#### Learning Reward Machines for Partially Observable Reinforcement Learning

Rodrigo Toro Icarte Ethan Waldie Margarita P. Castro

Toryn Q. Klassen Richard Valenzano Sheila A. McIlraith



AAAI SSS-23

# What is a Reward Machine (RM)?



\*Using Reward Machines for High-Level Task Specification and Decomposition in Reinforcement Learning by Toro Icarte et al. (ICML, 2018)

RMs are automata-based reward functions:

RMs are automata-based reward functions:

```
m = 0 \# global variable
2 def get_reward(s):
    if m == 0 and s.at("A"):
4
      m
    if m == 1 and s.at("B"):
5
      m
        = 2
    if m == 2 and s.at("C"):
        = 3
8
      m
    if m == 3 and s.at("D"):
9
      m
        = 0
10
      return 1
11
    return 0
```



RMs are automata-based reward functions:





... that allow for learning policies faster.









... but the RMs were **handcrafted**.

**1** Shows how to learn RMs from experiences (LRM).

- **1** Shows how to learn RMs from experiences (LRM).
- 2 Uses RMs as memory for partially observable RL.

- **1** Shows how to learn RMs from experiences (LRM).
- 2 Uses RMs as memory for partially observable RL.
- 3 Extends QRM to work under partial observability.

- **1** Shows how to learn RMs from experiences (LRM).
- 2 Uses RMs as memory for partially observable RL.
- 3 Extends QRM to work under partial observability.
- 4 Provides a theoretical and empirical analysis of LRM.

# The Cookie Domain
































































(+1 Reward)

















Solving the cookie domain requires memory!



Solving the cookie domain requires **memory**!  $\pi^*(a|o_t) \ll \pi^*(a|o_0, \cdots, o_t)$ 

## Partially Observable RL

The most popular approach:

Training LSTMs policies using a policy gradient method.

The most popular approach:

Training LSTMs policies using a policy gradient method.

... starves in the cookie domain.



# RMs as memory

If the agent can detect the color of the rooms  $(\Box, \Box, \Box, \Box)$ ,

If the agent can detect the color of the rooms  $(\Box, \Box, \Box, \Box)$ , and when it presses the button  $(\bigcirc)$ ,

If the agent can detect the color of the rooms  $(\Box, \Box, \Box, \Box)$ , and when it presses the button  $(\bigcirc)$ , eats a cookie  $(\bigcirc)$ ,

If the agent can detect the color of the rooms  $(\Box, \Box, \Box, \Box)$ , and when it presses the button  $(\bigcirc)$ , eats a cookie  $(\bigcirc)$ , and sees a cookie  $(\bigcirc)$ ,

If the agent can detect the color of the rooms  $(\Box, \Box, \Box, \Box)$ , and when it presses the button  $(\bigcirc)$ , eats a cookie  $(\bigcirc)$ , and sees a cookie  $(\bigcirc)$ , then:



... becomes a "perfect" memory for the cookie domain.











conditio	ns at	state <i>u</i> 0
if (□ ○)	$\rightarrow$	goto u <sub>1</sub>
else	$\rightarrow$	<b>goto</b> u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if (□ ○)	$\rightarrow$	goto u <sub>1</sub>
else	$\rightarrow$	<b>goto</b> u <sub>0</sub>



conditions at state $u_0$		
if $(\square O)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditions at state $u_0$		
if $(\Box \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>











conditions at state $u_0$		
if (□ ●)	$\rightarrow$	goto u <sub>1</sub>
else	$\rightarrow$	goto u <sub>0</sub>



conditions at state $u_0$		
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
<b>if</b> (□ ○)	$\rightarrow$	goto u <sub>1</sub>
else	$\rightarrow$	goto u <sub>0</sub>






conditions a	t sta	te $u_1$
if (□ or □ �) if (□ or □ �) else	$\rightarrow$ $\rightarrow$ $\rightarrow$	<b>goto</b> <i>u</i> <sub>2</sub> <b>goto</b> <i>u</i> <sub>3</sub> <b>goto</b> <i>u</i> <sub>1</sub>







































conditio	ns at	state <i>u</i> 3
if (□ ☉)	$\rightarrow$	goto u <sub>0</sub>
else	$\rightarrow$	<b>goto</b> u <sub>3</sub>































conditio	ns at	state <i>u</i> 3
if $(\square \odot)$ else	ightarrow	goto u <sub>0</sub> goto u <sub>3</sub>



conditio	ns at	state <i>u</i> 3
if $(\square \odot)$ else	ightarrow	goto u <sub>0</sub> goto u <sub>3</sub>



conditior	ns at	state u <sub>3</sub>
if ( <b>□</b> ☉) else	ightarrow	goto $u_0$ goto $u_3$







condition	ns at	state <i>u</i> 0
if (□ ○)	$\rightarrow$	goto u <sub>1</sub>
else	$\rightarrow$	<b>goto</b> u <sub>0</sub>



condition	ns at	state <i>u</i> 0
if $(\square O)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if ( <b>□ ○</b> ) else	$\rightarrow$ $\rightarrow$	<b>goto</b> $u_1$
else	$\rightarrow$	<b>goto</b> u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditio	ns at	state <i>u</i> 0
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



conditions at state $u_0$		
if $(\square \bigcirc)$ else	ightarrow	goto u <sub>1</sub> goto u <sub>0</sub>



Why is this a perfect memory?



Why is this a perfect memory?  $\pi^*(a|o_0, \cdots, o_t) = \pi^*(a|o_t, u_t)$ 



Why is this a perfect memory?  $\pi^*(a|o_0, \cdots, o_t) = \pi^*(a|o_t, u_t)$ Hard problem  $\xrightarrow{\text{RM}}$  Easy problem

# How to learn such RMs?

# Learning Reward Machines

Given a set of detectors (e.g.,  $\{\Box, \Box, \Box, \Box, O, \odot, \odot\}$ ) and traces  $\mathcal{T}$ ,
# Learning Reward Machines

Given a set of detectors (e.g.,  $\{\Box, \Box, \Box, \Box, \Theta, \odot, \odot\}$ ) and traces  $\mathcal{T}$ , learning RMs is a **discrete optimization** problem:

$$\begin{array}{ll} \underset{U,u_{0},\delta_{u},\delta_{r}}{\text{minimize}} & \sum_{i \in I} \sum_{t \in T_{i}} \log(|N_{x_{i,t},L(e_{i,t})}|) & (\text{LRM}) \\ \text{s.t.} & \langle U, u_{0},\delta_{u},\delta_{r} \rangle \in \mathcal{R}_{\mathcal{P}} & (1) \\ & |U| \leq u_{\max} & (2) \\ & x_{i,t} \in U & \forall i \in I, t \in T_{i} \cup \{t_{i}\} & (3) \\ & x_{i,0} = u_{0} & \forall i \in I & (4) \\ & x_{i,t+1} = \delta_{u}(x_{i,t},L(e_{i,t+1})) & \forall i \in I, t \in T_{i} & (5) \\ & N_{u,l} \subseteq 2^{2^{\mathcal{P}}} & \forall u \in U, l \in 2^{\mathcal{P}} & (6) \\ & L(e_{i,t+1}) \in N_{x_{i,t},L(e_{i,t})} & \forall i \in I, t \in T_{i} & (7) \end{array}$$

# Learning Reward Machines

Given a set of detectors (e.g.,  $\{\Box, \Box, \Box, \Box, \Theta, \odot, \odot\}$ ) and traces  $\mathcal{T}$ , learning RMs is a **discrete optimization** problem:

$$\begin{array}{ll} \underset{U,u_{0},\delta_{u},\delta_{r}}{\text{minimize}} & \sum_{i \in I} \sum_{t \in T_{i}} \log(|N_{x_{i,t},L(e_{i,t})}|) & (\text{LRM}) \\ \text{s.t.} & \langle U, u_{0}, \delta_{u}, \delta_{r} \rangle \in \mathcal{R}_{\mathcal{P}} & (1) \\ & |U| \leq u_{\max} & (2) \\ & x_{i,t} \in U & \forall i \in I, t \in T_{i} \cup \{t_{i}\} & (3) \\ & x_{i,0} = u_{0} & \forall i \in I & (4) \\ & x_{i,t+1} = \delta_{u}(x_{i,t}, L(e_{i,t+1})) & \forall i \in I, t \in T_{i} & (5) \\ & N_{u,l} \subseteq 2^{2^{\mathcal{P}}} & \forall u \in U, l \in 2^{\mathcal{P}} & (6) \\ & L(e_{i,t+1}) \in N_{x_{i,t},L(e_{i,t})} & \forall i \in I, t \in T_{i} & (7) \end{array}$$

... that we solved using **local search**.



# Results

Results



\*Note: The detectors were also given to the baselines.

# Discussion





Legend:	— A3C	— PPO	- LRM+DDQN
— $\epsilon$ -optimal	— ACER	- DDQN	LRM+DQRM

1) Would the LSTM-based baselines eventually learn to solve these domains?

Large Cookie Domain

Reward per 10,000 steps 10050'▲' 0  $\mathbf{2}$ 10 6 8 Training steps (in millions) Legend: A3C PPO LRM+DDQN  $\epsilon$ -optimal ACER DDQN LRM+DQRM



2) How does LRM relate to other methods that learn automata to aid RL agents?

2) How does LRM relate to other methods that learn automata to aid RL agents?

(Add extra traces if the RM is imperfect) Collect Optimization RL agent Proposed  $\approx \pi^*_{\theta}(a|o,u)$ RM Traces Model Local  $\langle o/w, 0 \rangle \langle \Box, 0 \rangle; \langle o/w, 0 \rangle$ 1) DDQN  $\langle \square \odot, 0 \rangle$ Search 2) DQRM U1 Uэ (🗖 🔍 🖓 ( 🗖 🛈, 1  $(\square \odot, 1)$ Un (o/w, 0) (many traces!)

2) How does LRM relate to other methods that learn automata to aid RL agents?

(Add extra traces if the RM is imperfect) Collect Optimization RL agent Proposed  $\approx \pi^*_{\theta}(a|o,u)$ Traces Model RM Local  $\langle o/w, 0 \rangle \langle \Box, 0 \rangle; \langle o/w, 0 \rangle$ 1) DDQN ⟨□ , 0⟩ Search 2) DQRM U1 U2 ( 🗖 🔍 d ( 🗖 🛈. 1  $(\square \odot, 1)$ Un  $\langle o/w, 0 \rangle$ (many traces!)

What's the learning objective?

2) How does LRM relate to other methods that learn automata to aid RL agents?

A. LRM tries to solve POMDPs

2) How does LRM relate to other methods that learn automata to aid RL agents?

 $\ensuremath{\textbf{A}}\xspace$  LRM tries to solve POMDPs

POMDPs are hard because:

• 
$$P(o_{t+1}|o_0, \cdots, o_t, a_t) \neq P(o_{t+1}|o_t, a_t)$$

• 
$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$

• As a result, 
$$\pi^*(a|o_t) \ll \pi^*(a|o_0,\cdots,o_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

A. LRM tries to solve POMDPs

POMDPs are hard because:

• 
$$P(o_{t+1}|o_0, \cdots, o_t, a_t) \neq P(o_{t+1}|o_t, a_t)$$

• 
$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$

• As a result, 
$$\pi^*(a|o_t) \ll \pi^*(a|o_0,\cdots,o_t)$$

Thus, LRM's learning objective is to find a machine such that:

• 
$$P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, u_t, a_t)$$

• 
$$P(r_{t+1}|o_0,\cdots,o_t,a_t) = P(r_{t+1}|o_t,u_t,a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

A. LRM tries to solve POMDPs

POMDPs are hard because:

$$P(o_{t+1}|o_0, \cdots, o_t, a_t) \neq P(o_{t+1}|o_t, a_t)$$
$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$

• As a result, 
$$\pi^*(a|o_t) \ll \pi^*(a|o_0, \cdots, o_t)$$

Thus, LRM's learning objective is to find a machine such that:

• 
$$P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, u_t, a_t)$$

• 
$$P(r_{t+1}|o_0,\cdots,o_t,a_t) = P(r_{t+1}|o_t,u_t,a_t)$$

**Result**:  $\pi^*(a_t|o_t, u_t)$  optimally solves the POMDP.

2) How does LRM relate to other methods that learn automata to aid RL agents?

A. LRM tries to solve POMDPs



In the cookie domain, LRM learns this RM because it holds that  $P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, u_t, a_t)$ 

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

NMRDPs are hard because:

• 
$$P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, a_t)$$

• 
$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$

• As a result, 
$$\pi^*(a|o_t) \ll \pi^*(a|o_0,\cdots,o_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

NMRDPs are hard because:

$$P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, a_t)$$
  

$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$
  

$$As a result, \pi^*(a|o_t) \ll \pi^*(a|o_0, \cdots, o_t)$$

Thus, their learning objective is to find the smallest machine such that:

• 
$$P(r_{t+1}|o_0, \cdots, o_t, a_t) = P(r_{t+1}|o_t, u_t, a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

NMRDPs are hard because:

$$P(o_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}|o_t, a_t)$$

$$P(r_{t+1}|o_0, \cdots, o_t, a_t) \neq P(r_{t+1}|o_t, a_t)$$

$$As a result, \pi^*(a|o_t) \ll \pi^*(a|o_0, \cdots, o_t)$$

Thus, their learning objective is to find the smallest machine such that:

• 
$$P(r_{t+1}|o_0,\cdots,o_t,a_t) = P(r_{t+1}|o_t,u_t,a_t)$$

**Result**:  $\pi^*(a_t|o_t, u_t)$  optimally solves the NMRDP.

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

These methods do not work in our domains because our domains are not NMRDPs.

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

These methods do not work in our domains because our domains are not NMRDPs.



In the cookie domain:

$$\mathsf{P}(o_{t+1}|o_0,\cdots,o_t,a_t) 
eq \mathsf{P}(o_{t+1}|o_t,a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

These methods do not work in our domains because our domains are not NMRDPs.



In the cookie domain:

$$P(o_{t+1}|o_0,\cdots,o_t,a_t) 
eq P(o_{t+1}|o_t,a_t)$$

And

$$P(r_{t+1}|o_0, \cdots, o_t, a_t) = P(r_{t+1}|o_t, a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

These methods do not work in our domains because our domains are not NMRDPs.



In the cookie domain:

$$P(o_{t+1}|o_0,\cdots,o_t,a_t) 
eq P(o_{t+1}|o_t,a_t)$$

And

$$P(r_{t+1}|o_0, \cdots, o_t, a_t) = P(r_{t+1}|o_t, a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

B. Methods that solve NMRDPs (e.g., JIRP, ISA, SRMI)

These methods do not work in our domains because our domains are not NMRDPs.



In the cookie domain:

$$P(o_{t+1}|o_0,\cdots,o_t,a_t) 
eq P(o_{t+1}|o_t,a_t)$$

And

$$P(r_{t+1}|o_0,\cdots,o_t,a_t) = P(r_{t+1}|o_t,a_t)$$

2) How does LRM relate to other methods that learn automata to aid RL agents?

**C.** Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

2) How does LRM relate to other methods that learn automata to aid RL agents?

**C.** Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

MDPs are usually easy:

• 
$$P(o_{t+1}, r_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}, r_{t+1}|o_t, a_t)$$

... but MDPs with sparse rewards are hard.

2) How does LRM relate to other methods that learn automata to aid RL agents?

**C.** Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

MDPs are usually easy:

• 
$$P(o_{t+1}, r_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}, r_{t+1}|o_t, a_t)$$

... but MDPs with sparse rewards are hard.

Thus, their learning objective is to find the smallest machine such that:

- It accepts traces that **can** be generated by interacting with the environment.
- It rejects traces that cannot be generated by interacting with the environment.

2) How does LRM relate to other methods that learn automata to aid RL agents?

**C.** Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

MDPs are usually easy:

• 
$$P(o_{t+1}, r_{t+1}|o_0, \cdots, o_t, a_t) = P(o_{t+1}, r_{t+1}|o_t, a_t)$$

... but MDPs with sparse rewards are hard.

Thus, their learning objective is to find the smallest machine such that:

- It accepts traces that **can** be generated by interacting with the environment.
- It rejects traces that cannot be generated by interacting with the environment.

Result: They learn a high-level model that is then used to encourage exploration.

2) How does LRM relate to other methods that learn automata to aid RL agents?

C. Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

These methods do not work in our domains because our domains are not MDPs.

2) How does LRM relate to other methods that learn automata to aid RL agents?

**C.** Methods that solve MDPs with sparse rewards (e.g., DeepSynth)

These methods do not work in our domains because our domains are not MDPs.



2) How does LRM relate to other methods that learn automata to aid RL agents?

#### Summary

There are three learning objectives for combining automata learning with RL:

- A. [LRM] Learn an RM that makes the whole problem Markovian.
- B. [JIRP] Learn the smallest DFA that makes the reward function Markovian.
- C. [DeepSynth] Learn a high-level model of the environment.

2) How does LRM relate to other methods that learn automata to aid RL agents?

#### Summary

There are three learning objectives for combining automata learning with RL:

- A. [LRM] Learn an RM that makes the whole problem Markovian.
- B. [JIRP] Learn the smallest DFA that makes the reward function Markovian.
- C. [DeepSynth] Learn a high-level model of the environment.

So, what's the right learning objective?
2) How does LRM relate to other methods that learn automata to aid RL agents?

## Summary

There are three learning objectives for combining automata learning with RL:

- A. [LRM] Learn an RM that makes the whole problem Markovian.
- B. [JIRP] Learn the smallest DFA that makes the reward function Markovian.
- C. [DeepSynth] Learn a high-level model of the environment.

So, what's the right learning objective?  $()')_/$ 

2) How does LRM relate to other methods that learn automata to aid RL agents?

## Summary

There are three learning objectives for combining automata learning with RL:

- A. [LRM] Learn an RM that makes the whole problem Markovian.
- B. [JIRP] Learn the smallest DFA that makes the reward function Markovian.
- C. [DeepSynth] Learn a high-level model of the environment.

So, what's the right learning objective?  $(v)_{/}$ 

Method	Cookie	Symbol	2-Keys
JIRP	$0.6\pm0.8$	$-31.1 \pm 13.5$	$2.0\pm1.4$
DeepSynth	$0.4\pm0.6$	$-30.4\pm14.0$	$2.4\pm1.5$
LRM (ours)	$\textbf{197.2}\pm\textbf{2.1}$	$\textbf{460.4} \pm \textbf{15.7}$	$\textbf{86.6} \pm \textbf{9.4}$

## **Concluding Remarks**

https://bitbucket.org/RToroIcarte/lrm

Thanks! :)



Rodrigo



Ethan



Toryn



Rick



Margarita



Sheila