# MACHINE LEARNING

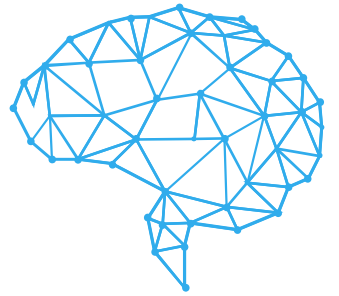## REINFORCEMENT LEARNING

# AGENDA

# AI

INTRODUCTION

# REINFORCEMENT LEARNING

## I N T R O D U C T I O N

There are many **problems**, where it is **very difficult** to **provide** the "**right**" **answers** to the algorithm, as we have been doing for supervised learning.

For example, if we have just built a **four-legged robot** and are **trying** to **program** it to **walk**, then initially we **have no idea** what the "**correct**" **actions** to take are to make it walk **(credit assignment problem).**

In **reinforcement learning**, we will **provide algorithms** only a **reward function**, to **indicate** when the **learning agent** is **doing well** or **poorly**.

---

**It will be the learning algorithm's job to figure out how to choose actions over time to maximize rewards**

# REINFORCEMENT LEARNING
## INTRODUCTION

**Many** of the **applications** of **reinforcement learning include**:

- **Autonomous helicopter flight.**
- **Robot Legged Locomotion.**
- **Cell-phone network routing.**
- **Marketing strategy selection.**
- **Factory control.**
- **Efficient web indexing.**
- **Traffic light control.**
- **Chemical reaction optimization.**
- **Personalized recommendations.**
- **Bidding and Advertisement.**

A **Markov decision process** is a **tuple**:

$$(S, A, \{P_{sa}\}, \gamma, R)$$

- $S$**:** the set of **states** (i.e. all positions and orientations of a helicopter).

- $A$: the set of **actions** (i.e. all directions in which you can push the helicopter's sticks).

- $P_{sa}$: state **transition probability distributions**. For **each state** $s \in S$ and action $a \in A$, $P_{sa}$ is a **distribution** over the **state space.** Thus, $P_{sa}$ **gives** the **distribution** over what **states** the **algorithm** will **transition,** if it takes **action** $a$ in **state** $s$.

$$\sum_{s'} P_{sa}(s') = 1$$

A **Markov decision process** is a **tuple**:

$$(S, A, \{P_{sa}\}, \gamma, R)$$

- $\gamma \in [0, 1)$: is called the **discount factor**.

- $R: S \ X \ A \rightarrow \mathbb{R}$: the **reward function**. (It can be written only as a function of $S$, $R: S \rightarrow \mathbb{R}$).

# REINFORCEMENT LEARNING

## MARKOV DECISION PROCESSES

**EXAMPLE:**

We will **define** a **robot navigation task**, in which you **have** a **robot** in a **grid environment** where **all gray cells** are **obstacles**.

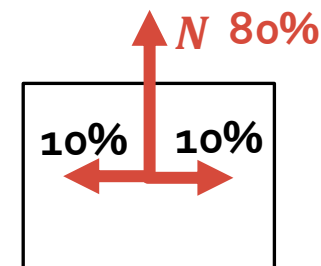The **reward** is **defined by** the **+1** and the **punishment** as **-1**.

The **robot can move** in **all cells except** for the **gray cells.**



11 states

STOCHASTIC ACTIONS

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

$N$ 80%

10%   10%

## MARKOV DECISION PROCESSES

**EXAMPLE:**

The **robot** can **transit from** one **state** to another **state** in **single steps.**

**When** the **robot reaches (4,3) or (4,2)** the **world ends.**

| | | | |
|---|---|---|---|
| $(1,3)$ | $(2,3)$ | $(3,3)$ | $(4,3)$ **+1** |
| $(1,2)$ | $(2,2)$ | $(3,2)$ 80% | $(4,2)$ **-1** |
| $(1,1)$ | $(2,1)$ | $(3,1)$ 10% ↑ 10% | $(4,1)$ |

**Transitions**

$$P_{s_t a_t}(s')$$

$$P_{(3,1)N}\big((3,2)\big) = 0.8$$

$$P_{(3,1)N}\big((4,1)\big) = 0.1$$

$$P_{(3,1)N}\big((2,1)\big) = 0.1$$

$$P_{(3,1)N}\big((3,3)\big) = 0$$

$$\vdots$$

**Rewards**

$$R(s)$$

$$R\big((4,3)\big) = +1$$

$$R\big((4,2)\big) = -1$$

$$R(s) = -0.2$$
**(fuel consumption)**

$$\vdots$$

The dynamics of a **Markov decision process** is as follows:

1. The **algorithm starts** in a **state** $s_0$ and **choose** an **action** $a_0 \in A$.

2. There is a **random transition from state** $s_0$ to **state** $s_1$ **drawn according** to $s_1 \sim P_{s_0 a_0}$.

3. The **process repeats indefinitely** where there is a **random selection** to **traverse** from state $s_t$ to state $s_{t+1}$, which is **drawn from** $s_{t+1} \sim P_{s_t a_t}$.

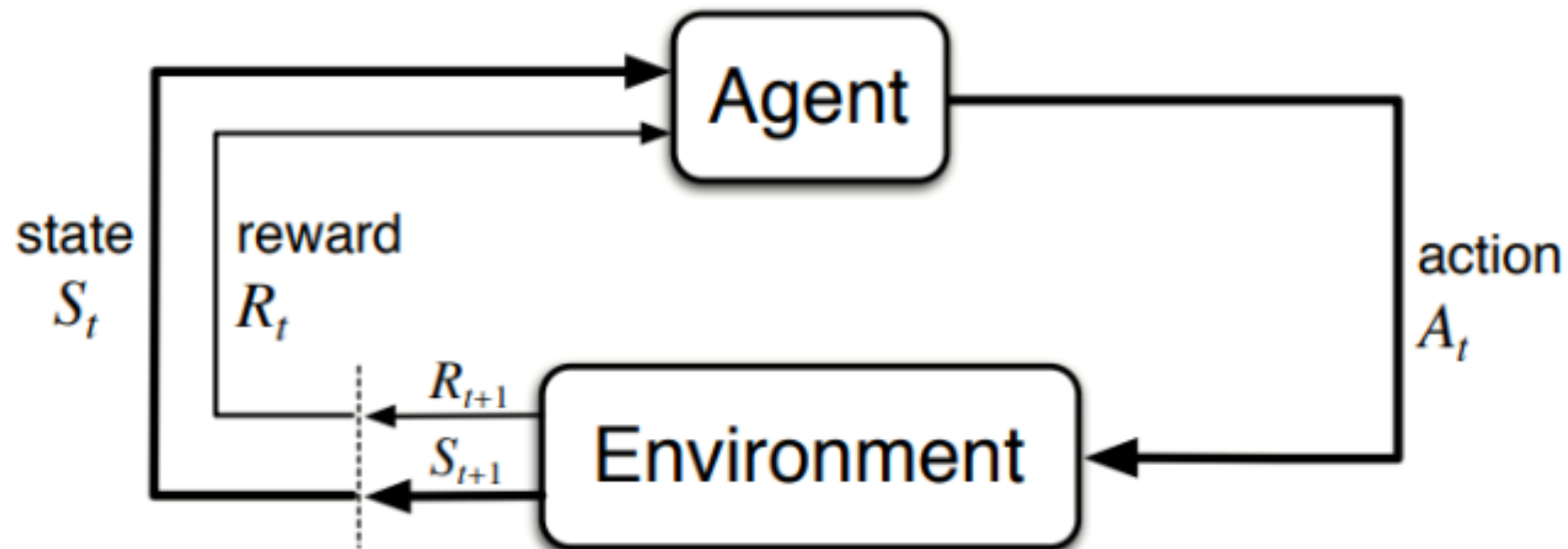$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots$$

4. The **total payoff** is be given by:

$$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots$$
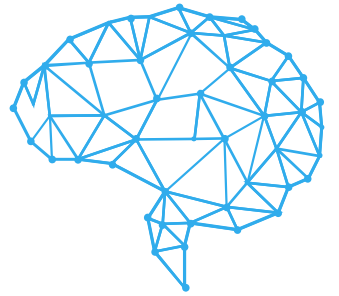
**(Rewards diminish as times goes by).**

# REINFORCEMENT LEARNING
## MARKOV DECISION PROCESSES

The **objective** in **reinforcement learning** is to **choose actions over time** to **maximize** the **expected value** of the **total payoff**:

$$E\big[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots\big]$$

The **reward** at **timestep** $t$ is **discounted** by a **factor** of $\gamma^t$ where $\gamma \in [0,1)$ is the **discount factor**.

Therefore, to **make** this **expectation large**, we would like to:

- **Obtain positive rewards** as **soon** as **possible**.

- **Postpone negative rewards,** as **long as possible**.

## POLICY AND VALUE FUNCTION

**POLICY**

A **policy** is **any function** $\pi : S \rightarrow A$ that **maps from states** to **actions (for every state, what action** is **recommended** to **take** in **that state).**

**EXAMPLE**: (**Optimal Policy**)

| | | | |
|---|---|---|---|
| $(1,3)$ → | $(2,3)$ → | $(3,3)$ → | **+1** |
| $(1,2)$ ↑ | $(2,2)$ | $(3,2)$ ↑ | **-1** |
| $(1,1)$ ↑ | $(2,1)$ ← | $(3,1)$ ← | $(4,1)$ ← |

**MDPs are very good in finding the best tradeoffs between actions (policies).**

## POLICY AND VALUE FUNCTION

**VALUE FUNCTION:**

For **any policy** $\pi$, we **define** the **value function** $V^\pi : S \rightarrow \mathbb{R}$ as the **expected total payoff** upon **starting** in **state** $s$ and **taking actions according** to the **policy** $\pi$.

$$V^\pi(s) = E\left[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \mid s_0 = s, \pi\right]$$

As a **side note**, this **notation** is **technically incorrect** because we **cannot condition** on $\pi$, which **does not represent** a **random variable.**

# POLICY AND VALUE FUNCTION

**EXAMPLE:** Value Function and Policy

**Very bad policy $\pi$**

| | | | |
|---|---|---|---|
| (1,3) → | (2,3) → | (3,3) → | **+1** |
| (1,2) ↓ | (2,2) | (3,2) → | **-1** |
| (1,1) → | (2,1) → | (3,1) ↑ | (4,1) ↑ |

**Value Function $V^{\pi}(s)$**

| | | | |
|---|---|---|---|
| $V^{\pi}((1,3))$<br>$0.52$ | $V^{\pi}((2,3))$<br>$0.73$ | $V^{\pi}((3,3))$<br>$0.77$ | **+1** |
| $V^{\pi}((1,2))$<br>$-0.90$ | | $V^{\pi}((3,2))$<br>$-0.82$ | **-1** |
| $V^{\pi}((1,1))$<br>$-0.88$ | $V^{\pi}((2,1))$<br>$-0.87$ | $V^{\pi}((3,1))$<br>$-0.85$ | $V^{\pi}((4,1))$<br>$-1$ |

$$V^{\pi}(s) = E\big[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \,|\, s_0 = s, \pi\big]$$

**Better to start at the top**

**BELLMAN EQUATIONS:**

We can **represent** the **value function** in the **following way**:

$$V^{\pi}(s) = E\big[R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \cdots \,|s_0 = s, \pi\big]$$

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

**Immediate reward
for starting in state** $s$

**Expected sum of future
discounted rewards.**
$V^{\pi}(s_1)$

Therefore, we can **write** the **equation** as a **recursive form** of **itself** ($s_0 \rightarrow s$ and $s_0 \rightarrow s'$) **using** the **definition** of **expected value** $E[X] = \sum x_i p_i$ ($s'$ is a **random variable**):

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

**BELLMAN EQUATIONS:**

**Given** a **fixed policy** $\pi$, its **value function** $V^\pi$ **satisfies Bellman equations**:

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

The **first term** is the **immediate reward** for **starting** in **state** $s$.

The **second term** gives the **expected sum** of **discounted rewards** obtained after the **first step** in the **MDP.**

# REINFORCEMENT LEARNING
## POLICY AND VALUE FUNCTION

**EXAMPLE:** Bellman Equation

Let us **assume** that we **start** at **state (3,1)** and the **policy** $\pi$ **indicates** to take a **step** in the **north direction** $\pi\big((3,1)\big) = N$.

The **Bellman Equation** is:

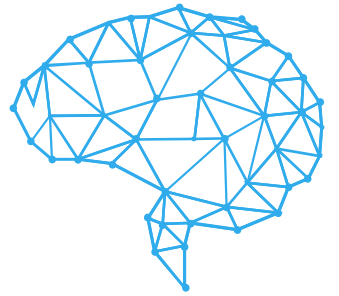$$V^{\pi}(s) = R(s) + \gamma \sum_{s' \in S} P_{s\pi(s)}(s') V^{\pi}(s')$$

| | | | |
|---|---|---|---|
| (1,3) | (2,3) | (3,3) | **+1** |
| (1,2) | (2,2) | (3,2) | **-1** |
| (1,1) | (2,1) | (3,1) | (4,1) |

$$V^{\pi}\big((3,1)\big) = R\big((3,1)\big) + \gamma\big[0.8V^{\pi}\big((3,2)\big) + 0.1V^{\pi}\big((4,1)\big) + 0.1V^{\pi}\big((2,1)\big)\big]$$

You will have **11 variables** and **11 equations**, **each** one **corresponding** to **each state**.

## POLICY AND VALUE FUNCTION

**OPTIMAL VALUE FUNCTION:**

The **optimal value function** will be the **best possible expected sum** of **discounted rewards** that can be attained **using any policy**.

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

The **Bellman representation** will be:

Expected future rewards if we take the BEST action $a$ in state $s$.

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

The **first term** above is the **immediate reward**. The **second term** is the **maximum over all actions** $a$ of the **expected future sum** of **discounted rewards** we'll get **upon after action** $a$.

**OPTIMAL POLICY:**

The **optimal policy** $\pi^*$ indicates the **best action** $a$ to **take** in **state** $s$, which **maximizes** the **expected future rewards**.

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

This **derives** in the **following expression**:

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

In conclusion, $\pi^*$ is the **optimal policy for all states** $s$.

**We can use the same policy $\pi^*$ no matter what the initial state of our MDP is.**

**OPTIMAL POLICY:**

To **find** the **optimal policy** $\pi^*(s)$, we would **need** to **obtain** the **best value function** $V^*(s')$ and **substitute it** in the **equation** we have just obtained:

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

The **problem** is **finding** $V^*(s')$ which corresponds to:

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

The **bottleneck** is that you would **need** to **solve** an **exponentially large number** of **systems** of **equations because** there is an **exponential number** of **policies** $\pi$.
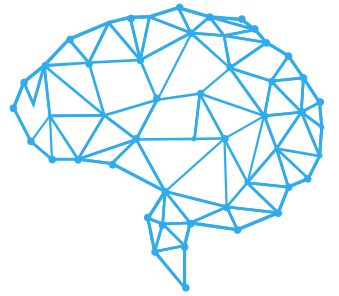
## VALUE ITERATION AND POLICY ITERATION

**VALUE ITERATION:**

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

2. **Repeat until convergence{**

   **For every state, update $V(s) := R(s) + \max\limits_{a \in A} \gamma \sum_{s' \in S} P_{sa}(s') V(s')$.**
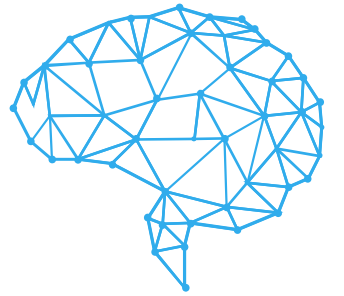
   **}**

3. **Substitute $V^*(s)$ in $\pi^*(s) = \arg\max\limits_{a \in A} \sum_{s' \in S} P_{sa}(s') V^*(s')$.**

The **algorithm repeatedly tries** to **update** the **estimated value function** using **Bellman Equations.** By **implementing this** you **ensure** that the **value function converges** to the **optimum.**

$$V(s) \to V^*(s)$$

## VALUE ITERATION AND POLICY ITERATION

**VALUE ITERATION:**

There are **two possible ways** of **performing** the **updates** in the **inner loop**.

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

1. **Synchronous update** (current estimate → new estimate):
   - **Compute** the **new values** of $V(s)$.
   - **Overwrite** the **old values with** the **new values**.

2. **Asynchronous updates**
   - **Loop over** the **states** (**predefined order**).
   - **Update** the **values** of $V(s)$ **one** at a **time**.

## VALUE ITERATION AND POLICY ITERATION

**EXAMPLE:** Value Iteration

**Running value iteration** on our **previous example**, we obtain the **following results** for $V^*(s)$:

$$\boxed{V^*(s)}$$

| | | | |
|---|---|---|---|
| $V^\pi((1,3))$ <br> 0.86 | $V^\pi((2,3))$ <br> 0.90 | $V^\pi((3,3))$ <br> 0.93 | **+1** |
| $V^\pi((1,2))$ <br> 0.82 | | $V^\pi((3,2))$ <br> 0.69 | **-1** |
| $V^\pi((1,1))$ <br> 0.78 | $V^\pi((2,1))$ <br> 0.75 | $V^\pi((3,1))$ <br> 0.71 | $V^\pi((4,1))$ <br> 0.49 |

# REINFORCEMENT LEARNING

## VALUE ITERATION AND POLICY ITERATION

**EXAMPLE:** Value Iteration

Let us see how we **obtain** $\pi^*(s)$ for a **single state:** $(3, 1)$

$$W = \sum_{s' \in S} P_{sa}(s')V(s') = (0.8 * 0.75) + (0.1 * 0.69) + (0.1 * 0.71) = 0.78$$

$$N = \sum_{s' \in S} P_{sa}(s')V(s') = (0.8 * 0.69) + (0.1 * 0.75) + (0.1 * 0.49) = 0.67$$

$$S = \sum_{s' \in S} P_{sa}(s')V(s') = (0.8 * 0.71) + (0.1 * 0.75) + (0.1 * 0.49) = 0.69$$

$$E = \sum_{s' \in S} P_{sa}(s')V(s') = (0.8 * 0.49) + (0.1 * 0.69) + (0.1 * 0.71) = 0.53$$

The action west maximizes our future rewards if we start at state (3,1).

$V^*(s)$

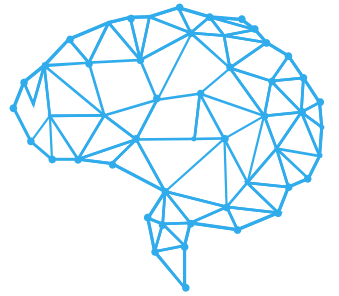| | | | |
|---|---|---|---|
| $V^\pi((1,3))$ <br> $0.86$ | $V^\pi((2,3))$ <br> $0.90$ | $V^\pi((3,3))$ <br> $0.93$ | **+1** |
| $V^\pi((1,2))$ <br> $0.82$ | | $V^\pi((3,2))$ <br> $0.69$ | **-1** |
| $V^\pi((1,1))$ <br> $0.78$ | $V^\pi((2,1))$ <br> $0.75$ | $V^\pi((3,1))$ <br> $0.71$ | $V^\pi((4,1))$ <br> $0.49$ |

## VALUE ITERATION AND POLICY ITERATION

**EXAMPLE:** Value Iteration

Applying the **previous procedure** to **all states,** we **obtain** our **optimal policy** (**best actions** $a = \pi(s)$ to take **at every state** $s$**).**

$$\boxed{\pi^*(s)}$$

| | | | |
|---|---|---|---|
| $(1, 3)$ → | $(2, 3)$ → | $(3, 3)$ → | **+1** |
| $(1, 2)$ ↑ | $(2, 2)$ | $(3, 2)$ ↑ | **-1** |
| $(1, 1)$ ↑ | $(2, 1)$ ← | $(3, 1)$ ← | $(4, 1)$ ← |

## VALUE ITERATION AND POLICY ITERATION

**POLICY ITERATION:**

1. **Initialize $\pi$ randomly.**

2. **Repeat until convergence{**

   **a) Let $V := V^\pi$ (Solve Bellman equations using policy $\pi$).**

   **b) For every state, compute $\pi(s) := R(s) \underset{a \in A}{\operatorname{argmax}} \sum_{s' \in S} P_{sa}(s') V(s')$.**
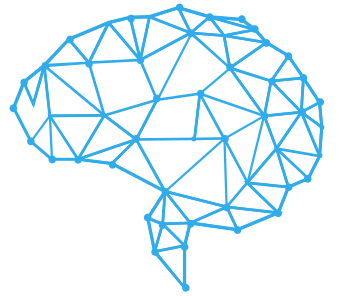
   **}**

The **policy $\pi$ found** in **step (b)** is also **called** the **policy** that is **greedy** with **respect** to $V$. By **implementing this** you **ensure** that the **value function converges** to the **optimum**.

$$V(s) \to V^*(s)$$
$$\pi(s) \to \pi^*(s)$$

## VALUE ITERATION AND POLICY ITERATION

**POLICY ITERATION VS VALUE ITERATION:**

| POLICY ITERATION | VALUE ITERATION |
|---|---|
| For **small MDPs** it is **very fast.** | For **small MDPs** it is **slower.** |
| For **MDPs** with **large state spaces** it is **computationally expensive** (large matrices) | For **MDPs** with **large state spaces** it is has **less computational expense**. |

**Value iteration** seems to be **used more often** than policy iteration.

In **many realistic problems**, we are **not given** state **transition probabilities** and **rewards explicitly**. Therefore, we need to **estimate them from data.**

More generally, in a Markov Decision Process $(S, A, \{P_{sa}\}, \gamma, R)$ the **states** $S$, the **actions** $A$ and the **discount factors** $\gamma$ are **almost** always **known.** The **discount factor** is **chosen based** on **how much tradeoff** do you want **between current** and **future rewards**.

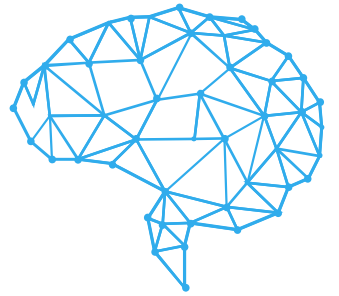We can **easily derive** the **maximum likelihood estimates** for the **state transition probabilities**:

$$P_{sa}(s') = \frac{\# \ times \ we \ took \ action \ a \ in \ state \ s \ and \ got \ to \ s'}{\# \ times \ we \ took \ action \ a \ in \ state \ s}$$

If there is the **case** that you have **never taken action** $a$ **in state** $s$ (0/0) we can **compute** the **uniform distribution over all states** :

$$P_{sa}(s') = \frac{1}{|S|}$$

**Putting together model learning** and **value iteration**, here is **one possible algorithm** for **learning** in an **MDP** with **unknown state transition probabilities**:

1. **Initialize $\pi$ randomly.**

2. **Repeat** {

   a) **Execute $\pi$** in the **MDP** for some **number** of **trials**

   b) **Using** the **accumulated experience** in the **MDP**, **update** our **estimates** for $P_{sa}$ (and $R$, **if applicable**).

   c) **Apply value iteration** with the **estimated state transition probabilities** and **rewards** to **get** a **new estimated value function $V$** (**use** this **update** to **initialize $V$** in **next iteration**).

   (d) **Update $\pi$** to be the **greedy policy** with **respect** to $V$**.**

   }