman equation using value iter

to get \check{V}

• Update $\pi(s) = \operatorname{argmax}_a \sum_{s'} P_{sa}(s') V(s')$

⋮

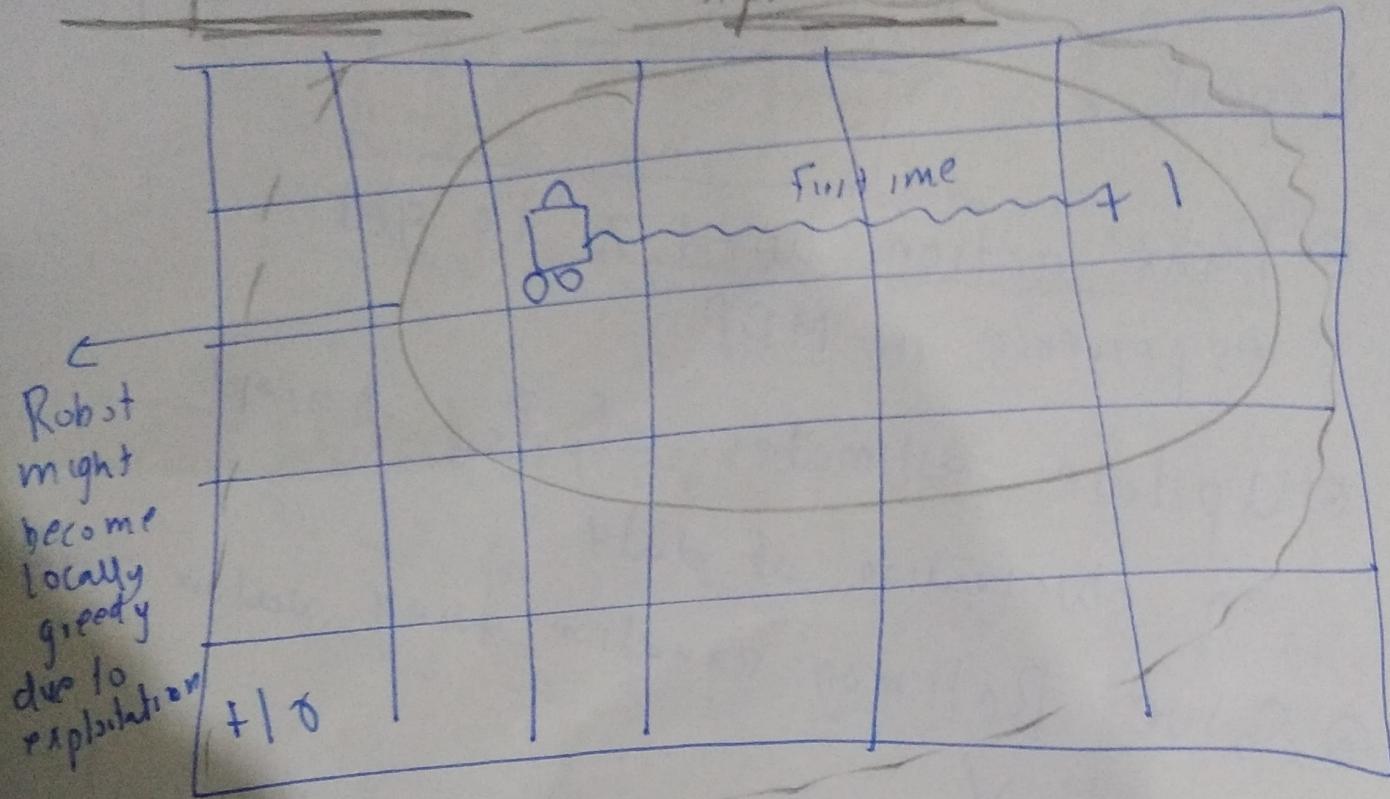
→ Sometimes reward function maybe unknown (May need to design it)

↳ Depends on action sometime

↳ Sometime random function of environment).

200

④ Exploration vs Exploitation



→ Exploration would enable robot to find +10.

→ How greedy should be the robot to maximize reward (should it explore MDP or only exploit the reward).

