#### **Interactively Shaping Agents via Human Reinforcement**

W. Bradley Knox 、 Peter Stone The University of Texas at Austin K-CAP-2009

## Introduction

- Background: As computational learning agents move into domains that incur real costs (e.g., autonomous driving or financial investment), it will be necessary to learn good policies without numerous high-cost learning trials.
- Goal: learn from human to speed learning

# Learning from Advice

Learning from Advice

- suggesting an action when a certain condition is true.
- [1]created a domain-specific natural language interface for giving advice to a reinforcement learner.
- general natural language recognition is unsolved
- Moreover, work still remains on how to embed advice into agents that learn from experience.

[1]G. Kuhlmann, P. Stone, R. Mooney, and J. Shavlik. Guiding a reinforcement learner with natural language advice: Initial results in RoboCup soccer

## Learning from demonstration

- Human provide trajectories
- Learn reward or policy

## Learning from Reinforcement

Provide human signal to shape the agent

[1] Their agent seeks to maximize the sum of human reinforcement and environmental reward.

A. Thomaz and C. Breazeal. Reinforcement Learning with Human Teachers: Evidence of Feedback and Guidance with Implications for Learning Performance

### Framework



# Algorithm

Algorithm 1	A general greedy TAMER algorithm		-51.
Require: Inp	ut: stepSize		-
1: ReinfModel.	init(stepSize)		
2: $\overrightarrow{s} \leftarrow \overrightarrow{0}$			
3: $\overrightarrow{f} \leftarrow \overrightarrow{0}$			
4: while true	do		
5: $h \leftarrow get H$	<pre>IumanReinfSincePreviousTimeStep()</pre>		
6: <b>if</b> $h \neq 0$	then		
7: error	$-h$ - $ReinfModel.predictReinf(\overrightarrow{f})$	柑刑百章	žf
8: ReinfA	$Iodel.update(\vec{f}, error)$	快生丈	F7
9: end if			
10: $\overrightarrow{s} \leftarrow get$	State Vec()		
11: $a \leftarrow argn$	$nax_a$ (ReinfModel. predict(getFeatures())	$\overrightarrow{s}, a)))$	
12: $\overrightarrow{f} \leftarrow get$	$Features(\vec{s}, a)$		动作选择
13: takeActio	n(a)		
14: wait for r	ext time step		
15: end while			_

## Credit assignment

$$c_{t} = P(event \ starting \ at \ t_{i} \ was \ targeted)$$
  
=  $P(target \ event \ between \ t_{i-1} \ and \ t_{i} \ sec. \ ago)$   
=  $\int_{t_{i-1}}^{t_{i}} f(x) dx$ 

$$\sum_{i=1}^{n} P(event \ starting \ at \ t_i \ was \ targeted) = 1$$

# Algorithm2

Require: Input: stepSize, windowSize,

Crediter.init(windowSize)		
$\overrightarrow{s} \leftarrow \overrightarrow{0}$		
$\vec{f} \leftarrow \vec{0}$		
$\overrightarrow{w} \leftarrow \overrightarrow{0}$		
while true do		
Crediter.updateTime(clockTime())		
$h \leftarrow getHumanReinfSincePreviousTimeStep()$	)	
if $h \neq 0$ then		
$\overrightarrow{credFeats} \leftarrow 0$		
for all $(\vec{f_t}, t) \in Crediter.historyWindow d$	do	
$c_t \leftarrow Crediter.assignCredit(t)$		ᆘᄴᅖᆘᆍᆠ
$\overrightarrow{credFeats} \leftarrow \overrightarrow{credFeats} + (c_t \times \overrightarrow{f_t})$		<b> </b>
end for		
$error \leftarrow h - (\ \overrightarrow{w} \cdot \overrightarrow{credFeats})$		
$\overrightarrow{w} \leftarrow \overrightarrow{w} + (stepSize \times error \times \overrightarrow{credFeats})$		
end if		
$\vec{s} \leftarrow getStateVec()$		
$a \leftarrow argmax_a ( \ \overrightarrow{w} \cdot (getFeatures(\overrightarrow{s}, a)))$		
$\overrightarrow{f} \leftarrow getFeatures(\overrightarrow{s}, a)$	— I	
takeAction(a)	动	作选择
$Crediter.updateWindow(\overrightarrow{f})$		
wait for next time step		
	$\begin{array}{l} \overrightarrow{s} \leftarrow \overrightarrow{0} \\ \overrightarrow{f} \leftarrow \overrightarrow{0} \\ \overrightarrow{w} \leftarrow \overrightarrow{0} \\ \hline \overrightarrow{m} \\ \overrightarrow{redFeats} \leftarrow 0 \\ \overrightarrow{fredFeats} \leftarrow \overrightarrow{0} \\ \overrightarrow{fredFeats} \leftarrow \overrightarrow{credFeats} + \overrightarrow{0} \\ \overrightarrow{fredFeats} \leftarrow \overrightarrow{credFeats} + (c_t \times \overrightarrow{f_t}) \\ \hline \overrightarrow{credFeats} \leftarrow \overrightarrow{credFeats} + (c_t \times \overrightarrow{f_t}) \\ \overrightarrow{end for} \\ \overrightarrow{error} \leftarrow h - (\overrightarrow{w} \cdot \overrightarrow{credFeats}) \\ \overrightarrow{w} \leftarrow \overrightarrow{w} + (stepSize \times error \times \overrightarrow{credFeats}) \\ \overrightarrow{w} \leftarrow \overrightarrow{w} + (stepSize \times error \times \overrightarrow{credFeats}) \\ \overrightarrow{f} \leftarrow getFeatures(\overrightarrow{s}, a) \\ \overrightarrow{f} \leftarrow getFeatures(\overrightarrow{s}, a) \\ \overrightarrow{f} \leftarrow getFeatures(\overrightarrow{s}, a) \\ \overrightarrow{f} \leftarrow argmax_a(\overrightarrow{w} \cdot (getFeatures(\overrightarrow{s}, a)))) \\ \overrightarrow{f} \leftarrow getFeatures(\overrightarrow{s}, a) \\ \overrightarrow{f} \leftarrow argmatewindow(\overrightarrow{f}) \\ \hline{wait for next time step} \end{array}$	$\begin{array}{llllllllllllllllllllllllllllllllllll$

23: end while

### Experiment-Tetris

 Table 1: Results of various Tetris agents.

Method	Mean Lines Cleared		Games
	at Game 3	at Peak	for Peak
TAMER	65.89	65.89	3
RRL-KBR [15]	5	50	120
Policy Iteration [2]	$\sim 0$ (no learning	3183	1500
[-]	until game 100)		
Genetic Algorithm [5]	$\sim$ 0 (no learning	586,103	3000
concere ingentinin [0]	until game 500)		
CE+BL [17]	$\sim 0$ (no learning	348,895	5000
	until game 100)		





Mountain Car

Mean Cumulative Reward in Mountain Car Cumulative Environmental Reward (Time to Goal) 0 -500 -1000 -1500 -2000 -2500 -3000 TAMER, Over 2nd and 3rd Runs Sarsa-3 meme -3500 Sarsa-20 ..... -4000 2 6 8 10 12 16 18 20 0 14 **Episode Number** 

## Conclusion

- works in the absence of an environmental reward function,
- reduces sample complexity
- is accessible to people who lack knowledge of computer science



#### 逆强化学习算法:最大熵逆强化学习



主动选初始状 态以及步长vs 随机选初始状 态和步长

步长相同 主动选初始状 态vs随机选初 始状态

X:轨迹数量 Y:reward loss



步长相同 主动选初始状 态vs随机选初 始状态

固定初始状态 vs固定选初始 状态