# REINFORCEMENT LEARNING

Replace Y with reward (Y is the supervision)

Maximize reward over time

In supervised learning i.i.d., goal was to do well in that example, however in R.L. the goal is to maximize reward over time

The concept of time makes R.L. special

Markov Decision Process (MDP)
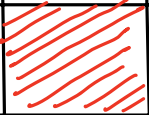
Tuple $(S, A, \{P_{sa}\}, V, R)$

$S$ - set of states

$A$ - set of actions

$P_{sa} \to$ Transition probabilities    $S_{next} \sim P_{s_{current} \cdot a_{current}}$

$V \in (0,1)$ - Discount factor

$R: S \times A \to \mathbb{R}$
    or
$R: S \to \mathbb{R}$  $\Big\} \to$ Reward Function

| | | | |
|---|---|---|---|
| 3 | -0.02 | -0.02 | -0.02 | +1 |
| 2 | | | -0.02 | -1 |
| 1 | | | | |
| | 1 | 2 | 3 | 4 |

$$S = \{(1,1),(1,2)\ldots\}$$

$$|S| = 11$$

$$A = \{N, S, W, E\}$$

$$P_{sa} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0.1 \rightarrow (1,4) \\ 0.1 \rightarrow (1,2) \\ 0.8 \rightarrow (2,3) \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

for $s = (1,3)$
$a = N$

(Assumption that 80% of times when $a = N$, it goes N, 10% west, 10% east)

$$R(S) \rightarrow \mathbb{R} \quad \text{or} \quad R(s,a) \rightarrow \mathbb{R}$$

↓ reward depends on state

↓ sometimes even if state is the same, action taken

$R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2) + \gamma^3 R(S_3) \longrightarrow$ maximize this

$$S_0 \xrightarrow{a_0} S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} S_3 \xrightarrow{a_3} \ldots \quad \text{Trial /}$$
Episode /
Trajectory

$$S_1 \sim P_{S_0 a_0} \qquad S_2 \sim P_{S_1 a_1}$$

[ Why is $P_{sa}$ a probability vector?

→ Suppose robot is placed in an environment, you have told it to go 100 cm, it may go 95/105 cm. Hence, there is stochasticity. Hence to account for it, we have probability]

$\gamma \rightarrow$ Discount factor. It incentivizes model to earn large positive rewards sooner, as it discounts reward. Another way to think, we push negative reward at the end so that it is highly discounted.

[In Finance, $\gamma$ is the interest rate, we want profits sooner & loss later]

Policy $\Pi: S \to A$

We want to learn a policy that maximizes value

Value: $V^{\pi}: S \to \mathbb{R}$

$$V^{\pi}(s) = \mathbb{E}\left[R(S_0) + \gamma R(S_1) + \gamma^2 R(S_2)\ldots \mid S_0 = S, \pi\right]$$

$\uparrow$
starting state

Rewards can be random because states are random due to stochastic nature of transition probabilities

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

( $S = S_0$ is not random )

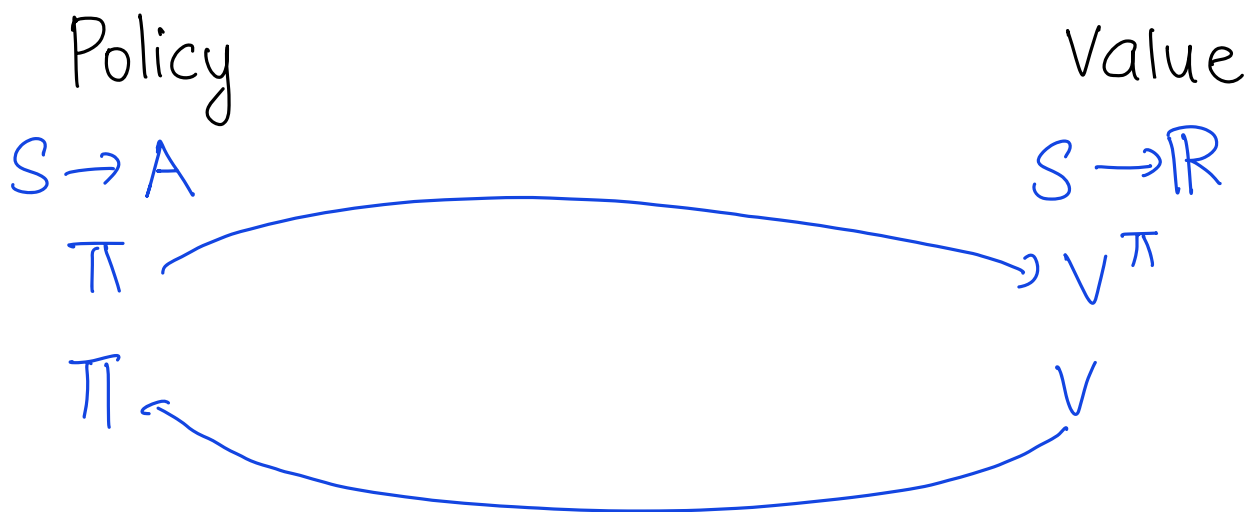$$V^{\pi}(S) = R(S) + \gamma \mathbb{E}(R(S_1) + \gamma R(S_2) \ldots)$$

$$V^{\pi}(S) = R(S) + \gamma \mathbb{E}(V^{\pi}(S_1))$$

$$= R(S) + \gamma \underbrace{\sum_{\underset{\text{future}}{S' \in S}} P_{S\pi(S)}(S') V^{\pi}(S')}$$

future state

Bellman Equation

If our goal was to maximize immediate reward, back to supervised learning

R.L. → focus on value. We look at long term reward

Policy and value duals of each other

Policy

$S \to A$

$\pi$

$\pi$

Value

$S \to \mathbb{R}$

$V^{\pi}$

$V$

the action is optimized to reach next state with highest value..

$$\pi \rightarrow V^\pi$$

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

$$V^\pi = \begin{bmatrix} - \\ - \\ - \\ \vdots \\ \vdots \end{bmatrix} \begin{array}{l} s_1 \quad V^\pi(s_1) \\ s_2 \quad V^\pi(s_2) \\ s_3 \quad \vdots \\ \quad \vdots \end{array}$$

$$V^\pi(s_1) = R(s_1) + \nu \sum_{s' \in S} P_{s\pi(s_1)}(s') V^\pi(s')$$

$$V^\pi(s_2) = R(s_2) + \nu \sum_{s' \in S} P_{s\pi(s_2)}(s') V^\pi(s')$$

$$\vdots$$

$$V^\pi(s_{|S|}) = R(s_{|S|}) + \nu \sum_{s' \in S} P_{s\pi(s_{|S|})}(s') V^\pi(s')$$

$$P^{\pi} = \begin{bmatrix} \underline{\quad\quad P_{S_1 \pi(S_i)} \quad\quad} \\ \vdots \\ \underline{\quad P_{S_i \pi(S_i)} \quad\quad} \end{bmatrix} \Bigg\} |S|$$

$$\underbrace{|S|}$$

$$\underbrace{V^{\pi}}_{|S|} = \underbrace{R}_{|S|} + \upsilon \underbrace{P^{\pi}}_{|S| \times |S|} \underbrace{V^{\pi}}_{|S|}$$

$$V^{\pi} = [I - \upsilon P^{\pi}]^{-1} R$$

# Optimal Value Function

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

$$V^*(s) = R(s) + \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

$$\pi^*(s) = \arg \max_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$$

# VALUE ITERATION

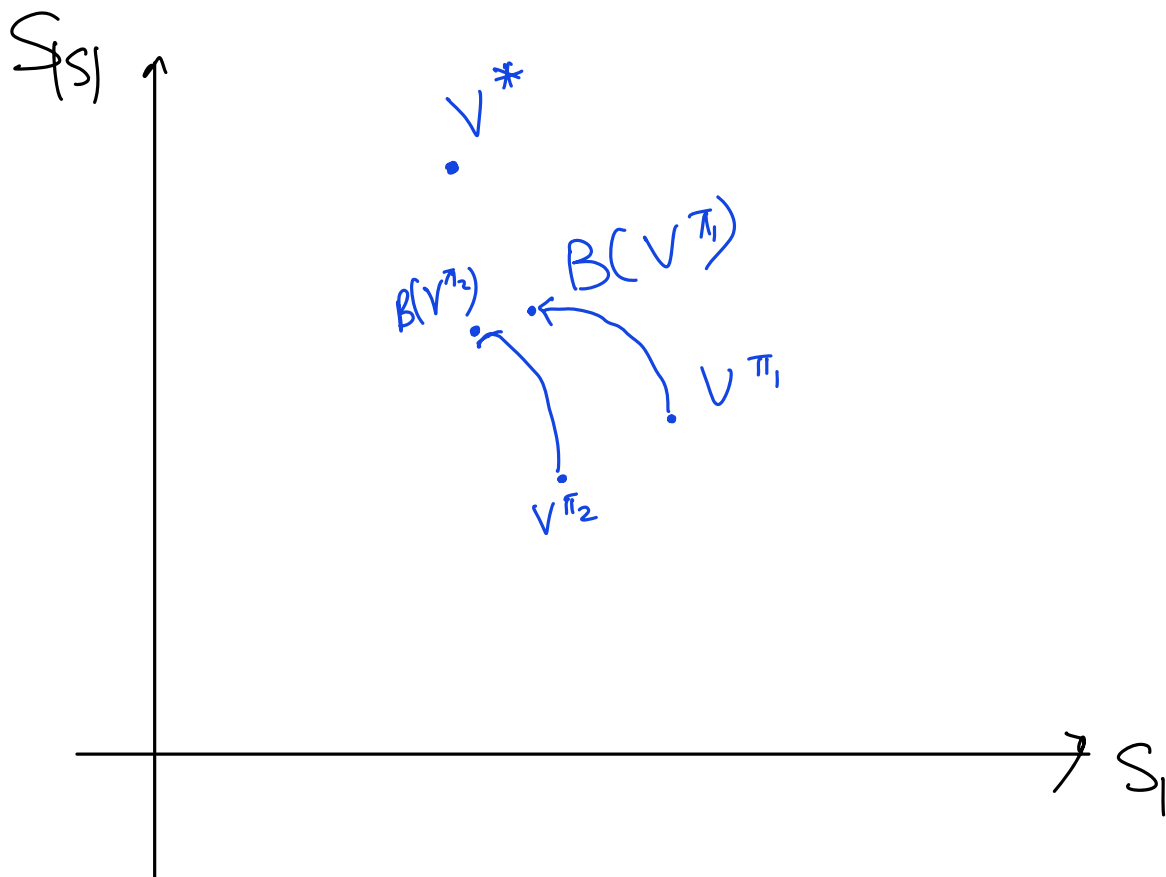Algorithm

1. For each state $s$, initialize $V(s) := 0$

2. Repeat until convergence

   For every state, update

   $$V(s) := R(s) + \max_{a \in A} \gamma \sum_{s'} P_{sa}(s') V(s')$$

   Bellman Backup Operator

$O(|S|^2 |A|)$



$V^\pi(S_1)$ is the projection of $S_1$

Bellman operator is contraction mapping.
Suppose you have 2 points in space.
You apply Bellman. The result will be

closer than the input. They all converge to a single point $\rightarrow$ fixed point. $V^*$

# POLICY ITERATION

1. Initialize $\pi$ randomly
2. Repeat until convergence $\leftarrow O(|S|^3)$
   (a) Set $V := V^\pi$
   (b) For each state $s$, set
   $$\pi(s) := \arg\max_{a \in A} \sum_{s'} P_{sa}(s') V(s')$$
   $\underbrace{\phantom{xxxxxx}} O(|A||S|)$

Policy: you will surely converge
Value: you get closer but may not reach final.

1. Policy $\pi(s) = a$

2. $V^{\pi}(s) = E[R(s_0) + \gamma R(s_1) \cdots | s_0 = s, \pi)$

$$= R(s) + \gamma E_{s' \sim P_{sa}}[V^{\pi}(s')]$$

$\gamma < 1$ makes $V^{\pi}(s)$ bounded
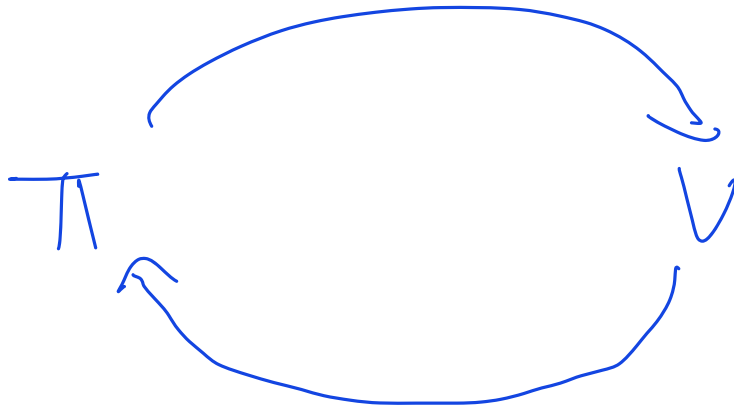
Optimal Value

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

$$= R(s) + \max_{a} \underbrace{E_{s' \sim P_{sa}}[V^*(s')]}_{\text{highest expected reward}}$$

Optimal Policy

$$\pi^*(s) = \arg\max_a \mathbb{E}_{s' \sim P_{sa}}[V^*(s')]$$

$$V^\pi = (I - P^\pi)^{-1} R \quad \text{(Policy evaluation)}$$

①

$$\pi(s) = \arg\max_a \mathbb{E}_{s \sim P_{sa}}\left[V(s')\right]$$
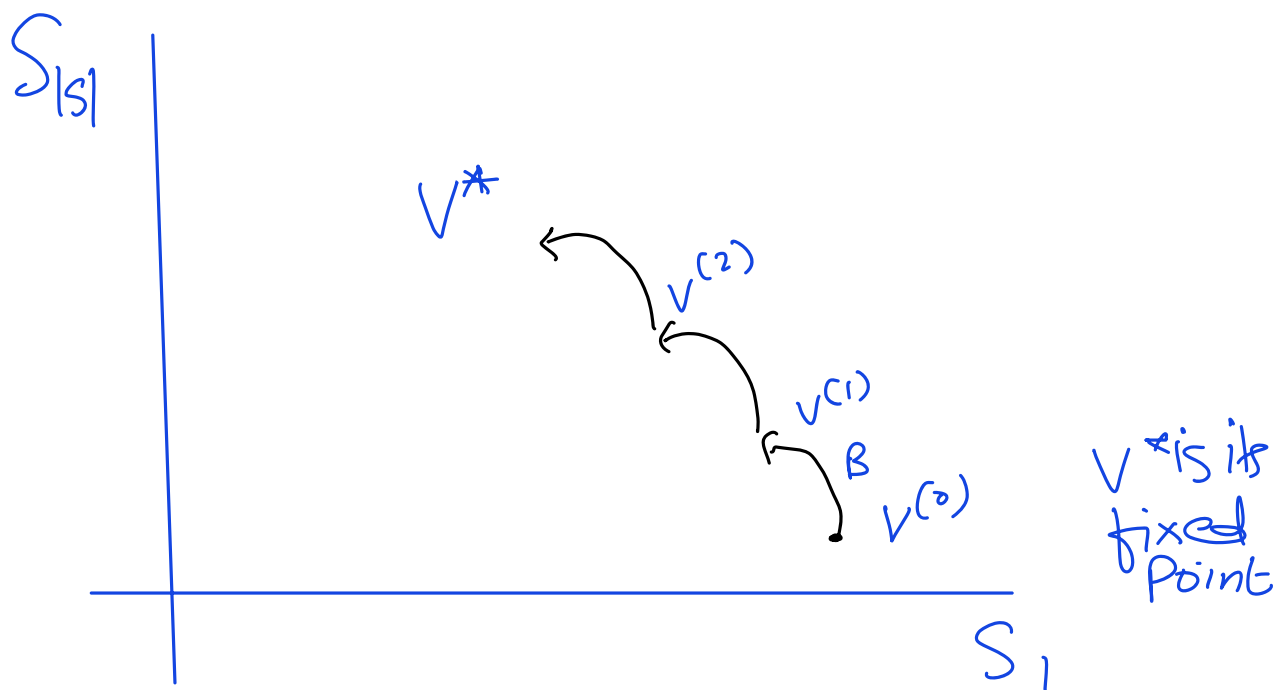
(Greedy Policy
wrt A)

# Value Iteration

Loop:

$$V^{(t+1)} = B(V^{(t)})$$

$$B(V'(s)) = R(s) + \max_a \mathbb{E}_{s' \sim P_{sa}}\left[V(s')\right]$$

$S_{|S|}$

$V^*$

$V^{(2)}$

$V^{(1)}$

$B$

$V^{(0)}$

$V^*$ is its fixed point

$S_1$

# Policy Iteration

Loop

1. $V \leftarrow \pi$ using Policy evaluation

2. $\pi \leftarrow V$ greed Policy

Policy we get in step 2 is different from that in 1

$P_{sa}$ is not given

Model based vs Model free

    Model $= P_{sa}$   Here model refers to environment

Learning Model

$$S_0^{(1)} \xrightarrow{a_0^{(1)}} S_1^{(1)} \cdots - - -$$

$$S_0^{(2)} \xrightarrow{a_0^{(2)}} S_1^{(2)} -\cdot - \cdot -$$

$$\vdots$$

$$\hat{P}_{sa}(s') = \frac{\text{\# of time we took action } a \text{ at state } s \ \& \text{ got to state } s'}{\text{\# of time we took action } a \text{ at state } s}$$

Sometimes $= \dfrac{0}{0} \Rightarrow$ Uniform