

# Artificial Intelligence, Reinforcement Learning & Sequential Decision Problems

(CIV6540 - Probabilistic Machine Learning for Civil Engineers)

Professor: James-A. Goulet

Département des génies civil, géologique et des mines  
Polytechnique Montréal



Chapter 14 – Goulet (2020)  
*Probabilistic Machine Learning for Civil Engineers*

MIT Press



Chapter 17 – Russell, S. and Norvig, P. (1995)  
*Artificial Intelligence, A modern approach*

Prentice-Hall



Sutton, R. S. and Barto, A. G. (2018)  
*Reinforcement learning: An introduction*  
MIT Press, 2nd edition

# Artificial intelligence – AI

**AI** : reproduction of an intelligent behavior by a machine.  
→ perception and dynamic interaction with its environment through the process of **taking actions** for achieving goals.

## Level 1 - Limited AI

Beyond human capacity in a specific tasks

### Behind limited AI:

Machine learning

### Behind machine learning:

Probability theory / Decision Theory  
/ Linear Algebra / Optimization/...



[Google images]

## Reinforcement Learning

AI → Reinforcement Learning → Sequential Decisions



$$\mathcal{S} = \left\{ \text{NES Cartridge} \right\}$$

$$\mathcal{A} = \left\{ \text{NES Controller} \right\}$$

$$s_t = \left\{ \text{NES Screen} \right\}$$

$$r(s_t) = \{ \text{fct}(\cdot) \}$$

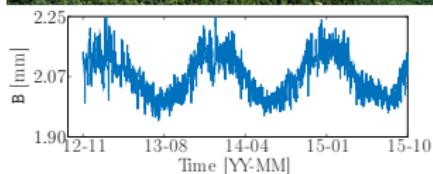
$$p(s_{t+1}|s_t, a) = \text{Unknown...}$$

**Reinforcement learning** →

**Long term** expected utility to take an action  $a$  if we are in a **state**  $s$

[<https://youtu.be/QVyu9oVyh9Q>, Google images]

Anomaly detection & Structural Health Monitoring



## Reinforcement learning →

$S = \{\text{SHM system}\}$

$$\mathcal{A} = \{\text{alarm}, \neg \text{alarm}\}$$

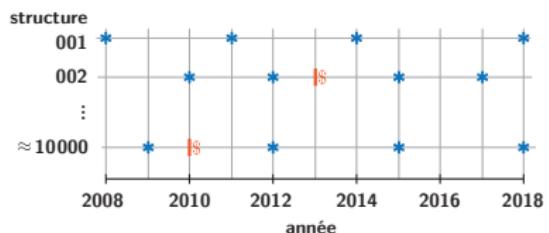
$s_t = \{\mathbf{x}_t\}$  (e.g. SSM, M7)

$$r(s_t, a_t) = \{\text{fct}(s_t, a_t)\}$$

$p(s_{t+1}|s_t, a)$  = unknown...

**Long term** expected utility to take an action  $a$  if we are in a state  $s$

# Network-Scale Infrastructure Maintenance Planning



$\mathcal{S}$  = {Infrastructure network}

$$\mathcal{A} = \{\text{wait, repair, replace}\}$$

$s_t = \{\%A, \%B, \%C, \%D\}$

$$r(s_t) = \{\text{fct}(s_t)\}$$

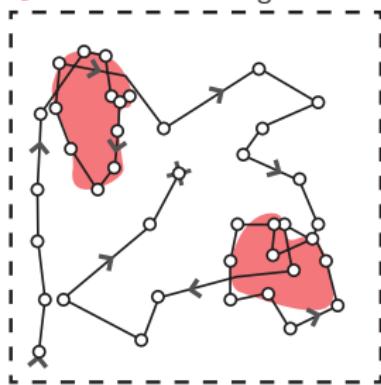
$$p(s_{t+1} | s_t, a) = \text{unknown...}$$

## Reinforcement learning →

**Long term** expected utility to take an action  $a$  if we are in a state  $s$

# Contaminated site characterization

- In-situ observation
- ▣ Site to be decontaminated
- Contaminated regions



$$\mathcal{S} = \{\text{Contaminated site}\}$$

$$\mathcal{A} = \{\text{Sample } (\Delta x, \Delta y), \text{ stop}\}$$

$$s_t = \{[\text{Hg}]_s\}$$

$$r(s_t, a_t) = \{\text{fct}(s_t, a_t)\}$$

$$p(s_{t+1}|s_t, a) = \text{unknown...}$$

**Reinforcement learning** →

**Long term** expected utility to take an  
**action  $a$**  if we are in a **state  $s$**

## Potential applications of RL in civil engineering

## Emerald ash borer – Tree cutting policy



$$\begin{aligned}
 \mathcal{S} &= \{\text{Tree population}\} \\
 \mathcal{A} &= \{\text{cut, wait}\} \\
 s_t &= \{[\text{insects}]_{\mathcal{S}}\} \\
 r(s_t) &= \{\text{fct}(s_t)\} \\
 p(s_{t+1}|s_t, a) &= \text{unknown...}
 \end{aligned}$$

## Reinforcement learning →

**Long term** expected utility to take an action  $a$  if we are in a **state**  $s$

Professor: I-A Goulet

[Google images]

© Polytechnique Montréal

# Reinforcement learning

When we can explicitly represent  $p(s_{t+1}|s_t, a)$

→ **Model-based RL: MDP/POMDP, etc...**

When we cannot represent  $p(s_{t+1}|s_t, a)$  explicitly

→ **Model-free RL: Q-learning, SARSA, etc...**

# Module #10 Outline

## Intro

## Nomenclature Markov Decision Process, MDP Q-Learning

### Topics organization

#### Background

- 1 Revision probability & linear algebra 
- 2 Probability distributions  

#### Machine Learning Basics

- 0 Introduction 
- 3 Bayesian Estimation  $p(A|B) = \frac{p(B|A)p(A)}{p(B)}$
- 4 MCMC sampling & Newton  

#### Supervised learning

- 5 Regression 
- 6 Classification  
- 7 LSTM networks for time series 

#### Unsupervised learning

- 7 State-space model for time-series 

#### Decision Making & RL

- 8 Decision Theory 
- 9 AI & Sequential decision problems 



# Section Outline

---

## Nomenclature

2.1 Variables

2.2 Example #1 - Infrastructure Maintenance Planning

---

# RL sequential decision problem – Nomenclature

States:  $s \in \mathbb{R}$  or  $s_j \in \mathcal{S} = \{s_1, s_2, \dots, s_S\}$

continuous

discrete

Actions:  $a_i \in \mathcal{A} = \{a_1, a_2, \dots, a_A\}$

Transition model:  $p(s'|s_j, a_i) = \Pr(s' = s_k | s_j, a_i), \forall k = 1 : S$

Reward:  $r(s, a, s') \in \mathbb{R}, (\equiv r(s, a), \equiv r(s))$

**Policy:**  $\pi = \{a(s_1), a(s_2), \dots, a(s_S)\}$

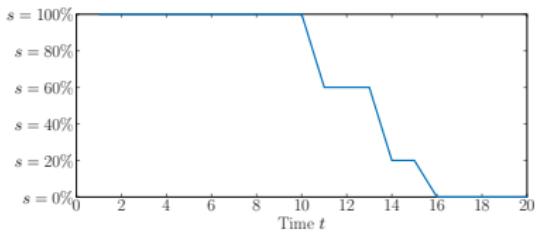
## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Definitions

Bridges are inspected every two years:



$$s \in \mathcal{S} = \left\{ \begin{array}{l} 100\% \\ 80\% \\ 60\% \\ 40\% \\ 20\% \\ 0\% \end{array} \right\}$$



# Data

$$D = \left\{ \begin{array}{l} \{s_1, s_2, \dots, s_T\}_1 \\ \{s_1, s_2, \dots, s_T\}_2 \\ \vdots \\ \{s_1, s_2, \dots, s_T\}_N \end{array} \right\}_{N \times T}$$

[Google images]

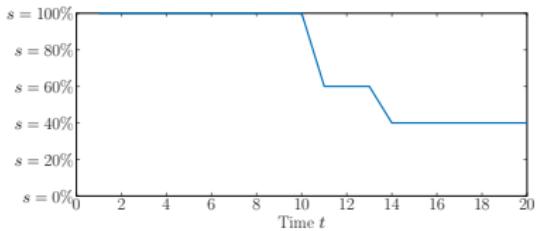
## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Definitions

Bridges are inspected every two years:



$$s \in \mathcal{S} = \left\{ \begin{array}{l} 100\% \\ 80\% \\ 60\% \\ 40\% \\ 20\% \\ 0\% \end{array} \right\}$$



## Data

$$D = \left\{ \begin{array}{l} \{s_1, s_2, \dots, s_T\}_1 \\ \{s_1, s_2, \dots, s_T\}_2 \\ \vdots \\ \{s_1, s_2, \dots, s_T\}_N \end{array} \right\}_{N \times T}$$

[Google images]

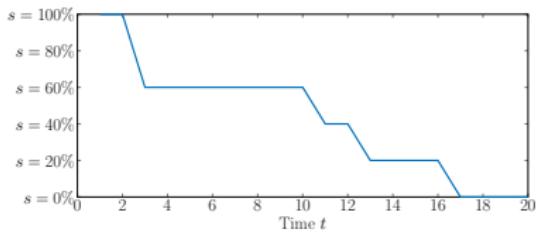
## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Definitions

Bridges are inspected every two years:



$$s \in \mathcal{S} = \left\{ \begin{array}{l} 100\% \\ 80\% \\ 60\% \\ 40\% \\ 20\% \\ 0\% \end{array} \right\}$$



## Data

$$\mathcal{D} = \left\{ \begin{array}{l} \{s_1, s_2, \dots, s_T\}_1 \\ \{s_1, s_2, \dots, s_T\}_2 \\ \vdots \\ \{s_1, s_2, \dots, s_T\}_N \end{array} \right\}_{N \times T}$$

[Google images]

# Infrastructure Maintenance - Actions

$$a \in \mathcal{A} = \{\text{do nothing, maintain, replace}\}$$

## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Model

Employ the dataset  $\mathcal{D}$  to build a **Markovian model**

$$p(s_{t+1}|s_t, a) \equiv p(s_{t+1}|s_{1:t}, a)$$

$$\hat{\Pr}(s_{t+1} = k | s_t = i, a = \text{do nothing}) = \frac{\#\{s_t = i, s_{t+1} = k\}}{\#\{s_t = i\}}$$

$$p(s_{t+1}|s_t = 1, a = \text{do nothing}) = \begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \end{bmatrix}_{s_{t+1}}$$

## Example #1 - Infrastructure Maintenance Planning

# Infrastructure Maintenance - Model

Employ the dataset  $\mathcal{D}$  to build a **Markovian model**

$$p(s_{t+1}|s_t, a) \equiv p(s_{t+1}|s_{1:t}, a)$$

$$\hat{\Pr}(s_{t+1} = k | s_t = i, a = \text{do nothing}) = \frac{\#\{s_t = i, s_{t+1} = k\}}{\#\{s_t = i\}}$$

$$p(s_{t+1}|s_t = 2, a = \text{do nothing}) = \begin{bmatrix} 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \end{bmatrix}$$

## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Model

Employ the dataset  $\mathcal{D}$  to build a **Markovian model**

$$p(s_{t+1}|s_t, a) \equiv p(s_{t+1}|s_{1:t}, a)$$

$$\hat{\Pr}(s_{t+1} = k | s_t = i, a = \text{do nothing}) = \frac{\#\{s_t = i, s_{t+1} = k\}}{\#\{s_t = i\}}$$

$$p(s_{t+1}|s_t, a = \text{do nothing}) = \begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}_{S \times S}$$

## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Model (cont.)

$$p(s_{t+1}|s_t, a = \text{maintain}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}_{S \times S}$$

$$p(s_{t+1}|s_t, a = \text{replace}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}_{S \times S}$$

## Example #1 - Infrastructure Maintenance Planning

## Infrastructure Maintenance - Rewards

**Rewards** for being in a state  $s$  or for taking an action  $a$

$$\begin{aligned} r(S) &= \overbrace{\text{AADTF}10^5 \text{ users} \cdot 365 \text{ days} \cdot 3\$/\text{user} \cdot \text{capacity}(S)}^{\text{AATF}} \\ &= \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} M\$ \end{aligned}$$

$$\text{capacity}(S) = \left\{ \underbrace{1}_{s=100\%}, \underbrace{1}_{80\%}, \underbrace{1}_{60\%}, \underbrace{0.90}_{40\%}, \underbrace{0.75}_{20\%}, \underbrace{0}_{0\%} \right\}$$

$$r(A) = \left\{ \underbrace{0}_{\text{do nothing}}, \underbrace{-5}_{\text{repair}}, \underbrace{-20}_{\text{replace}} \right\} M\$$$

$$r(s, a, s') = r(s, a) = r(s) + r(a)$$

## Illustrative example - Optimal actions

Given that the transition model is **Markovian**, the goal of the **Markovian Decision Process** is to identify the **optimal actions** a to be taken for each state  $s$  so that a *policy*

$$\pi^*(\mathcal{S}) = \{a(s_1), a(s_2), \dots, a(s_S)\}$$



# Section Outline

---

## Markov Decision Process, MDP

- 3.1 Utility over time
  - 3.2 Value Iteration
  - 3.3 Policy Iteration
  - 3.4 Partially Observable MDP
-

# Planning horizon

## Finite planning horizon:

The rewards and utilities are considered over a fixed period of time.

The optimal policy is **non-stationary**;  $\pi_t^*(s)$  depends on  $t$ .

## Infinite planning horizon:

The rewards and utilities are considered over an infinite period.

The optimal policy is **stationary**;  $\pi^*(s)$  does not depend on  $t$ .

We will only study problems with an **infinite planning horizon**

# Utility for an infinite planning horizon

For a discount factor  $\gamma \in (0, 1[$ , the discounted sum of rewards over an infinite planning horizon is finite (note: interest rate  $= \frac{1}{\gamma} - 1$ )

$$\begin{aligned}\mathbb{U}(s_{t=0}, s_{t=1}, \dots, s_{t=\infty}) &= r(s_0) + \gamma r(s_1) + \gamma^2 r(s_2) + \dots \\ &= \sum_{t=0}^{\infty} \gamma^t r(s_t) \\ &\leq \frac{\max_{s \in \mathcal{S}} r(s)}{1 - \gamma}\end{aligned}$$

⚠ The issue is that we do not know  $s_1, s_2, \dots$   
we only know  $\underbrace{p(s_{t+1}|s_t, a) \equiv p(s'|s, a)}_{\text{transition model}}$

# Expected utility for an infinite planning horizon

Given  $p(s'|s, a)$  and  $r(s)$ , we can compute the expected utility

## Expected utility

$$\bar{U}(s) \equiv \bar{U}(s, \pi) \equiv \mathbb{E}[\mathbb{U}(s, \pi)] = r(s) + \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^t r(S_t) \right]$$

## Optimal policy

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \bar{U}(s', \pi^*))$$

⚠  $\pi^*$  appears on both sides of the equality → **iterative resolution**

# Utility for an infinite planning horizon – Bellman

**Bellman equation:** “*The utility of a state is the immediate reward for that state plus the expected discounted utility of the next state given that the optimal action is taken*”

$$\begin{aligned}\overline{\mathbb{U}}(s) &= \max_{a \in \mathcal{A}} \mathbb{E} [r(s, a, s') + \gamma \overline{\mathbb{U}}(s')] \\ &= \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}(s'))\end{aligned}$$

## Bellman update & Value iteration

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}^{(i)}(s'))$$

## Algorithm 1: Value iteration

- 1 define: states  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}$ ;
- 2 define: transition model  $p(s'|s, a)$ ;
- 3 define: rewards  $r(s, a, s')$ , discount  $\gamma$ , convergence tol.  $\eta$ ;
- 4 initialize:  $\overline{\mathbb{U}}'(s) = \overline{\mathbb{U}}(s) = 0$ ;
- 5 **while**  $\overline{\mathbb{U}}(s) = 0$  or  $|\overline{\mathbb{U}}'(s) - \overline{\mathbb{U}}(s)| \geq \eta$  **do**
- 6      $\overline{\mathbb{U}}(s) \leftarrow \overline{\mathbb{U}}'(s)$ ;
- 7     **for**  $s \in \mathcal{S}$  **do**
- 8         **for**  $a \in \mathcal{A}$  **do**
- 9              $\overline{\mathbb{U}}(s, a) = \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}(s'))$ ;
- 10          $\pi^*(s) \leftarrow a^* = \arg \max_{a \in \mathcal{A}} \overline{\mathbb{U}}(s, a)$ ;
- 11          $\overline{\mathbb{U}}'(s) \leftarrow \overline{\mathbb{U}}(s, a^*)$ ;
- 12 **return:**  $\overline{\mathbb{U}}(s) = \overline{\mathbb{U}}'(s)$ ,  $\pi^*(s)$ ;

# Infrastructure Maintenance - Summary

$$s \in \mathcal{S} = \left\{ \begin{array}{l} 100\% \\ 80\% \\ 60\% \\ 40\% \\ 20\% \\ 0\% \end{array} \right\}$$

$$a \in \mathcal{A} = \{\text{Do nothing, Maintain, Replace}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{M\$}$$

$$r(s, a, s') = r(s) + r(a)$$

$$p(s_{t+1}|s_t, a = \text{Do nothing}) = \begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}_{s \times s}$$

$$p(s_{t+1}|s_t, a = \text{Maintain}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}_{s \times s}$$

$$p(s_{t+1}|s_t, a = \text{Replace}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}_{s \times s}$$

Value Iteration

# Infrastructure Maintenance - Value iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing, Maintain, Replace}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

$$p(s'|s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Iteration: 1 ,  $\forall s \in \mathcal{S}$ 

$$\overline{U}_0(\mathcal{S}) = \{0, 0, 0, 0, 0, 0\} \text{ M\$}$$

$$\overline{U}_1(s) = \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s'|s, a_1) \cdot (r(s, a_1) + \gamma \overline{U}_0(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s'|s, a_2) \cdot (r(s, a_2) + \gamma \overline{U}_0(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s'|s, a_3) \cdot (r(s, a_3) + \gamma \overline{U}_0(s')) \right\}$$

$$\overline{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \overline{U}^{(2)}(1) = 223 \text{ M\$}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Value Iteration

# Infrastructure Maintenance - Value iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing, Maintain, Replace}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

$$p(s'|s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Iteration: 2 ,  $\forall s \in \mathcal{S}$ 

$$\bar{U}_1(\mathcal{S}) = \{110, 211, 309, 393, 459, 440\} \text{ M\$}$$

$$\bar{U}_2(s) = \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s'|s, a_1) \cdot (r(s, a_1) + \gamma \bar{U}_1(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s'|s, a_2) \cdot (r(s, a_2) + \gamma \bar{U}_1(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s'|s, a_3) \cdot (r(s, a_3) + \gamma \bar{U}_1(s')) \right\}$$

$$\bar{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \bar{U}^{(2)}(1) = 223 \text{ M\$}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Value Iteration

# Infrastructure Maintenance - Value iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing, Maintain, Replace}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

$$p(s'|s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Iteration: 3 ,  $\forall s \in \mathcal{S}$ 

$$\bar{U}_2(\mathcal{S}) = \{223, 329, 430, 511, 572, 550\} \text{ M\$}$$

$$\bar{U}_3(s) = \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s'|s, a_1) \cdot (r(s, a_1) + \gamma \bar{U}_2(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s'|s, a_2) \cdot (r(s, a_2) + \gamma \bar{U}_2(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s'|s, a_3) \cdot (r(s, a_3) + \gamma \bar{U}_2(s')) \right\}$$

$$\bar{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \bar{U}^{(2)}(1) = 223 \text{ M\$}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

## Value Iteration

Value iteration -  $\bar{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \bar{U}^{(2)}(1) = 223 \text{ M\$}$

$$\bar{U}^{(0)}(\mathcal{S}) = \{ \overbrace{0}^{s=100\%}, \overbrace{0}^{80\%}, \overbrace{0}^{60\%}, \overbrace{0}^{40\%}, \overbrace{0}^{20\%}, \overbrace{0}^{0\%} \} \text{ M\$}$$

$$\bar{U}^{(1)}(\mathcal{S}) = \{ 110, 211, 309, 393, 459, 440 \} \text{ M\$}$$

$$\bar{U}^{(2)}(\mathcal{S}) = \{ 223, 329, 430, 511, 572, 550 \} \text{ M\$}$$

$$\vdots$$

$$\bar{U}^{(536)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$\bar{U}^{(537)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$r(\mathcal{S}) = \{ 109.5, 109.5, 109.5, 98.6, 82.1, 0 \} \text{ M\$}$$

$$r(\mathcal{A}) = \{ 0, -5, -20 \} \text{ M\$}$$

$$p(s' | s, \text{do nothing}) =$$

$$\begin{bmatrix} \color{red}{0.95} & \color{red}{0.03} & \color{red}{0.02} & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\bar{U}^{(2)}(1) \leftarrow r(1) + \max_{a \in \mathcal{A}} \left( \sum_{s' \in \mathcal{S}} p(s' | s, a) \cdot (r(a) + \gamma \bar{U}^{(1)}(s')) \right)$$

$$\leftarrow \underbrace{109.5 \text{ M\$}}_{r(s=1)} + \max_{a \in \mathcal{A}} \left\{ \begin{array}{l} \color{red}{0.95} \cdot (0.97 \cdot 110 \text{ M\$}) + \color{red}{0.03} \cdot (0.97 \cdot 211 \text{ M\$}) + \color{red}{0.02} \cdot (0.97 \cdot 309 \text{ M\$}) = \boxed{113.5 \text{ M\$}} \\ \color{blue}{1} \cdot (-5 + 0.97 \cdot 110 \text{ M\$}) = 101.7 \text{ M\$} \\ \color{blue}{1} \cdot (-20 + 0.97 \cdot 110 \text{ M\$}) = 86.7 \text{ M\$} \end{array} \right.$$

$$= 109.5 \text{ M\$} + 113.5 \text{ M\$}$$

$$= 223 \text{ M\$} \quad (\pi^*(s=1) = a^* = \text{do nothing})$$

Value Iteration

Value iteration -  $\overline{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \overline{U}^{(2)}(1) = 223 \text{ M\$}$

$$\overline{U}^{(0)}(\mathcal{S}) = \{ \underbrace{0}_{s=100\%}, \underbrace{0}_{80\%}, \underbrace{0}_{60\%}, \underbrace{0}_{40\%}, \underbrace{0}_{20\%}, \underbrace{0}_{0\%} \} \text{ M\$}$$

$$\overline{U}^{(1)}(\mathcal{S}) = \{ 110, 211, 309, 393, 459, 440 \} \text{ M\$}$$

$$\overline{U}^{(2)}(\mathcal{S}) = \{ 223, 329, 430, 511, 572, 550 \} \text{ M\$}$$

$$\vdots$$

$$\overline{U}^{(536)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$\overline{U}^{(537)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$r(\mathcal{S}) = \{ 109.5, 109.5, 109.5, 98.6, 82.1, 0 \} \text{ M\$}$$

$$r(\mathcal{A}) = \{ 0, -5, -20 \} \text{ M\$}$$

$$p(s' | s, \text{repair}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\overline{U}^{(2)}(1) \leftarrow r(1) + \max_{a \in \mathcal{A}} \left( \sum_{s' \in \mathcal{S}} p(s' | s, a) \cdot (r(a) + \gamma \overline{U}^{(1)}(s')) \right)$$

$$\leftarrow \overbrace{109.5 \text{ M\$}}^{r(s=1)} + \max_{a \in \mathcal{A}} \left\{ \begin{array}{l} 0.95 \cdot (0.97 \cdot 110 \text{ M\$}) + 0.03 \cdot (0.97 \cdot 211 \text{ M\$}) + 0.02 \cdot (0.97 \cdot 309 \text{ M\$}) = 113.5 \text{ M\$} \\ 1 \cdot (-5 + 0.97 \cdot 110 \text{ M\$}) = 101.7 \text{ M\$} \\ 1 \cdot (-20 + 0.97 \cdot 110 \text{ M\$}) = 86.7 \text{ M\$} \end{array} \right.$$

$$= 109.5 \text{ M\$} + 113.5 \text{ M\$}$$

$$= 223 \text{ M\$} \quad (\pi^*(s=1) = a^* = \text{do nothing})$$

## Value Iteration

Value iteration -  $\overline{U}^{(1)}(1) = 109.5 \text{ M\$} \rightarrow \overline{U}^{(2)}(1) = 223 \text{ M\$}$

$$\overline{U}^{(0)}(\mathcal{S}) = \{ \underbrace{0}_{s=100\%}, \underbrace{0}_{80\%}, \underbrace{0}_{60\%}, \underbrace{0}_{40\%}, \underbrace{0}_{20\%}, \underbrace{0}_{0\%} \} \text{ M\$}$$

$$\overline{U}^{(1)}(\mathcal{S}) = \{ 110, 211, 309, 393, 459, 440 \} \text{ M\$}$$

$$\overline{U}^{(2)}(\mathcal{S}) = \{ 223, 329, 430, 511, 572, 550 \} \text{ M\$}$$

$$\vdots$$

$$\overline{U}^{(536)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$\overline{U}^{(537)}(\mathcal{S}) = \{ 3640, 3635, 3630, 3615, 3592, 3510 \} \text{ M\$}$$

$$r(\mathcal{S}) = \{ 109.5, 109.5, 109.5, 98.6, 82.1, 0 \} \text{ M\$}$$

$$r(\mathcal{A}) = \{ 0, -5, -20 \} \text{ M\$}$$

$$p(s' | s, \text{replace}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

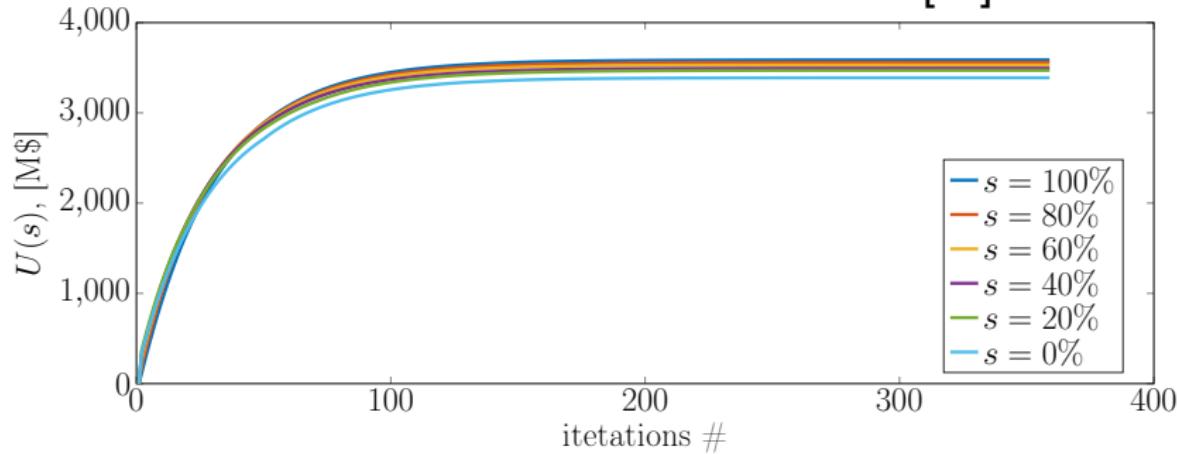
$$\overline{U}^{(2)}(1) \leftarrow r(1) + \max_{a \in \mathcal{A}} \left( \sum_{s' \in \mathcal{S}} p(s' | s, a) \cdot (r(a) + \gamma \overline{U}^{(1)}(s')) \right)$$

$$\leftarrow \overbrace{109.5 \text{ M\$}}^{r(s=1)} + \max_{a \in \mathcal{A}} \left\{ \begin{array}{l} 0.95 \cdot (0.97 \cdot 110 \text{ M\$}) + 0.03 \cdot (0.97 \cdot 211 \text{ M\$}) + 0.02 \cdot (0.97 \cdot 309 \text{ M\$}) = 113.5 \text{ M\$} \\ 1 \cdot (-5 + 0.97 \cdot 110 \text{ M\$}) = 101.7 \text{ M\$} \\ 1 \cdot (-20 + 0.97 \cdot 110 \text{ M\$}) = 86.7 \text{ M\$} \end{array} \right.$$

$$= 109.5 \text{ M\$} + 113.5 \text{ M\$}$$

$$= 223 \text{ M\$} \quad (\pi^*(s=1) = a^* = \text{do nothing})$$

# Infrastructure Maintenance - Value iteration [m]



$$\overline{U}(S) = \{3640, 3635, 3630, 3615, 3592, 3510\} \text{ M\$}$$

$$\pi^*(S) = \{\underbrace{\text{Do nothing}}_{s=100\%}, \underbrace{\text{Maintain}}_{s=80\%}, \underbrace{\text{Maintain}}_{s=60\%}, \underbrace{\text{Maintain}}_{s=40\%}, \underbrace{\text{Replace}}_{s=20\%}, \underbrace{\text{Replace}}_{s=0\%}\}$$

[CIV\_ML/Value\_iteration\_bridge.m]

# Bellman update & Value iteration

## Value iteration

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}^{(i)}(s'))$$

The max operation is computationally expensive if  $\mathcal{A}$  is large  
**solution: policy iteration**

A policy  $\pi = \{a(s_1), a(s_2), \dots, a(s_S)\}$  defines an action to be taken for each state  $s \in \mathcal{S}$

## Policy iteration

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \sum_{s' \in \mathcal{S}} p(s'|s, \underbrace{\pi_i(s)}_{a^*}) \cdot (\underbrace{r(s, \underbrace{\pi_i(s)}_{a^*}, s')}_{a^*} + \gamma \overline{\mathbb{U}}^{(i)}(s'))$$

---

## Algorithm 2: Policy iteration

---

- 1 define: states  $s \in \mathcal{S}$ , actions  $a \in \mathcal{A}$ ;
- 2 define: transition model  $p(s'|s, a)$ ;
- 3 define: rewards  $r(s, a, s')$ , discount  $\gamma$ ;
- 4 initialize:  $\overline{\mathbb{U}}(s) = 0$ ,  $\pi'(s) = \pi(s) \in \mathcal{A}$ ;
- 5 **while**  $\pi'(s) \neq \pi(s)$  **do**
- 6    $\pi(s) = \pi'(s)$ ;
- 7    $\overline{\mathbb{U}}(s) = \text{Policy eval.}(\overline{\mathbb{U}}(s), \pi(s), p(s'|s, a), r(s, a, s'), \gamma)$ ;
- 8   **for**  $s \in \mathcal{S}$  **do**
- 9      $\pi'(s) \leftarrow \arg \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}(s'))$ ;
- 10 **return:**  $\overline{\mathbb{U}}(s)$ ,  $\pi^*(s) = \pi'(s)$ ;

---

Policy Iteration

# Infrastructure Maintenance - Policy iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing (1), Maintain (2), Replace (3)}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

Iteration: 1 ,  $\forall s \in \mathcal{S}$ 

$$p(s' | s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\pi_0(\mathcal{S}) = \{1, 1, 1, 1, 1, 1\}$$

$$\bar{U}_1(\mathcal{S}) = \{2063, 1290, 768, 455, 197, 0\} \text{ M\$}$$

$$\pi_1(s) = \arg \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s' | s, a_1) (r(s, a_1) + \gamma \bar{U}_1(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s' | s, a_2) (r(s, a_2) + \gamma \bar{U}_1(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s' | s, a_3) (r(s, a_3) + \gamma \bar{U}_1(s')) \right\}$$

$$\pi_5(\mathcal{S}) = \pi_6(\mathcal{S}) = \pi^*(\mathcal{S}) = \{1, 2, 2, 2, 3, 3\}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Policy Iteration

# Infrastructure Maintenance - Policy iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing (1), Maintain (2), Replace (3)}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

Iteration: 2 ,  $\forall s \in \mathcal{S}$ 

$$p(s' | s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\pi_1(\mathcal{S}) = \{2, 2, 3, 3, 3, 3\}$$

$$\bar{U}_2(\mathcal{S}) = \{3483, 3483, 3468, 3457, 3441, 3359\} \text{ M\$}$$

$$\pi_2(s) = \arg \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s' | s, a_1) (r(s, a_1) + \gamma \bar{U}_2(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s' | s, a_2) (r(s, a_2) + \gamma \bar{U}_2(s')) \right. \\ \left. \sum_{s' \in \mathcal{S}} p(s' | s, a_3) (r(s, a_3) + \gamma \bar{U}_2(s')) \right\}$$

$$\pi_5(\mathcal{S}) = \pi_6(\mathcal{S}) = \pi^*(\mathcal{S}) = \{1, 2, 2, 2, 3, 3\}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Policy Iteration

# Infrastructure Maintenance - Policy iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing (1), Maintain (2), Replace (3)}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

Iteration: 3 ,  $\forall s \in \mathcal{S}$

$$p(s' | s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\pi_2(\mathcal{S}) = \{1, 1, 2, 2, 3, 3\}$$

$$\bar{U}_3(\mathcal{S}) = \{3622, 3606, 3603, 3588, 3576, 3494\} \text{ M\$}$$

$$\pi_3(s) = \arg \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s' | s, a_1) (r(s, a_1) + \gamma \bar{U}_3(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s' | s, a_2) (r(s, a_2) + \gamma \bar{U}_3(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s' | s, a_3) (r(s, a_3) + \gamma \bar{U}_3(s')) \right\}$$

$$\pi_5(\mathcal{S}) = \pi_6(\mathcal{S}) = \pi^*(\mathcal{S}) = \{1, 2, 2, 2, 3, 3\}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

Policy Iteration

# Infrastructure Maintenance - Policy iteration

$$s \in \mathcal{S} = \{100\%, 80\%, 60\%, 40\%, 20\%, 0\% \}$$

$$a \in \mathcal{A} = \{\text{Do nothing (1), Maintain (2), Replace (3)}\}$$

$$r(\mathcal{S}) = \{109.5, 109.5, 109.5, 98.6, 82.1, 0\} \text{ M\$}$$

$$r(\mathcal{A}) = \{0, -5, -20\} \text{ M\$}$$

$$\gamma = 0.97$$

Iteration: 6 ,  $\forall s \in \mathcal{S}$ 

$$p(s' | s, a) =$$

$$\begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\pi_5(\mathcal{S}) = \{1, 2, 2, 2, 3, 3\}$$

$$\overline{U}_4(\mathcal{S}) = \{3640, 3635, 3630, 3615, 3592, 3510\} \text{ M\$}$$

$$\pi_6(s) = \arg \max_{a \in \mathcal{A}} \left\{ \sum_{s' \in \mathcal{S}} p(s' | s, a_1) (r(s, a_1) + \gamma \overline{U}_5(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s' | s, a_2) (r(s, a_2) + \gamma \overline{U}_5(s')) \right. \\ \left. + \sum_{s' \in \mathcal{S}} p(s' | s, a_3) (r(s, a_3) + \gamma \overline{U}_5(s')) \right\}$$

$$\pi_5(\mathcal{S}) = \pi_6(\mathcal{S}) = \pi^*(\mathcal{S}) = \{1, 2, 2, 2, 3, 3\}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

## Policy Iteration

Policy eval. -  $\overline{U}^{(1)}(1) = 2063 \text{ M\$} \rightarrow \overline{U}^{(2)}(1) = 3483 \text{ M\$}$

$$\overline{U}^{(1)}(\mathcal{S}) = \{2063, 1290, 768, 455, 197, 0\} \text{ M\$}, \quad \pi^{(1)}(\mathcal{S}) = \{2, 2, 3, 3, 3, 3\}$$

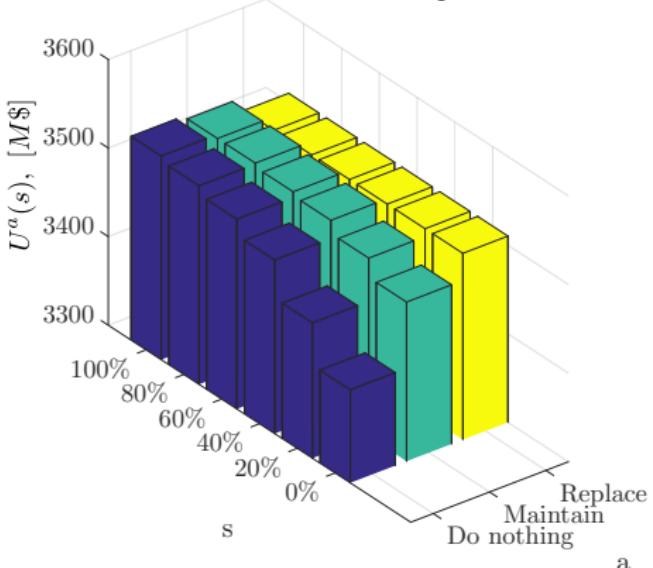
Policy Iteration loop  
 Policy eval. loop  
 State  
 $\overline{U}^{(2)(1)}(1) = \overbrace{1 \cdot (109.5 \text{ M\$} - 5 + 0.97 \cdot 2063 \text{ M\$})}^{\pi^{(1)}(1)=2, \text{ maintain}} = 2106 \text{ M\$}$   
 $\overline{U}^{(2)(2)}(1) = 1 \cdot (109.5 \text{ M\$} - 5 + 0.97 \cdot 2106 \text{ M\$}) = 2147 \text{ M\$}$   
 $\vdots$   
 $\overline{U}^{(2)(350)}(1) = 1 \cdot (109.5 \text{ M\$} - 5 + 0.97 \cdot 3483 \text{ M\$}) = 3483 \text{ M\$}.$   
 $\overline{U}^{(2)}(1) = \overline{U}^{(2)(351)}(1) = 1 \cdot (109.5 \text{ M\$} - 5 + 0.97 \cdot 3483 \text{ M\$}) = 3483 \text{ M\$}.$

Policy iteration -  $\pi^{(1)}(1) = 2 \rightarrow \pi^{(2)}(1) = 1$

$$\overline{\mathbb{U}}^{(2)}(\mathcal{S}) = \{3483, 3483, 3468, 3457, 3441, 3359\} \text{ M\$}$$

$$\begin{aligned}\pi^{(2)}(1) &= \arg \max_{a \in \mathcal{A}} \left\{ \begin{array}{l} \sum_{s' \in \mathcal{S}} p(s'|s, a_1) \cdot (r(s, a_1) + \gamma \overline{\mathbb{U}}_2(s')) \\ \sum_{s' \in \mathcal{S}} p(s'|s, a_2) \cdot (r(s, a_2) + \gamma \overline{\mathbb{U}}_2(s')) \\ \sum_{s' \in \mathcal{S}} p(s'|s, a_3) \cdot (r(s, a_3) + \gamma \overline{\mathbb{U}}_2(s')) \end{array} \right\} \\ &= \arg \max_{a \in \mathcal{A}} \left\{ \begin{array}{l} 109.5 + 0 \text{ M\$} + 0.97(0.95 \cdot 3483 \text{ M\$} + 0.03 \cdot 3483 \text{ M\$} + 0.02 \cdot 3468 \text{ M\$}) = \boxed{3488 \text{ M\$}} \\ 109.5 - 5 \text{ M\$} + 0.97 \cdot (1 \cdot 3483 \text{ M\$}) \\ 109.5 - 20 \text{ M\$} + 0.97 \cdot (1 \cdot 3483 \text{ M\$}) \end{array} \right\} \\ &= 3483 \text{ M\$} \\ &= 3468 \text{ M\$} \\ &= 1.\end{aligned}$$

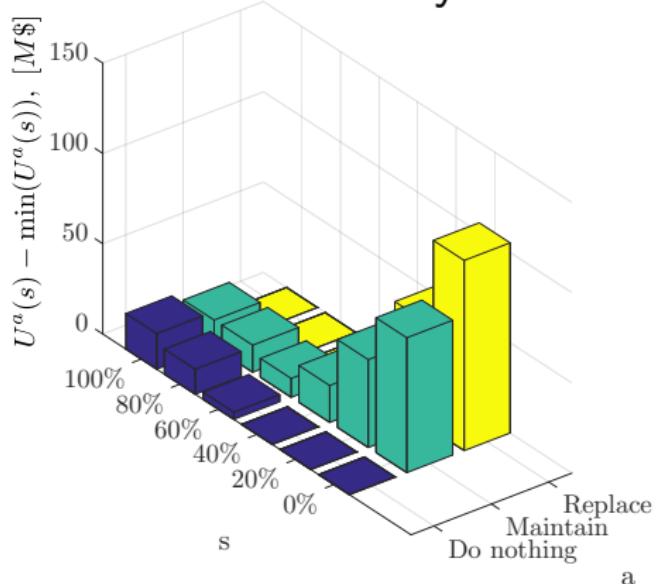
# Infrastructure Maintenance - Policy iteration [m]



$$\pi^*(\mathcal{S}) = \underbrace{\text{Do nothing}}_{s=100\%}, \underbrace{\text{Maintain}}_{s=80\%}, \underbrace{\text{Maintain}}_{s=60\%}, \underbrace{\text{Maintain}}_{s=40\%}, \underbrace{\text{Replace}}_{s=20\%}, \underbrace{\text{Replace}}_{s=0\%}$$

[CIV\_ML/Policy\_iteration\_bridge.m]

# Infrastructure Maintenance - Policy iteration [m]



$$\pi^*(S) = \underbrace{\text{Do nothing}}_{s=100\%}, \underbrace{\text{Maintain}}_{s=80\%}, \underbrace{\text{Maintain}}_{s=60\%}, \underbrace{\text{Maintain}}_{s=40\%}, \underbrace{\text{Replace}}_{s=20\%}, \underbrace{\text{Replace}}_{s=0\%}$$

[CIV\_ML/Policy\_iteration\_bridge.m]

# Policy v.s. Value iteration

$\pi$ : Policy iteration  $\equiv \pi$ : Value iteration

Time: Policy iteration  $\ll$  Time: Value iteration

# Partially Observable MDP (POMDP)

One limitation of MDP is that at a time  $t$ , you must know exactly in what state (s) you are.

When this is not the case → **POMDP**

(Note: computationally more demanding than MDP)



Chapter 17 – Russell, S. and Norvig, P. (1995).

*Artificial Intelligence, A modern approach.* Prentice-Hall.

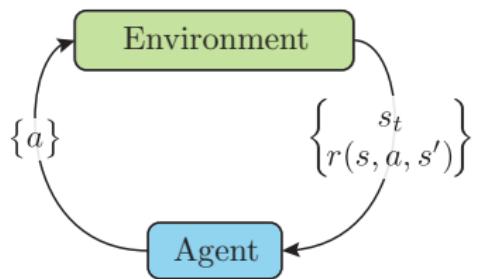
# Model-based v.s. Model-free

When we can explicitly represent  $p(s_{t+1}|s_t, a)$

→ **Model-based RL: MDP/POMDP, etc...**

When we cannot represent  $p(s_{t+1}|s_t, a)$  explicitly

→ **Model-free RL: Q-learning, SARSA, etc...**



Interaction between the **agent** and the **environment** over multiple **episodes**



# Section Outline

---

## Q-Learning

- 4.1 Context
  - 4.2 Temporal difference learning
  - 4.3 Q-learning
  - 4.4 Example – Infrastructure Maintenance
  - 4.5 Function Approximation
-

# Expected Utility

**Utility**  $\leftarrow$  sum of discounted reward

$$\mathbb{U}(s, \pi) = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}), \quad a_t = \pi(s_t)$$

**Expected Utility**

$$\mathbb{E}[\mathbb{U}(s, \pi)] \equiv \overline{\mathbb{U}}(s, \pi) \equiv \overline{\mathbb{U}}(s)$$

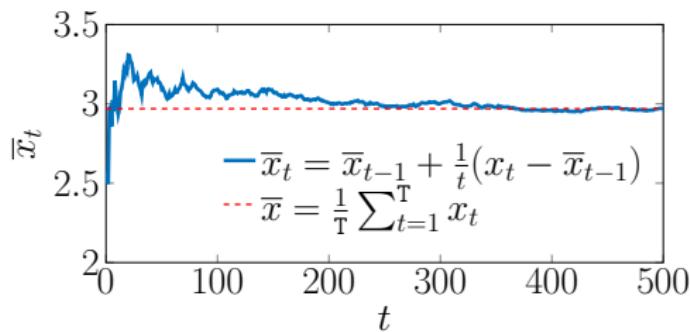
$\overline{\mathbb{U}}(s, \pi) \leftarrow$  average over realizations of  $\mathbb{U}(s, \pi)$

# Running average

The average  $\bar{x}_T$  of a set  $\{x_1, x_2, \dots, x_T\}$  can be calculated incrementally by following

$$\bar{x}_t = \bar{x}_{t-1} + \frac{1}{t}(x_t - \bar{x}_{t-1})$$

e.g.: a set of  $T = 500$   
realizations of  
 $x_t : X \sim \mathcal{N}(x; 3, 1)$



## Temporal difference learning

## Temporal difference (TD) learning

## Running average

$$\bar{x}_t = \bar{x}_{t-1} + \frac{1}{t}(x_t - \bar{x}_{t-1})$$

## Expected Utility

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \overline{\mathbb{U}}^{(i)}(s) + \frac{1}{N(s)} (\mathbb{U}^{(i)}(s) - \overline{\mathbb{U}}^{(i)}(s))$$

## TD learning

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \overline{\mathbb{U}}^{(i)}(s) + \alpha \underbrace{\left( r(s, a, s') + \gamma \overline{\mathbb{U}}^{(i)}(s') - \overline{\mathbb{U}}^{(i)}(s) \right)}_{\text{TD-target}}$$

Learning rate –  $\alpha$ 

$$\left. \begin{array}{l} \sum_{\substack{N(s,a)=1 \\ \infty}}^{\infty} \alpha(N(s,a)) = \infty \\ \sum_{\substack{N(s,a)=1 \\ \infty}}^{\infty} \alpha^2(N(s,a)) < \infty \end{array} \right\} \text{e.g. } \alpha(N(s,a)) = \frac{c}{c + N(s,a)}, c \in \mathbb{R}^+$$

# Action-utility function

$\mathbb{Q}(s, a)$ : Action-utility function such that

$$\overline{\mathbb{U}}(s) \equiv \overline{\mathbb{U}}(s, \pi^*) = \max_{a \in \mathcal{A}} \mathbb{Q}(s, a)$$

$$\pi^*(s) = \arg \max_{a \in \mathcal{A}} \mathbb{Q}(s, a)$$

## Temporal difference Q-learning (off-policy)

$$\mathbb{Q}^{(i+1)}(s, a) \leftarrow \mathbb{Q}^{(i)}(s, a) + \alpha(r(s, a, s') + \gamma \max_{a' \in \mathcal{A}} \mathbb{Q}^{(i)}(s', a') - \mathbb{Q}^{(i)}(s, a))$$

# Exploration-exploitation tradeoff

When transiting from state  $s$  to state  $s'$ :

- you cannot always rely on the current  $\pi^*(s) \rightarrow \text{local maximum}$
- you cannot always act randomly  $\rightarrow \text{poor performance}$

An  **$\epsilon$ -greedy agent** selects an action following

$$a = \begin{cases} a_i : p(a) = \frac{1}{A}, \forall a, & \text{with } \Pr(\epsilon(N(s))) \quad \text{Random action} \\ \arg \max_{a \in \mathcal{A}} Q(s, a), & \text{with } 1 - \Pr(\epsilon(N(s))) \quad \text{Greedy action,} \end{cases}$$

where  $\epsilon(N(s)) \rightarrow 0$  as  $N(s) \rightarrow \infty$

---

### Algorithm 3: Temporal difference Q-learning

---

```

1 define:  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ ,  $r(s, a, s')$ ,  $\gamma$ ,  $\epsilon$ ,  $\alpha$ 
2 initialize:  $\mathbb{Q}^{(0)}(s, a)$ ,  $N(s, a) = 0$ ,  $\forall \{s \in \mathcal{S}, a \in \mathcal{A}\}$ ,  $i = 0$ 
3 for episodes  $e \in \{1 : E\}$  do
4   initialize:  $s$ 
5   for time  $t \in \{1 : T\}$  do
6      $u : U \sim \mathcal{U}(0, 1)$ 
7      $a = \begin{cases} a_i : p(a) = \frac{1}{A}, \forall a, & \text{if } u < \epsilon(N(s)) \\ \arg \max_{a \in \mathcal{A}} \mathbb{Q}^{(i)}(s, a), & \text{if } u \geq \epsilon(N(s)) \end{cases}$  Random Greedy
8     observe:  $r(s, a, s')$ ,  $s'$ 
9      $\mathbb{Q}^{(i+1)}(s, a) \leftarrow \mathbb{Q}^{(i)}(s, a) + \alpha(N(s, a)) (r(s, a, s') + \gamma \max_{a' \in \mathcal{A}} \mathbb{Q}^{(i)}(s', a') - \mathbb{Q}^{(i)}(s, a))$ 
10     $s = s'$ ,  $N(s, a) = N(s, a) + 1$ ,  $i = i + 1$ 
11 return:  $\overline{\mathbb{Q}}(s) = \max_{a \in \mathcal{A}} \mathbb{Q}^{(i)}(s, a)$ 
12       $\pi^*(s) = \arg \max_{a \in \mathcal{A}} \mathbb{Q}^{(i)}(s, a)$ 

```

---

## Example – Infrastructure Maintenance

## Infrastructure Maintenance – Episodes

Learning rate:  $\alpha(N(s, a)) = \frac{70}{70+N(s,a)}$

$\epsilon$ -Greedy:  $\epsilon(N(s)) = \frac{70}{70+N(s)}$

$Q^{(0)}(s, a) = \$0 M, \forall \{s \in \mathcal{S}, a \in \mathcal{A}\}$

**Environment:**   $\equiv p(s'|s, a)$

500 episodes, each comprising 100 steps

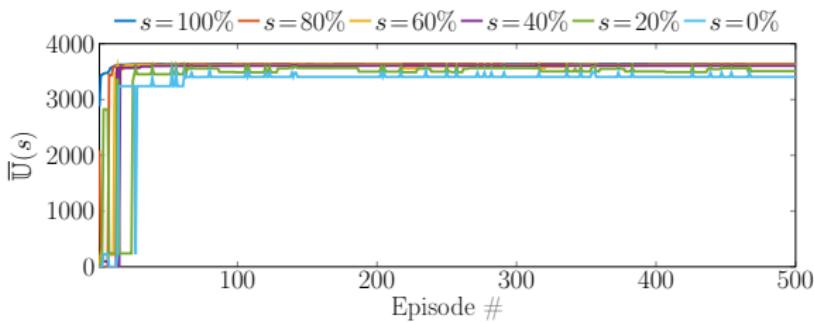
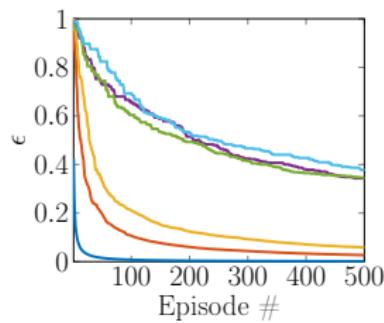
$$p(s'|s, a) = \begin{bmatrix} 0.95 & 0.03 & 0.02 & 0 & 0 & 0 \\ 0 & 0.9 & 0.05 & 0.03 & 0.02 & 0 \\ 0 & 0 & 0.8 & 0.12 & 0.05 & 0.03 \\ 0 & 0 & 0 & 0.7 & 0.25 & 0.05 \\ 0 & 0 & 0 & 0 & 0.6 & 0.4 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

$$Q^{(0)}(s, a) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

$$Q^{(0)}(s, a) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

## Example – Infrastructure Maintenance

## Infrastructure Maintenance – Q-Learning



$$\bar{U}^{(50K)}(\mathcal{S}) = \{3640, 3635, 3630, 3615, 3592, 3510\} \$M.$$

$$\pi^*(\mathcal{S}) = \underbrace{\text{Do nothing}}_{s=100\%}, \underbrace{\text{Maintain}}_{s=80\%}, \underbrace{\text{Maintain}}_{s=60\%}, \underbrace{\text{Maintain}}_{s=40\%}, \underbrace{\text{Replace}}_{s=20\%}, \underbrace{\text{Replace}}_{s=0\%}$$

# Q-learning with function approximation



$$\begin{aligned}\mathcal{A} &= \{\text{[radio]}\} \\ \mathcal{S} &= \{\text{[block]}\} \\ r(\mathcal{S}) &= \{\text{fct}(\text{[block]})\} \\ p(s_{t+1}|s_t, a) &= \text{unknown...} \\ \overline{U}(s) &\rightarrow Q(s, a) \\ &\quad (\text{Neural Network})\end{aligned}$$



Chapter 21 – Russell, S. and Norvig, P. (1995)  
*Artificial Intelligence, A modern approach*. Prentice-Hall.



Sutton, R. S. and Barto, A. G. (2018)  
*Reinforcement learning: An introduction*, MIT Press, 2nd edition

[Google images]

## Summary

**AI** : reproduction of an intelligent behavior by a machine.  
→ perception and dynamic interaction with its environment through the process of **making rational decisions**.

States:  $s \in \mathbb{R}$  or  $s_j \in \mathcal{S} = \{s_1, s_2, \dots, s_S\}$

Actions:  $a_i \in \mathcal{A} = \{a_1, a_2, \dots, a_A\}$

$$\text{Model: } p(s' | s_j, a_j) = \Pr(s' = s_k | s_j, a_j), \forall k = 1 : S$$

Reward:  $r(s, a, s') \in \mathbb{R}$

Policy:  $\pi = \{a(s_1), a(s_2), \dots, a(s_S)\}$

When we can explicitly represent  $p(s_{t+1}|s_t, a)$   
→ Multi-Head RL-MDP/PGMDP

## → Model-based RL: MDP/POMDP, etc...

When we cannot represent  $p(s_{t+1}|s_t, a)$  explicitly

## → Model-free RL: Q-learning, SARSA, etc...

### Infinite planning horizon:

The rewards and utilities are considered over an infinite period of time.

The optimal policy is **stationary**;  $\pi^*(s)$  does not depend on  $t$ .

$$\mathbb{Q}^{(i+1)}(s, a) \leftarrow \mathbb{Q}^{(i)}(s, a) + \alpha(\mathbb{N}(s, a)) (r(s, a, s') + \gamma \max_{a' \in \mathcal{A}} \mathbb{Q}^{(i)}(s', a')) - \mathbb{Q}^{(i)}(s, a))$$

**Bellman equation:** "The utility of a state is the immediate reward for that state plus the expected discounted utility of the next state given that the optimal action is taken"

$$\overline{U}(s) = \max_{a \in \mathcal{A}} \sum_{s' \in S} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{U}(s'))$$

## Bellman update & Value iteration

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \max_{a \in \mathcal{A}} \sum_{s' \in S} p(s'|s, a) \cdot (r(s, a, s') + \gamma \overline{\mathbb{U}}^{(i)}(s'))$$

## Policy iteration

$$\overline{\mathbb{U}}^{(i+1)}(s) \leftarrow \sum_{s' \in S} p(s'|s, \pi_i(s)) \cdot (r(s, \pi_i(s), s') + \gamma \overline{\mathbb{U}}^{(i)}(s'))$$

$\pi$ : Policy iteration  $\equiv \pi$ : Value iteration  
 Time: Policy iteration  $\ll$  Time: Value iteration

## TD Q-learning: