

Introduction to Multi-Agent Programming

11. Learning in Multi-Agent Systems (Part A)

SDP, MDPs, Value Iteration, Policy Iteration, RL

Alexander Kleiner, Bernhard Nebel

Contents

- Introduction
- Sequential decision problems
- Markov decision processes
 - Value Iteration & Policy Iteration
 - Reinforcement Learning (RL)

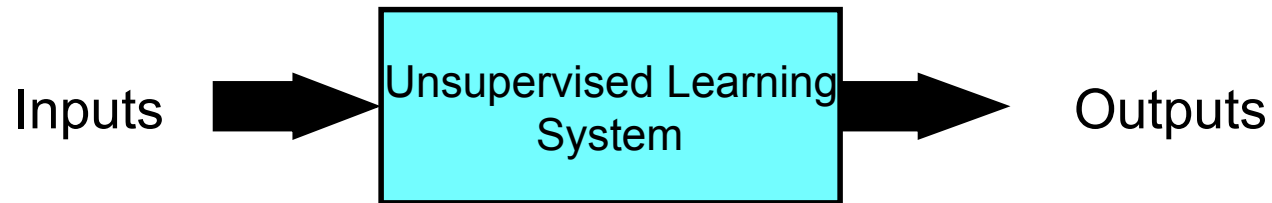
Introduction

- The importance of learning in MAS:
 - Agents are typically deployed in complex domains, i.e., dynamic domains with large state spaces, and uncertainty of action execution
 - Sometimes impossible to prepare agents for **any** situation
- Learning methods can be used to
 - enable the agent to do rich decisions based on little experience (**generalization**)
 - enable the agent to change its behavior online according to changes in the world (**adaption**)
- However, machine learning suffers under the “**curse of dimensionality**”
 - Exponential growth of the state space with an increasing number of state variables
 - Exponential growth of action space with an increasing number of action (In MAS even harder)

Different Types Of Learning feedback

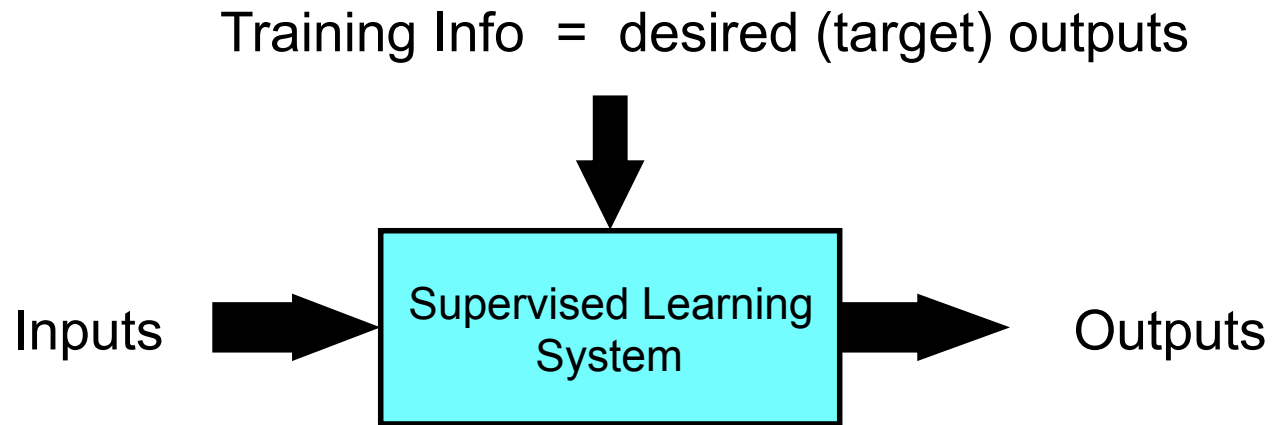
- The learning **feedback** indicates the performance level achieved so far
- The following learning feedbacks are distinguished:
 - **Supervised** learning (teacher)
 - **Reinforcement** learning (critic)
 - **Unsupervised** learning (observer)

Unsupervised Learning



Example: clustering of texts on the Internet according to counted word frequencies

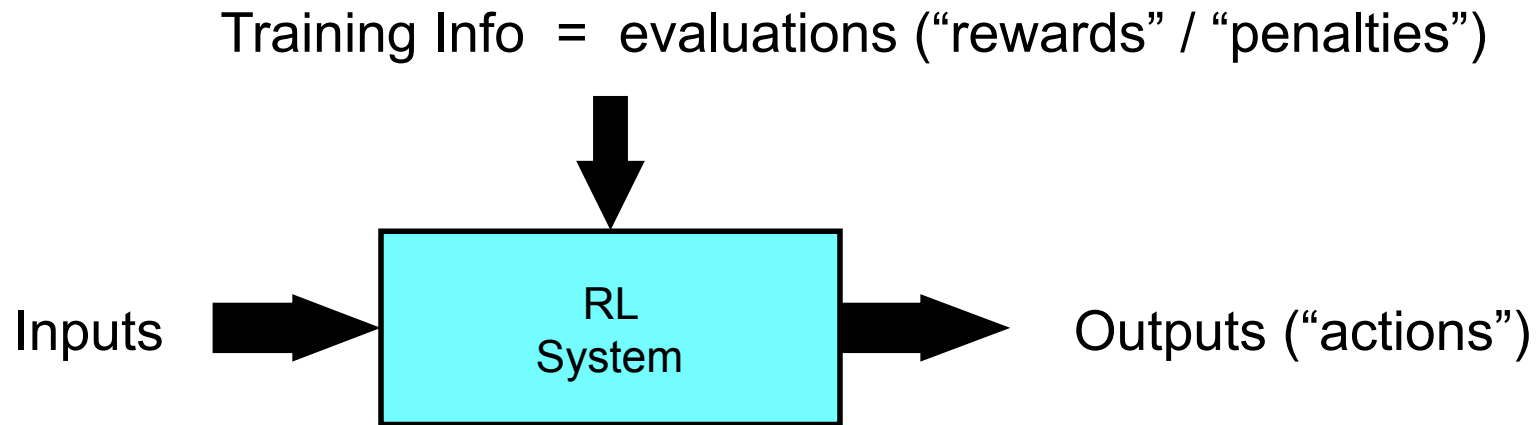
Supervised Learning



Error = (target output – actual output)

Example: detecting faces in images

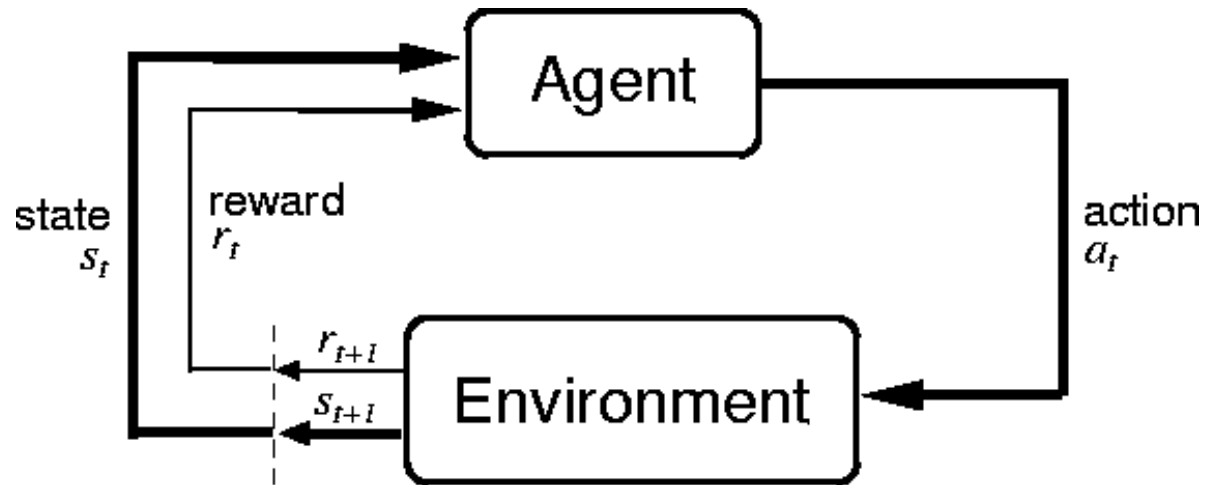
Reinforcement Learning



Objective: get as much reward as possible

Example: robot driving without collisions

The Agent-Environment Interface



Agent and environment interact at discrete time steps: $t = 0, 1, 2, \dots$

Agent observes state at step t : $s_t \in S$

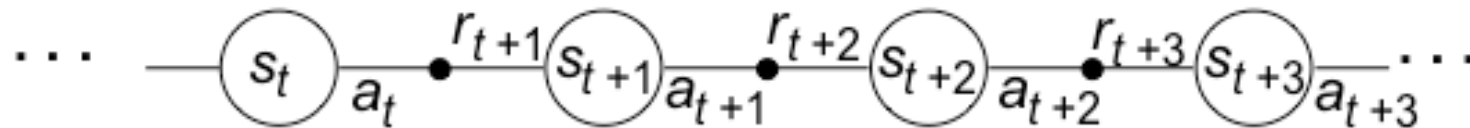
produces action at step t : $a_t \in A(s_t)$

gets resulting reward: $r_{t+1} \in \mathfrak{R}$

and resulting next state: s_{t+1}

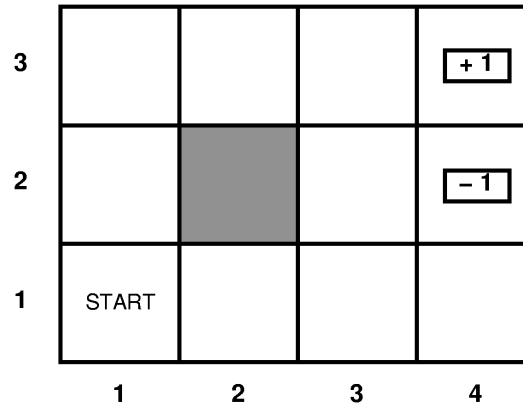
The Credit-Assignment Problem

- The problem of properly **assigning feedback** for an overall performance change to each of the system activities that **contributed** to that change



- Which actions were **invariant**, which were **important**?
- Can be **decomposed** into two sub-problems:
 - The inter-agent CAP
 - Assignment of credit for an overall performance change to the external actions of the agents
 - The intra-agent CAP
 - Assignment of credit for a particular external action of an agent to its **internal** modules

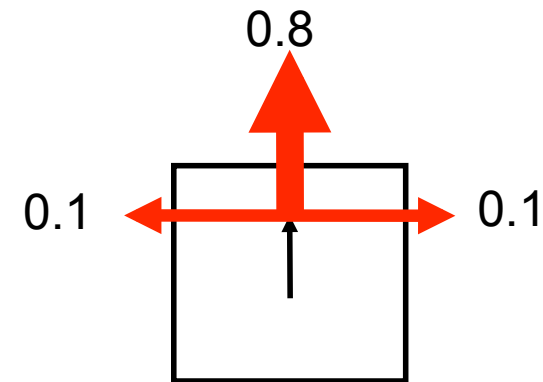
Sequential Decision Problems (1)



- Beginning in the start state the agent must choose an action at each time step.
- The interaction with the environment terminates if the agent reaches one of the goal states (4, 3) (reward of +1) or (4,2) (reward -1). Each other location has a reward of -.04.
- In each location the available actions are Up, Down, Left, Right.

Sequential Decision Problems (2)

- **Deterministic version:** All actions always lead to the next square in the selected direction, except that moving into a wall results in no change in position.
- **Stochastic version:** Each action achieves the intended effect with probability 0.8, but the rest of the time, the agent moves at right angles to the intended direction.



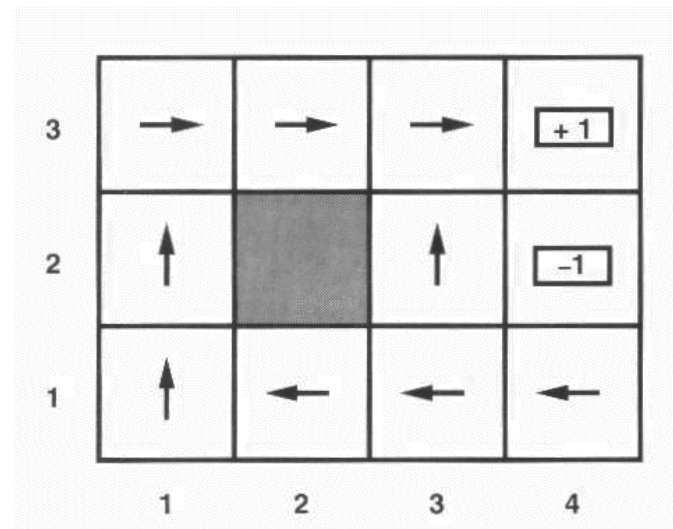
Markov Decision Problem (MDP)

- Given a set of actions A , a set of states S in an accessible, stochastic environment, an **MDP** is defined by
 - Initial state S_0
 - Transition Model $T(s,a,s')$
 - Reward function $R(s)$
- **Transition model:** $T(s,a,s')$ is the probability that state s' is reached, if action a is executed in state s .
- **Policy:** Complete mapping π that specifies for each state s which action $\pi(s)$ to take.
- **Wanted:** The ***optimal policy*** π^* is the policy that maximizes the expected utility.

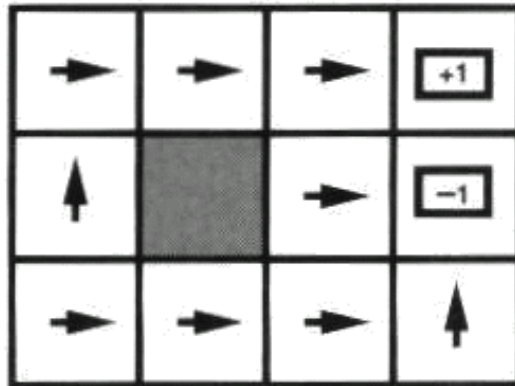
Optimal Policies (1)

- Given the optimal policy, the agent uses its **current percept** that **tells** it its **current state**.
- It then **executes** the **action** $\pi^*(s)$.
- We obtain a simple reflex agent that is computed from the information used for a utility-based agent.

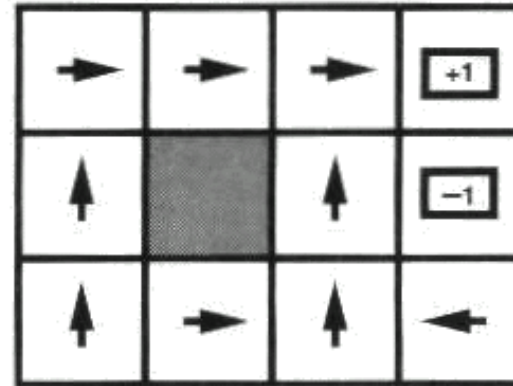
Optimal policy for our MDP
when $R(s) = -0.4$ for non-
terminals:



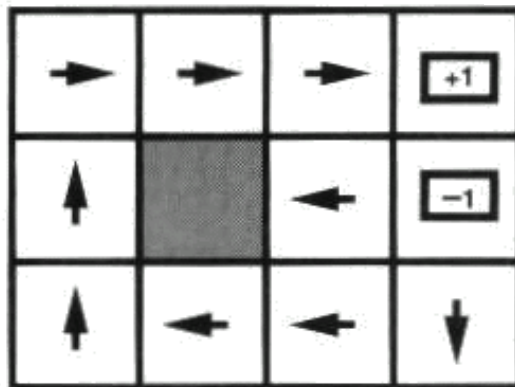
Optimal Policies (2)



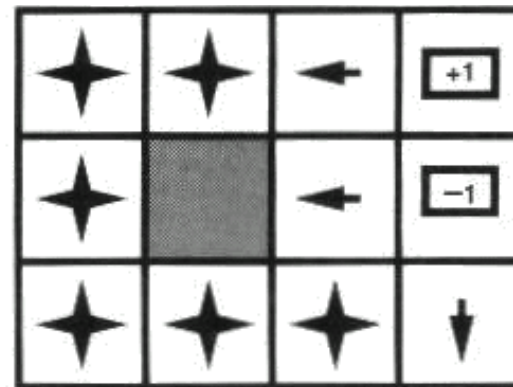
$$R(s) \leq -1.6248$$



$$-0.4278 < R(s) < -0.085$$



$$-0.0221 < R(s) < 0$$



$$0 < R(s)$$

How to compute optimal policies?

Finite and Infinite Horizon Problems

- Performance of the agent is measured by the sum of rewards for the states visited.
- To determine an optimal policy we will first calculate the utility of each state and then use the state utilities to select the optimal action for each state.
- The result depends on whether we have a finite or infinite horizon problem.
- Utility function for state sequences: $U_h([s_0, s_1, \dots, s_n])$
- Finite horizon: $U_h([s_0, s_1, \dots, s_{N+k}]) = U_h([s_0, s_1, \dots, s_N])$ for all $k > 0$.
- For **finite horizon** problems the optimal policy depends on the horizon N .
- In **infinite horizon** problems the optimal policy only depends on the current state.

Assigning Utilities to State Sequences

- For **finite horizon** problems utilities for each state can be computed by summing-up rewards of each state:
 - $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + R(s_1) + R(s_2) + \dots$
- For **infinite horizon** problems utilities have to be computed by discounting future rewards:
 - $U_h([s_0, s_1, s_2, \dots]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$
- The term $\gamma \in [0:1[$ is called the **discount factor**.
- With **discounted rewards** the utility of an infinite state sequence is always finite. The discount factor expresses that future rewards have less value than current rewards.

Utilities of States

- The utility of a state depends on the utility of the state sequences that follow it.
- Let $U^\pi(s)$ be the utility of a state under policy π .
- Let s_t be the state of the agent after executing π for t steps. Thus, the utility of s under π is

$$U^\pi(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid \pi, s_0 = s \right]$$

- The true utility $U(s)$ of a state is $U^{\pi^*}(s)$.
- $R(s)$ is the short-term reward for being in s and $U(s)$ is the long-term total reward from s onwards.

Choosing Actions using the Maximum Expected Utility Principle

The agent simply chooses the action that maximizes the expected utility of the subsequent state:

$$\pi(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U(s')$$

The utility of a state is the immediate reward for that state plus the expected discounted utility of the next state, assuming that the agent chooses the optimal action:

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

Example

The utilities of the states in our 4x3 world with $\gamma=1$ and $R(s)=-0.04$ for non-terminal states:

3	0.812	0.868	0.918	+1
2	0.762		0.660	-1
1	0.705	0.655	0.611	0.388
	1	2	3	4

Which action would an optimal agent choose here?

Bellman-Equation

- The equation

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

is also called the **Bellman-Equation**.

- In our 4x3 world the equation for the state (1,1) is

$$U(1,1) = -0.04 + \gamma \max\left\{ \begin{array}{ll} 0.8 U(1,2) + 0.1 U(2,1) + 0.1 U(1,1), & (Up) \\ 0.9 U(1,1) + 0.1 U(1,2), & (Left) \\ 0.9 U(1,1) + 0.1 U(2,1), & (Down) \\ 0.8 U(2,1) + 0.1 U(1,2) + 0.1 U(1,1) \} & (Right) \end{array} \right.$$

→ Given the numbers for the optimal policy, Up is the optimal action in (1,1).

Value Iteration (1)

An algorithm to calculate an optimal strategy.

Basic Idea: Calculate the utility of each state. Then use the state utilities to select an optimal action for each state.

How to calculate the utility of each state?

The bellman equation can be used to build as system of n equations for n states

However, due to the transition model and the therefore required max operator, the system is non-linear

→ Solution can not be computed in closed form (can only be done for deterministic problems)

Value Iteration (2)

Iterative Procedure

Solution:

We can apply an **iterative approach** in which we replace the **equality** of the bellman equation by an **assignment**:

$$U(s) \leftarrow R(s) + \gamma \max_a \sum_{s'} T(s, a, s') U(s')$$

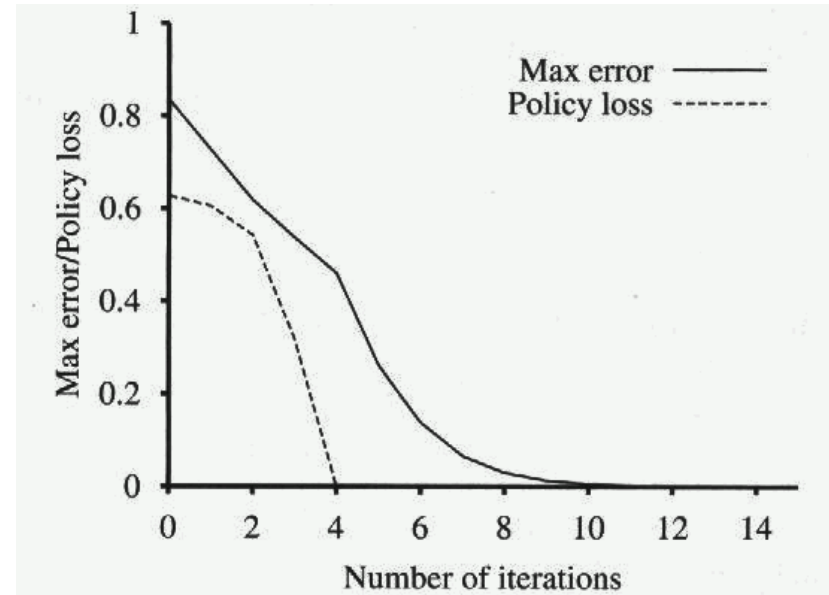
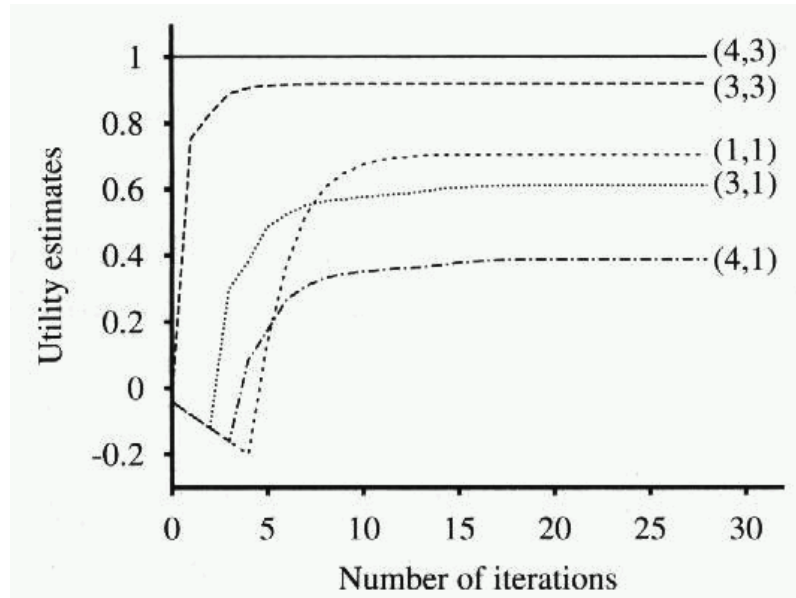
The Value Iteration Algorithm

```
function VALUE-ITERATION(mdp,  $\epsilon$ ) returns a utility function
  inputs: mdp, an MDP with states  $S$ , transition model  $T$ , reward function  $R$ , discount  $\gamma$ 
            $\epsilon$ , the maximum error allowed in the utility of any state
  local variables:  $U$ ,  $U'$ , vectors of utilities for states in  $S$ , initially zero
                     $\delta$ , the maximum change in the utility of any state in an iteration

  repeat
     $U \leftarrow U'$ ;  $\delta \leftarrow 0$ 
    for each state  $s$  in  $S$  do
       $U'[s] \leftarrow R[s] + \gamma \max_a \sum_{s'} T(s, a, s') U[s']$ 
      if  $|U'[s] - U[s]| > \delta$  then  $\delta \leftarrow |U'[s] - U[s]|$ 
  until  $\delta < \epsilon(1 - \gamma)/\gamma$ 
  return  $U$ 
```

It can be shown that value iteration converges

Application Example



In practice the policy often becomes optimal before the utility has converged.

Policy Iteration

- Value iteration computes the optimal policy even at a stage when the utility function estimate has not yet converged.
- If one action is better than all others, then the exact values of the states involved need not to be known.
- Policy iteration alternates the following two steps beginning with an initial policy π_0 :
 - **Policy evaluation**: given a policy π_t , calculate $U_t = U^{\pi_t}$, the utility of each state if π_t were executed.
 - **Policy improvement**: calculate a new maximum expected utility policy π_{t+1} according to

$$\pi_{t+1}(s) = \operatorname{argmax}_a \sum_{s'} T(s, a, s') U(s')$$

The Policy Iteration Algorithm

```
function POLICY-ITERATION(mdp) returns a policy
  inputs: mdp, an MDP with states S, transition model T
  local variables: U, U', vectors of utilities for states in S, initially zero
                    $\pi$ , a policy vector indexed by state, initially random

  repeat
     $U \leftarrow$  POLICY-EVALUATION( $\pi$ , U, mdp)
    unchanged?  $\leftarrow$  true
    for each state s in S do
      if  $\max_a \sum_{s'} T(s, a, s') U[s'] > \sum_{s'} T(s, \pi[s], s') U[s']$  then
         $\pi[s] \leftarrow \operatorname{argmax}_a \sum_{s'} T(s, a, s') U[s']$ 
        unchanged?  $\leftarrow$  false
  until unchanged?
  return P
```

Reinforcement Learning

- Learning from **interaction** with an external environment or other agents
- Goal-oriented learning
- Learning and making observations are **interleaved**
- Process is modeled as MDP or variants

Key Features of RL

- Learner is not told which **actions** to take
- Possibility of **delayed reward** (sacrifice short-term gains for greater long-term gains)
- **Model-free**: Models are learned online, i.e., have not to be defined in advance!
- Trial-and-Error search
- The need to **explore** and **exploit**

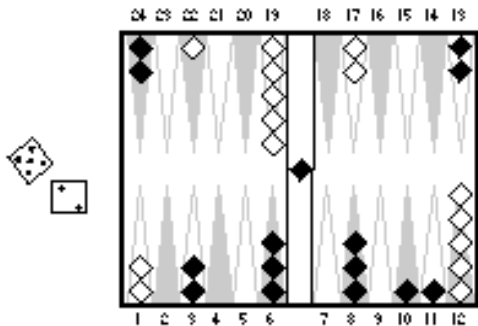
Some Notable RL Applications

- TD-Gammon: Tesauro
 - world's best backgammon program
- Elevator Control: Crites & Barto
 - high performance down-peak elevator controller
- Dynamic Channel Assignment: Singh & Bertsekas, Nie & Haykin
 - high performance assignment of radio channels to mobile telephone calls
- ...

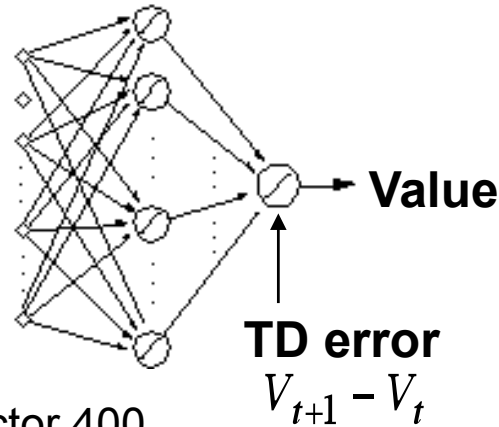
Some Notable RL Applications

TD-Gammon

Tesauro, 1992–1995



Effective branching factor 400



Action selection
by 2–3 ply search

Start with a random network

Play very many games against self

Learn a value function from this simulated experience

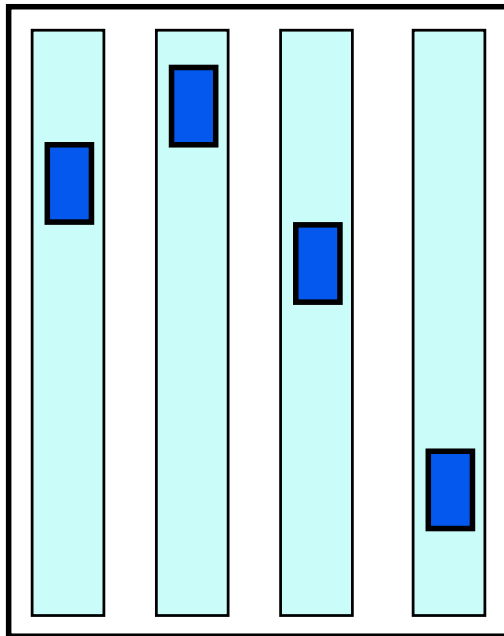
This produces arguably the best player in the world

Some Notable RL Applications

Elevator Dispatching

Crites and Barto, 1996

10 floors, 4 elevator cars



STATES: button states; positions, directions, and motion states of cars; passengers in cars & in halls

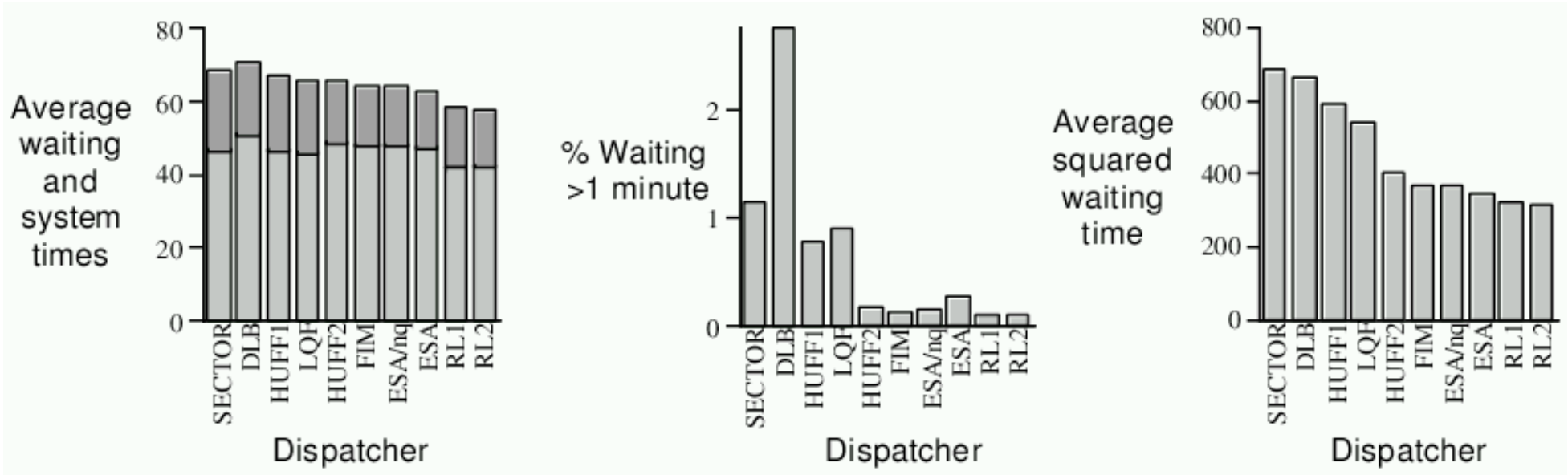
ACTIONS: stop at, or go by, next floor

REWARDS: roughly, -1 per time step for each person waiting

Conservatively about 10^{22} states

Some Notable RL Applications

Performance Comparison Elevator Dispatching



Q-Learning (1)

- Very **common** Reinforcement Learning method
- Maintains a table of **Q-values**
 - $Q(s,a)$ – “what is the outcome of action a in state s ”?
- Since values are with respect to states **and** actions, no explicit transition model T needed
- Updates are performed with a step size parameter in order to prevent value **overwriting** during different traces
- Converges to the optimum Q-values with probability 1

Q-Learning (2)

- At time t the agent performs the following steps:
 - Observe the **current** state s_t
 - Select and **perform** action a_t
 - Observe the **subsequent** state s_{t+1}
 - Receive **immediate** payoff r_t
 - **Adjust** Q-value for state s_t

Q-Learning (3)

Update and Selection

- Update function:

$$Q_{k+1}(s_t, a_t) := (1 - \alpha) Q_k(s_t, a_t) + \alpha \left[R(s_t, a_t) + \gamma \max_{a \in A} Q_k(s_{t+1}, a_{t+1}) \right]$$

- Where k denotes the version of the Q function, and α denotes a learning **step size parameter** that should decay over time
- Intuitively, actions can be **selected** by:

$$\pi(s_t) = \underset{a \in A}{\operatorname{argmax}} Q(s_t, a)$$

Q-Learning (4)

Algorithm

Initialise $Q(s, a)$ arbitrary for all $s \in S$ and $a \in A$

Repeat

select best action a_t with the greedy policy:

$$a_t = \pi(s_t) = \underset{a \in A}{\operatorname{argmax}} Q(s_t, a)$$

apply a_t in the world and observe s_{t+1} and immediate reward r_t :

$$s_t \rightarrow s_{t+1}$$

$$r_t$$

adapt the value function for state s_t

$$Q_{k+1}(s_t, a_t) := (1 - \alpha) Q_k(s_t, a_t) + \alpha \left[r_t + \gamma \max_{a \in A} Q_k(s_{t+1}, a_{t+1}) \right]$$

Until ($Q_{k+1} - Q_k < \varepsilon$) or (s is terminal)

The Exploration/Exploitation Dilemma

- Suppose you form estimates

$$Q_t(a) = Q^*(a) \quad \text{action value estimates}$$

- The greedy action at time t is:

$$a_t^* = \operatorname{argmax}_a Q_t(a)$$

$$a_t = a_t^* \Rightarrow \text{exploitation}$$

$$a_t \neq a_t^* \Rightarrow \text{exploration}$$

- You can't exploit all the time; you can't explore all the time
- You can never stop exploring; but you should always reduce exploring

e-Greedy Action Selection

- Greedy action selection:

$$a_t = a_t^* = \arg \max_a Q_t(a)$$

- e-Greedy:

$$a_t = \begin{cases} a_t^* & \text{with probability } 1 - \varepsilon \\ \text{random action} & \text{with probability } \varepsilon \end{cases}$$

- Continuously decrease of ε during each episode necessary!

→ the simplest way to try to balance exploration and exploitation