# REINFORCEMENT LEARNING

## FUNDAMENTALS
## +
## APPLICATIONS

# WEEK 3

---

## OPTIMALITY

+

## DYNAMIC PROGRAMMING

$$v(s) = \mathbb{E}[R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s]$$

$$q(s,a) = \mathbb{E}[R_{t+1} + \gamma q(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

# BELLMAN EXPECTATION EQUATIONS



$$v_\pi(s) = \sum \pi(a|s) \cdot q_\pi(s,a)$$

$$q_\pi(s,a) = \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$v_\pi(s) = \sum \pi(a|s) \cdot \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$q_\pi(s,a) = \sum p(s'|s,a) \cdot [r_{s',a} + \gamma \sum \pi(a'|s') \cdot q_\pi(s',a')]$$

# OPTIMAL VALUE FUNCTIONS

$$v_*(s) \equiv \max_\pi v_\pi(s)$$

$$q_*(s,a) \equiv \max_\pi q_\pi(s,a)$$

$$\pi' \geq \pi \text{ iff } v_{\pi'}(s) \geq v_{\pi}(s) \; \forall \; s$$

For any Markov Decision Process,

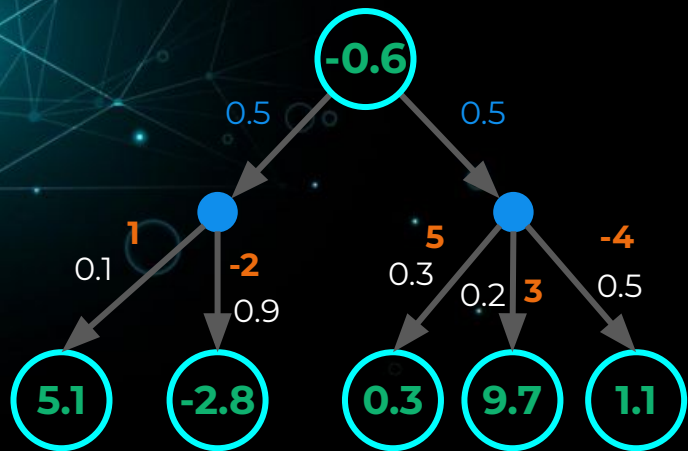- There exists an optimal policy $\pi_*$ that is better than or equal to all other policies, $\pi_* \geq \pi$.
- There is always a deterministic optimal policy $\pi_*(a|s) = \text{argmax} q_*(s,a)$
- There can be more than one optimal policy, but all optimal policies share the same optimal value functions

$$v_{\pi^*}(s) = v_*(s)$$
$$q_{\pi^*}(s,a) = q_*(s,a)$$

# STATE-VALUE FUNCTION

Bellman Expectation Equation (Uniform Random Policy)

$$v_\pi(s) = \sum_a \pi(a|s) \cdot \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$



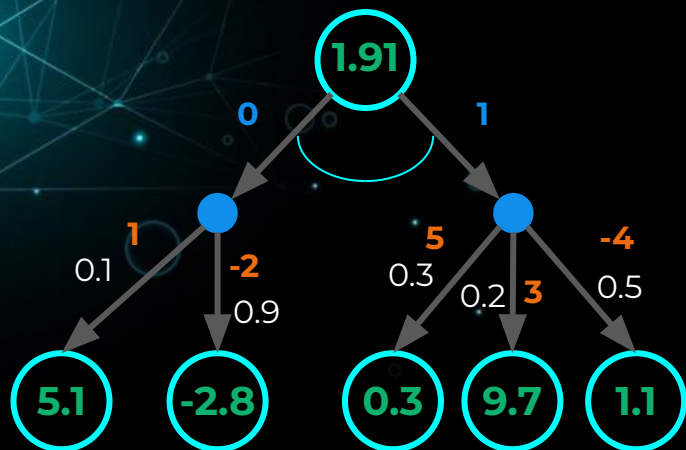$( 0.5 \cdot 0.1 \cdot [1 + 0.7 \cdot 5.1] )$
$+$
$( 0.5 \cdot 0.9 \cdot [-2 + 0.7 \cdot -2.8] )$
$+$
$( 0.5 \cdot 0.3 \cdot [5 + 0.7 \cdot 0.3] )$
$+$
$( 0.5 \cdot 0.2 \cdot [3 + 0.7 \cdot 9.7] )$
$+$
$( 0.5 \cdot 0.5 \cdot [-4 + 0.7 \cdot 1.1] )$
$=$
$-0.6$

# STATE-VALUE FUNCTION

## Bellman Optimality Equation

$$q_\pi(s,a) = \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

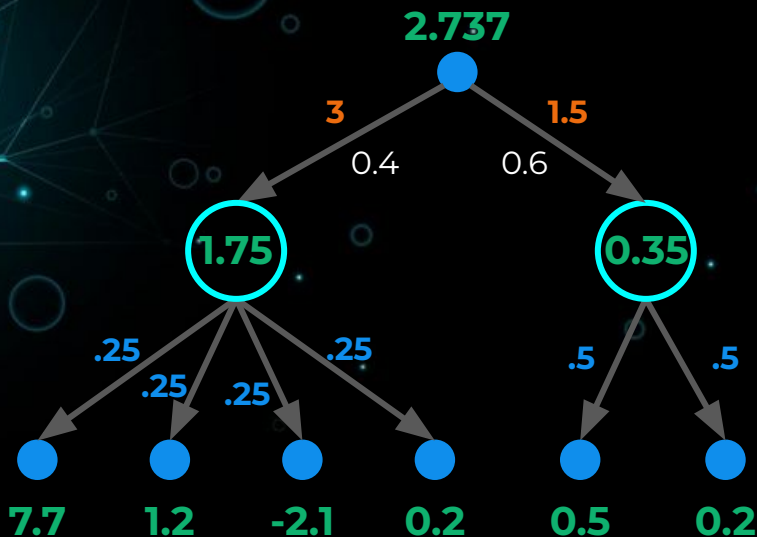$$v_*(s) = \max_a \sum p(s',r|s,a) \cdot [r + \gamma v_*(s')]$$



$$( 1 \cdot 0.3 \cdot [5 + 0.7 \cdot 0.3] )$$
$$+$$
$$( 1 \cdot 0.2 \cdot [3 + 0.7 \cdot 9.7] )$$
$$+$$
$$( 1 \cdot 0.5 \cdot [-4 + 0.7 \cdot 1.1] )$$
$$=$$
$$1.906$$

# ACTION-VALUE FUNCTION

Bellman Expectation Equation (Uniform Random Policy)

$$q_\pi(s,a) = \sum p(s'|s,a) \cdot [r_{a,s'} + \gamma \sum \pi(a'|s') \cdot q_\pi(s',a')]$$

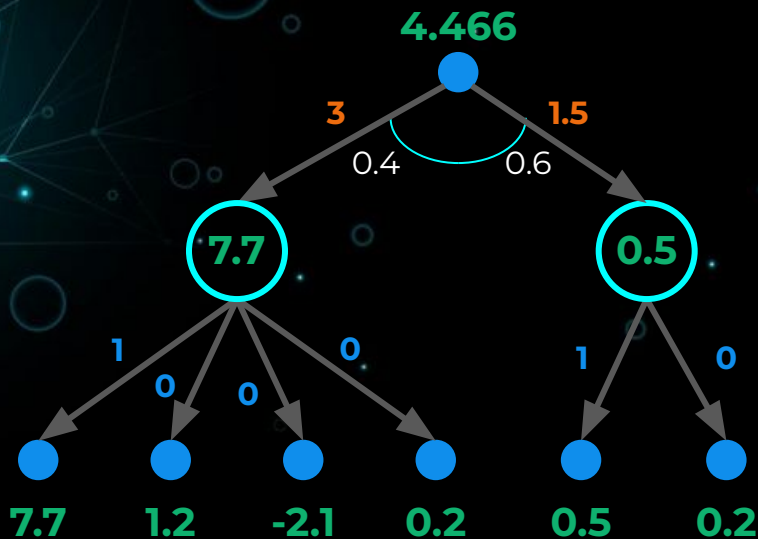2.737

3     1.5

0.4     0.6

1.75     0.35

.25    .25    .25    .25     .5    .5

7.7    1.2    -2.1    0.2    0.5    0.2

( 0.4 · [3 + 0.7 · 1.75] )
+
( 0.6 · [1.5 + 0.7 · 0.35] )
=
2.737

# ACTION-VALUE FUNCTION

## Bellman Optimality Equation

$$q_*(s,a) = \sum p(s'|s,a) \cdot [r_{a,s'} + \gamma \max_{a'} q_*(s',a')]$$



$(\ 0.4 \cdot [3 + 0.7 \cdot 7.7]\ )$
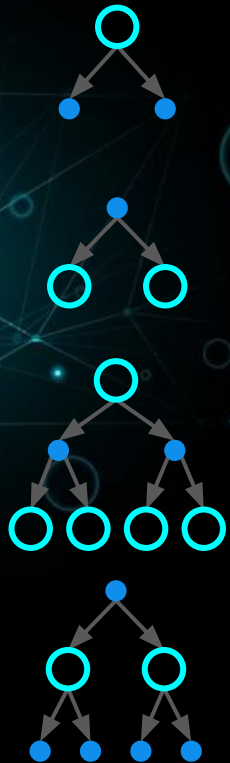
$+$

$(\ 0.6 \cdot [1.5 + 0.7 \cdot 0.5]\ )$

$=$

**4.466**

# OPTIMALITY

## Solving the Bellman Optimality Equation

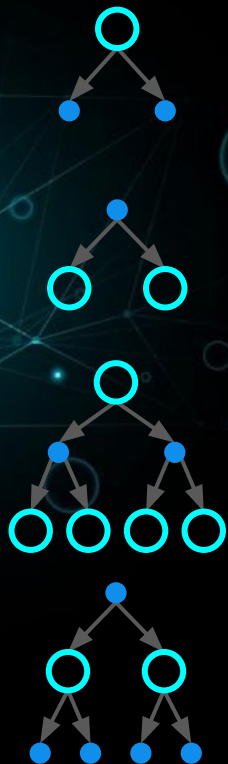# BELLMAN EXPECTATION EQUATIONS

$$v_\pi(s) = \sum \pi(a|s) \cdot q_\pi(s,a)$$

$$q_\pi(s,a) = \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$v_\pi(s) = \sum \pi(a|s) \cdot \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$q_\pi(s,a) = \sum p(s'|s,a) \cdot [r_{s',a} + \gamma \sum \pi(a'|s') \cdot q_\pi(s',a')]$$

# DYNAMIC PROGRAMMING

$$v_\pi(s) = \mathbf{max}\, q_\pi(s,a)$$

$$q_\pi(s,a) = \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$v_\pi(s) = \mathbf{max} \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$$

$$q_\pi(s,a) = \sum p(s'|s,a) \cdot [r_{s',a} + \gamma\,\mathbf{max}\, q_\pi(s',a')]$$

# BELLMAN OPTIMALITY EQUATIONS

$$v_*(s) = \max \sum p(s',r|s,a) \cdot [r + \gamma v_*(s')]$$

$$q_*(s,a) = \sum p(s'|s,a) \cdot [r_{a,s'} + \gamma \max q_*(s',a')]$$

# POLICY EVALUATION

## (PREDICTION)

MDP + $\pi$ go in,
$v_\pi$ comes out

# POLICY EVALUATION

Prediction

If I follow this (FIXED) policy in this environment, how much reward can I expect starting in each state?

# POLICY EVALUATION

Prediction

$$V_0 \rightarrow V_1 \rightarrow V_2 \rightarrow V_3 \rightarrow \ldots \rightarrow V_\pi$$

# POLICY EVALUATION

## Prediction

Environment
Dynamics
+
Way of Behaving

→

How much reward
will I experience if I
behave that way, in
that environment?

→

How good each state
and action are if you
behave that way in
that environment.

# POLICY EVALUATION

The Algorithm : State-Value Function

Input: $\pi$, the MDP ($\mathcal{S}$, $\mathcal{P}$, $\mathcal{R}$, $\gamma$, $\mathcal{A}$), small positive number $\theta$
Initialize: An array $V(s) = 0 \; \forall \; s$, $\Delta = \theta+1$

---

while $\Delta > \theta$:
    $\Delta = 0$
    For $s$ in $\mathcal{S}$:
        $v \leftarrow V(s)$
        $V(s) \leftarrow \sum \pi(a|s) \cdot \sum p(s',r|s,a) \cdot [r + \gamma v_\pi(s')]$
        $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
return $V$

$$v_{k+1} = \mathcal{R} + \gamma \mathcal{P}_\pi v_k \quad \text{!!!}$$

# POLICY EVALUATION

## The Algorithm : Action-Value Function

Input: $\pi$, the MDP ($\mathcal{S}$, $\mathcal{P}$, $\mathcal{R}$, $\gamma$, $\mathcal{A}$), small positive number $\theta$

Initialize: An array $Q(s,a)$ = 0 $\forall$ $s$, $\Delta$ = $\theta$+1

while $\Delta > \theta$:
$\quad \Delta = 0$
$\quad$ For $s$ in $\mathcal{S}$ and $a$ in $\mathcal{A}$:
$\qquad q \leftarrow Q(s,a)$
$\qquad Q(s,a) \leftarrow \sum p(s',r|s,a) \cdot [r + \gamma \sum \pi(a'|s') \cdot Q(s',a')]$
$\qquad \Delta \leftarrow \max(\Delta, |q - Q(s,a)|)$
return $Q$

# POLICY IMPROVEMENT

## (CONTROL)

MDP + π + V *go in,*
π' ≥ π *comes out*

# POLICY IMPROVEMENT

Intuition

- We want our agent to learn and improve
- We need a better policy
- We need a way to turn our policy into a better policy

# POLICY IMPROVEMENT

Control

| Environment Dynamics + Value Function | → | WHAT HAPPENS IF WE'RE AS GREEDY AS POSSIBLE?!?!?! | → | How good each state and action *can possibly be* given the environment and its value function. |

# POLICY IMPROVEMENT

Control

MDP
$\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma, \mathcal{A} \rangle$

→

Policy
Improvement
Algorithm

→

$v_*$

$q_*$

# POLICY IMPROVEMENT

The Algorithm : State-Value Function

POLICY_STABLE = TRUE

---

for $s$ in $\mathcal{S}$:
        old_action = $\pi(s)$
        $\pi(s) \leftarrow \text{argmax} \sum p(s',r|s,a) \cdot [r + \gamma V(s')]$
        if old_action != $\pi(s)$:
                POLICY_STABLE = FALSE
if POLICY_STABLE:
        return $V$, $\pi$
else:
        policy_evaluation($\pi$)

# POLICY IMPROVEMENT

The Algorithm : Action-Value Function

POLICY_STABLE = TRUE

---

```
for s in 𝒮:                    a
    old_action = π(s)
    π(s) ← argmax Q(s,a)
    if old_action != π(s):
        POLICY_STABLE = FALSE
if POLICY_STABLE:
    Return Q, π
else:
    policy_evaluation(π)
```

# POLICY ITERATION
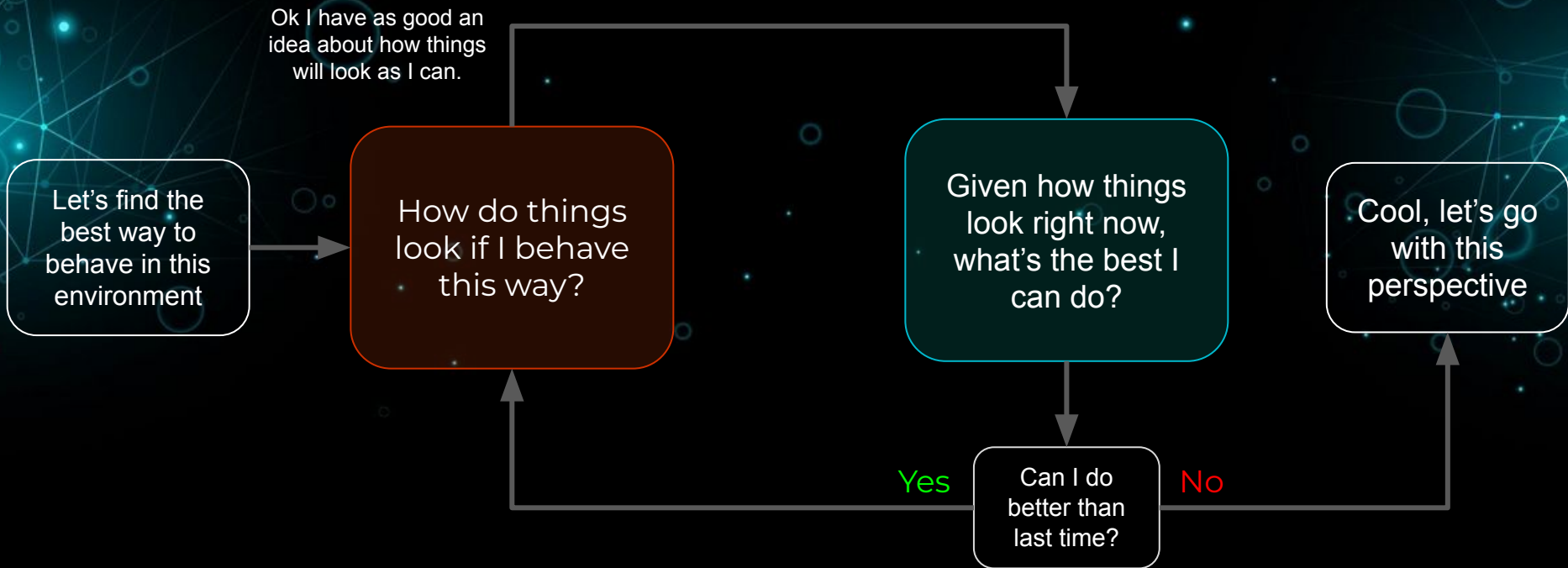
## (PREDICTION + CONTROL)

MDP + π go in,
π* and v* come out

# POLICY ITERATION
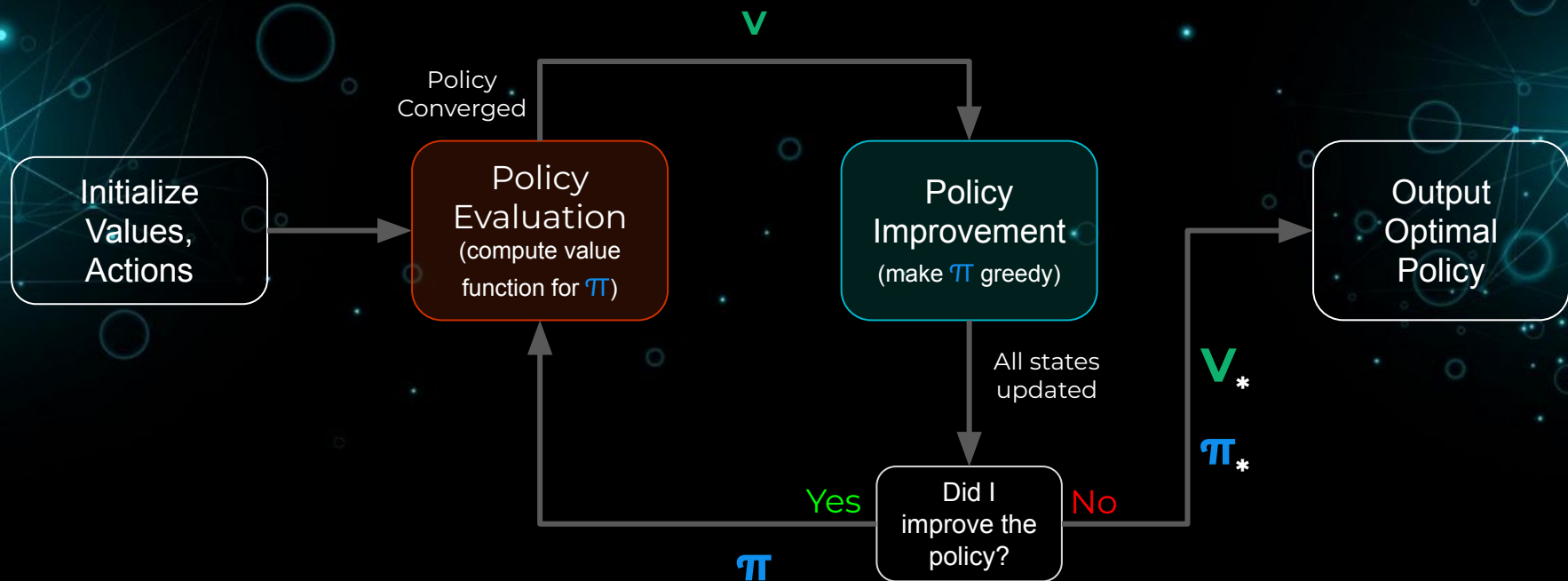
## Generalized Policy Iteration

Ok I have as good an idea about how things will look as I can.

Let's find the best way to behave in this environment

How do things look if I behave this way?

Given how things look right now, what's the best I can do?

Cool, let's go with this perspective

Yes   Can I do better than last time?   No

# POLICY ITERATION

## The Algorithm

# VALUE ITERATION

Policy Iteration with 1 (one) Evaluation Sweep

MDP *goes in,*
$\pi_*$ *comes out*

*MAGIC !*

# VALUE ITERATION

## The Algorithm

Input: The MDP ($\mathcal{S}$, $\mathcal{P}$, $\mathcal{R}$, $\gamma$, $\mathcal{A}$), small positive number $\theta$
Initialize: An array $\mathbf{V}_a$ = 0 $\forall$ $s$, $\Delta = \theta + 1$

---

while $\Delta > \theta$:
    $\Delta = 0$
    For $s$ in $\mathcal{S}$:
        $\mathbf{v} \leftarrow \mathbf{V}(s)$
        $\mathbf{V}(s) \leftarrow \max \sum p(s',r|s,a) \cdot [r + \gamma \vee(s')]$
        $\Delta \leftarrow \max(\Delta, |\mathbf{v} - \mathbf{V}(s)|)$
set $\pi(s) \leftarrow \text{argmax} \sum p(s',r|s,a) \cdot [r + \gamma \vee(s')]$

Return $\pi \approx \pi_*$