Anna Helena Reali Costa Escola Politécnica, Universidade de São Paulo Computer Engineering Department

### TRANSFER OF KNOWLEDGE IN REINFORCEMENT LEARNING





# Suppose we <u>know</u> how to ride a tricycle



Suppose we want to <u>learn</u> to ride a bicycle





Workshop eScience - Anna Reali (POLI / USP)





### How would knowing to ride a tricycle help in learning to ride a bike?



Workshop eScience - Anna Reali (POLI / USP)



### How would knowing to ride a tricycle help in learning to ride a bike?











Can we design agents able to learn from experience and transfer knowledge across different problems to improve their learning performance?



### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions

### Outline

### Reinforcement Learning (RL)

- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions

### Learning

 Learning modifies the agent's decision mechanisms to improve performance





### Learning Agent

- Design of a learning agent depends on what feedback is available
- Type of feedback:
   Supervised learning: correct answers/output for each example/input (classification, regression)
   Unsupervised learning: no feedback is given (clustering)

Reinforcement learning: occasional feedback



## Learning Agent

- Design of a learning agent depends on what feedback is available
- Type of feedback:
  - Supervised learning: correct answers/output for each example/input (classification, regression) Unsupervised learning: no feedback is given (clustering)

Reinforcement learning: occasional feedback



# **Reinforcement Learning (RL)**

- Feedback is delayed, occasional
- Time really matters (sequential, non i.i.d data)
- Agent's actions affect the subsequent data it receives
- Trial and error learning (via experiences)
   Task: Learn from this delayed reward to choose sequences of actions that produce the greatest cumulative reward



### Sequential decision-making problems



<position, speed>





Workshop eScience - Anna Reali (POLI / USP)





### **Markov Decision Problem**

- A Markov Decision Process (MDP) is:
   -Set of states S
  - -Set of actions A
  - -Dynamics P(s'|s,a) probability of transition
  - -Reward r(s,a)
- Solution: a policy,  $\pi$ : S  $\rightarrow$  A, that maximizes

$$Q^{\pi}(s,a) = E\left[\sum_{t=0}^{\infty} \gamma^{t} r(s_{t}) \mid \pi\right]$$

### **RL** method

- Key idea: updating the utility value
   Q(s,a) using the experience sequences
  - A trade off when choosing action between:
  - its immediately good (Exploitation)

#### ×

its long term good (Exploration)

### **RL** methods

- Initialize Q(s,a) arbitrarily
- Observe the current state s
- Do until a stop condition is reached:
  - select an action **a** and execute it in **s**
  - receive immediate reward *r*
  - observe the new state  $\boldsymbol{s'}$
  - update value function Q(s,a) and policy  $\pi$
  - S ← S'

### E.g.: Robot in a room



### E.g.: Robot in a room



### But...

- RL often requires many samples
- It takes too long to learn...
- The solution can usually only be applied to one specific task in a fixed setting



### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions









### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions

### **RL** methods

- Initialize Q(s,a)
- Observe the current state s
- Do until a stop condition is reached:
  - select an action a and execute it in s
  - receive immediate reward r
  - observe the new state s'
  - update value function Q(s,a) and policy  $\boldsymbol{\pi}$

- s ← s'

### Heuristically Accelerated Learning – HAL

- **Proposal:** use heuristics in the choice of actions  $\pi(s) = \begin{cases} \arg \max_{a} \hat{Q}(s,a) + \xi H_t(s,a)^{\beta} & prob. = \varepsilon, \\ a_{random} & prob. = 1 - \varepsilon \end{cases}$ 
  - Maintain convergence guarantee
    Good heuristics: fast convergence
    With bad heuristics performance degrades,
    ... but the process recovers!!

### Heuristically Accelerated Learning – HAL

Proposal: use heuristics in the choice of

actions

$$\pi(s) = \begin{cases} \arg \max_{a} \hat{Q}(s,a) + \xi H_t(s,a)^{\beta} \quad prob. = \varepsilon, \\ a_{random} \quad prob. = 1 - \varepsilon \end{cases}$$



# Results in soccer game with 2 players, both learning



### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions

### **Generalizing experiences**

- The idea: improving the representation language
  - If we use a more powerful representation,
     we can generalize states and actions across
     tasks (and, therefore, generalize policies)



### **Representation: predicates**



### Ground state:

- inRoom(r<sub>1</sub>)
- seeDoor(d<sub>3</sub>)
- seeAdjCorridor(c<sub>1</sub>)

### **Representation: predicates**



### **Abstract state:**

- inRoom(X)
- seeDoor(Y)
- seeAdjCorridor(Z)
   Use of variables

1 abstract state comprises a set of ground states

**Relational MDP** 

### Abstraction $\rightarrow$ Generalization



- Enables state aggregation
- Reduced space
- Enables transfer learning

### Transfer

• Tasks described by the same predicates





### **Agent Architecture**

Simultaneous 2-layer reinforcement learning (S2L-RL)



### **Agent Architecture**

Simultaneous 2-layer reinforcement learning (S2L-RL)



### **Results** Navigation task



KOGA, FREIRE, COSTA. *IEEE Trans. on Cybernetics* 2015



### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Generalizing the Experience
- Initializing the value function
- Conclusions

### **RL algorithms**

- Initialize Q(s,a)
- Observe the current state s
- Do until a stop condition is reached:
  - select an action *a* and execute it in *s*
  - receive immediate reward *r*
  - observe the new state  $\boldsymbol{s'}$
  - update value function Q(s,a) and policy  $\pi$
  - S ← S'

# PITAM: Probabilistic Inter-TAsk Mappings

- Autonomously Define an Inter-Task Mapping
- Use of Object-Oriented MDPs
- Weight Q-value according to PITAM weight and initiate Q-table



• Classes: "Types" of objects

### **OO-MDP**

- Attributes: set of attributes composes a class
- Objects: Entities in the environment that belong to a class and have valuations for each attribute
   Object *id* Attributes



Mapping Prey to Gold, and Predator to Miner

Goldmine



Mapping Prey to Gold, and Predator to Miner

Goldmine



Mapping Prey to Gold, and Predator to Miner

Goldmine



Mapping Prey to Gold, and Predator to Miner

Goldmine



### Results

Source Task: Goldmine Domain (simpler) Target Task: Predator-Prey

Regular Learning: RL-learning from scratch QManualMapping: Hand-coded Mapping using Q-value reuse (Taylor et al., JMLR 2007)

PITAM: our proposal

Steps to conclude the Task:

	Regular	QManualMapping	PITAM
offset	76.21	51.22	30.48
generalization	45.67	45.92	29.48

SILVA, F.L.; COSTA, AHR. TIRL2017.



### Outline

- Reinforcement Learning (RL)
- Transfer in RL
- Improving the Exploration Strategy
- Spreading the Experience
- Initializing the value function
- Conclusions



### **Some Real-world Applications**

- Controlling Gene Regulatory Networks
   BRACIS 2016
- Energy Management Systems for Smart Homes
  - » IJCAI 2015
- Recommender Systems
  - » JBCS2015
- Coordination of EV Charging
  - On-going work

**>>** 

### **Future Work**

- Explore methods of approximation to handle realworld problems
- Explore languages to represent problems and solutions
- Automatic attribute discovery to better represent problems and transfer of knowledge
- Automatic discovery of similarity between tasks



### Thanks!!

Anna Helena Reali Costa Escola Politécnica Universidade de São Paulo anna.reali@usp.br









