

## RL Tutorial

Soar Workshop 32 – Nate Derbinsky

While waiting...

1. Make sure you have internet access
2. Download Soar Tutorial package v9.3.2  
[code.google.com/p/soar/wiki/SoarTutorial](http://code.google.com/p/soar/wiki/SoarTutorial)
3. Download Graphviz  
[www.graphviz.org](http://www.graphviz.org)
4. Download Eclipse (with at least Java)  
[www.eclipse.org](http://www.eclipse.org)
5. Download tutorial support files  
[web.eecs.umich.edu/~nlderbin/workshop32](http://web.eecs.umich.edu/~nlderbin/workshop32)

19 June 2012

RL Tutorial

1

## Setting Expectations

**Not** a tutorial on Reinforcement Learning (RL)

Reinforcement Learning: An Introduction

Richard S. Sutton, Andrew G. Barto

Topics

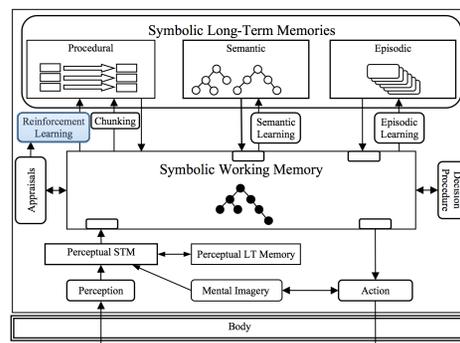
- RL as a learning mechanism in Soar
- Agent design
- Advanced issues
- Additional resources

19 June 2012

RL Tutorial

2

## Soar 9



19 June 2012

RL Tutorial

3

## Procedural Learning Mechanisms

### Chunking

- Converts *deliberation* in substates into *reaction* via rule compilation

- Creates new rules

### Reinforcement Learning

- *Tunes* operator numeric preferences to reflect expectation of reward

- Updates existing rules

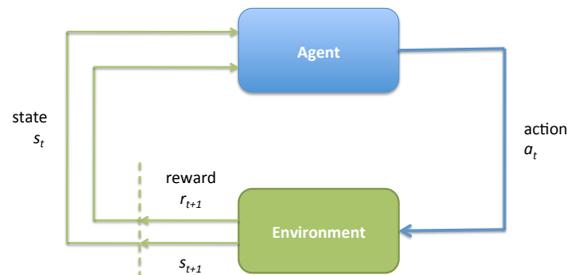
19 June 2012

RL Tutorial

4

## Reinforcement Learning Cycle

Goal: learn an action-selection policy such as to maximize expected receipt of future reward



19 June 2012

RL Tutorial

5

## Left-Right Demo

1. Soar Java Debugger
2. Source left-right agent  
Agents/left-right/left-right.soar



19 June 2012

RL Tutorial

6

## Left-Right Demo

*Script*

1. srand 5041231
2. step
3. run 1 -p
4. click: op\_pref tab
  - note numeric indifferents
5. print left-right\*rl\*left
6. print left-right\*rl\*right
7. run
  - note movement direction
8. print left-right\*rl\*left
9. print left-right\*rl\*right
10. init-soar
11. Repeat from #2 (~5 times)

19 June 2012

RL Tutorial

7

## Left-Right: Takeaways

Reinforcement learning changes rules in procedural memory

- Changes are persistent
- Change affects numeric indifferent preferences, which in turn affects the selection of operators
- Change is in the direction of the underlying reward signal (will discuss this more shortly)

19 June 2012

RL Tutorial

8

## Components of RL

- RL rules
- Reward representation
- Learning

19 June 2012

RL Tutorial

9

## RL Rules

The RL mechanism maintains Q-values for state-operator pairs in specially formulated rules, identified by syntax

- RHS must be a single action, asserting a single numeric indifferent preference with a constant value

```

sp {left-right*rl*left      sp {left-right*rl*right
  (state <s> ^name left-right  (state <s> ^name left-right
    ^operator <op> +)          ^operator <op> +)
  (<op> ^name move            (<op> ^name move
    ^dir left)                ^dir right)
--> (<s> ^operator <op> = 0)    --> (<s> ^operator <op> = 0)
}                                }

```

19 June 2012

RL Tutorial

10

## Pop Quiz!

Which of the following are valid RL rules?

✘ sp {rule1  
 (state <s> ^name quiz  
     ^operator <op> +)  
 (<op> ^name op-name)  
 -->  
 (<s> ^operator <op> = 1, >)}

✘ sp {rule2  
 (state <s> ^name quiz  
     ^operator <op> +)  
 (<op> ^name op-name)  
 -->  
 (<s> ^operator <op> = 0.8)  
 (write |debugging|)}

☺ sp {rule3  
 (state <s> ^name left-right  
     ^operator <op> +)  
 (<op> ^name move  
     ^dir right)  
 -->  
 (<s> ^operator <op> = -1.4)}

✘ sp {rule4  
 (state <s> ^name left-right  
     ^operator <op> +)  
 (<op> ^name move  
     ^dir right)  
 -->  
 (<s> ^operator <op> = <one>)}

19 June 2012

RL Tutorial

11

## Left-Right Demo

*Focus: RL Rules*

1. Soar Java Debugger
2. Source left-right agent  
 Agents/left-right/left-right.soar
3. print --full --rl
4. run
5. print --full --rl
6. print --rl
7. Attempt sourcing rules from previous slide
  - Check with print --rl

19 June 2012

RL Tutorial

12

## Reward Representation

Soar creates a reward-link structure on each state in WM

- Soar Java Debugger
  1. step 5
  2. print --exact (\* ^reward-link \*)

Reward is recognized by syntax

```
(<reward-link> ^reward <r>)
(<r> ^value [val])
```

- [val] must be a numeric constant (integer or float)
- The reward-link is **not** directly modified by IO
- The reward-link is **not** automatically cleaned
- Reward is collected at the beginning of each *decide* phase
- Reward on a state's reward-link pertains to that state (more on this later)
- Multiple reward values are summed by default

19 June 2012

RL Tutorial

13

## Reward Rule Examples

```
sp {left-right*reward*left      sp {left-right*reward*right
  (state <s> ^name left-right    (state <s> ^name left-right
    ^location left              ^location right
    ^reward-link <rl>)}         ^reward-link <rl>)}
-->                               -->
(<rl> ^reward <r>)                (<rl> ^reward <r>)
(<r> ^value -1)                   (<r> ^value 1)
}
```



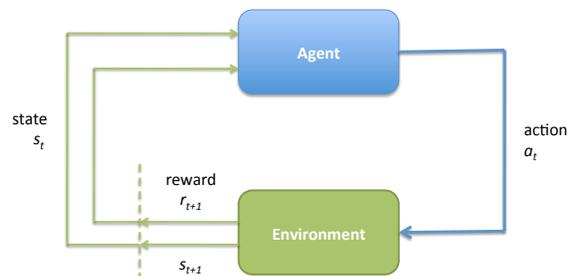
19 June 2012

RL Tutorial

14

## Reinforcement Learning Cycle

Goal: learn an action-selection policy such as to maximize expected receipt of future reward



19 June 2012

RL Tutorial

15

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d					
d+1					

19 June 2012

RL Tutorial

16

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>				
d+1					

19 June 2012

RL Tutorial

17

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>			
d+1					

19 June 2012

RL Tutorial

18

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>	select operator <sub>d</sub>		
d+1					

19 June 2012

RL Tutorial

19

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>	select operator <sub>d</sub>		initiate external action(s)
d+1					

19 June 2012

RL Tutorial

20

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>	select operator <sub>d</sub>		initiate external action(s)
d+1	state <sub>d+1</sub> reward <sub>d+1</sub>				

19 June 2012

RL Tutorial

21

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>	select operator <sub>d</sub>		initiate external action(s)
d+1	state <sub>d+1</sub> reward <sub>d+1</sub>	evaluate operators <sub>d+1</sub>			

19 June 2012

RL Tutorial

22

## RL Cycle in Soar

Decision	Input	Propose	Decide	Apply	Output
d	state <sub>d</sub>	evaluate operators <sub>d</sub>	select operator <sub>d</sub>		initiate external action(s)
d+1	state <sub>d+1</sub> reward <sub>d+1</sub>	evaluate operators <sub>d+1</sub>	select operator <sub>d+1</sub> update policy <sub>d</sub>		

19 June 2012

RL Tutorial

23

## RL Updates

- Takes place during *decide* phase, after operator selection
- For all RL rule instantiations (*n*) that supported the *last* selected operator

$$\text{value}_{d+1} = \text{value}_d + (\delta_d/n)$$

Where, roughly...

$$\delta_d = \alpha[\text{reward}_{d+1} + \gamma(q_{d+1}) - \text{value}_d]$$

Where...

- $\alpha$  is a parameter (learning rate)
- $\gamma$  is a parameter (discount rate)
- $q_{d+1}$  is dictated by learning policy
  - On-policy (SARSA): value of selected operator
  - Off-policy (Q-learning): value of operator with maximum selection probability

19 June 2012

RL Tutorial

24

## Left-Right Demo

*Focus: RL Updates*

1. Soar Java Debugger
2. Source left-right agent  
Agents/left-right/left-right.soar
3. watch --rl
4. run
5. print s1
6. print --depth 2 rl
7. rl --stats  $\delta = 0.3[1+0.9(0)-0] = 0.3$
8. rl  $\delta = 0.3[1+0.9(0)-0.3] = 0.21$

19 June 2012

RL Tutorial

25

## Controlling RL in Soar

Get/Set a parameter:

- rl [-g|--get] <name>
- rl [-s|--set] <name> <value>

Try changing learning rate to 1 (deterministic)

1. rl
2. rl --set learning-rate 1
3. rl

19 June 2012

RL Tutorial

26

## RL Agent Design

1. RL rules
2. Reward rules
3. Enable RL

Baseline: left-right-start.soar

[web.eecs.umich.edu/~nlderbin/workshop32/rl](http://web.eecs.umich.edu/~nlderbin/workshop32/rl)

19 June 2012

RL Tutorial

27

## 1. RL Rules

*left-right-1.soar*

- a) Add RL rules for moving left/right (slide 10)
  - Check via print --rl
- b) Remove indifferent preference from operator proposal rules
  - Note: this will cause a tie impasse if (a) is not done correctly

19 June 2012

RL Tutorial

28

## 2. Reward Rules

*left-right-2.soar*

Add rules to add reward to reward-link, dependent upon the location to which the agent decided to move (slide 14)

19 June 2012

RL Tutorial

29

## 3. Enable RL

*left-right-3.soar*

By default, RL is **disabled**. To enable, set the parameter (slide 26):

```
rl --set learning on
```

Note: the agent will function without this step, but rules will not update, and hence will act identically to the baseline agent.

19 June 2012

RL Tutorial

30

## Advanced Issues

- Gaps in rule coverage
- RL in substates
- Exploration policies
- Rule generation

19 June 2012

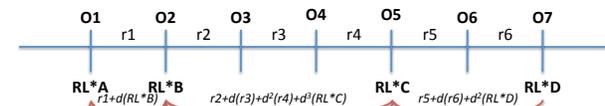
RL Tutorial

31

## Gaps in Rule Coverage

**Gap:** one or more contiguous decision cycles during which no RL rules fire

By default, Soar will automatically propagate RL updates over gaps, where rewards are discounted with respect to the length of the gap

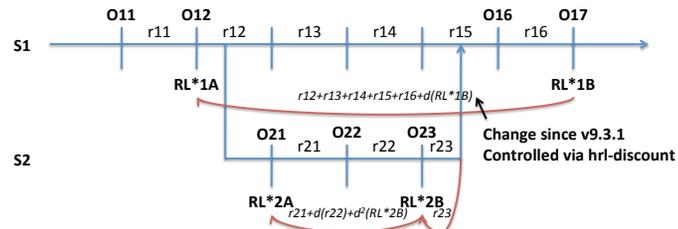


19 June 2012

RL Tutorial

32

## RL in Substates



- Rewards are collected independently on each state
- Rewards at a superstate are attributed to the last RL-supported, selected operator

19 June 2012

RL Tutorial

33

## Exploration Policies

There are numerous policies for selecting operators probabilistically (see Manual)

- Deterministic (first, last)
- Softmax (default)
- Epsilon Greedy (change since v9.3.1 – doesn't switch!)
- Boltzmann

Pertinent commands (see Manual)

- indifferent-selection
- predict
- select

19 June 2012

RL Tutorial

34

## Rule Generation

The number of RL rules required for an agent is usually unfeasibly large to write by hand

### Options

- gp command
- Rule templates
- Chunking

19 June 2012

RL Tutorial

35

## gp Command

- Similar syntax to sp command
- Generates rules, at source time, according to the cross-product of a fixed number of dimensions, each with a fixed domain
  - Very fast, potentially creates a lot more rules than will be encountered by the agent

```
gp {example*gp
  (state <s> ^name left-right      > Total: 36 productions sourced.
    ^operator <op> +                > print
    ^d1 [ a b c ]                  example*gp*36
    ^d2 [ 1 2 3 ]                  example*gp*35
  (<op> ^name move
    ^dir [ left right up down ])
  -->
  (<s> ^operator <op> = 0)
}
```

19 June 2012

RL Tutorial

36

## Left-Right Demo

*Focus: gp (left-right-3-gp.soar)*

Replace RL rules (slide 28) with a single gp command

- Source agent
- Verify rules: `print --full --rl`

19 June 2012

RL Tutorial

37

## Rule Templates

- Allows Soar to dynamically generate new RL rules based on a predefined pattern as the agent encounters novel states
- Similar to syntax of RL rules
  - Requires `:template` flag
    - Note: without will be an RL rule that matches multiple states
  - Numeric indifferent preference value can be a variable
- Behavior
  - Rather than firing, creates new RL rules from instantiations
  - Significantly slower than gp

19 June 2012

RL Tutorial

38

## Left-Right Demo

*Focus: template (left-right-3-template.soar)*

Replace RL rules (slide 28) with a single rule template

1. Source agent
2. Verify rules
  - `print --full --rl`
  - `print --full --template`
3. run
4. Verify rules
  - `print --full --rl`
  - `print --full --template`

19 June 2012

RL Tutorial

39

## Chunking

*“Since RL rules are regular productions, they can be learned by chunking just like any other production. This method is more general than using the gp command or rule templates, and is useful if the environment state consists of arbitrarily complex relational structures that cannot be enumerated.”*

- Soar Manual

19 June 2012

RL Tutorial

40

## Additional Resources

- Documentation
- Demo agents
- Readings

19 June 2012

RL Tutorial

41

## Documentation

### Manual and Tutorial Documentation/

### Additional Topics

- RL update details
- Eligibility traces
- Additional learning agent walkthrough (water-jug)
- Usage: commands, parameters, statistics, etc.

...

19 June 2012

RL Tutorial

42

## Demo Agents

### Agents/

- Left-Right
- Water Jug RL
  - Discussed in Soar-RL Tutorial

19 June 2012

RL Tutorial

43

## Select Readings

*code.google.com/p/soar/wiki/Publications*

- 2005
  - Soar-RL: Integrating Reinforcement Learning with Soar
    - Shelley Nason, John E. Laird (Cognitive Systems Research)
- 2007
  - The Importance of Action History in Decision Making and Reinforcement Learning
    - Yongjia Wang, John E. Laird (ICCM)
- 2008
  - A Computational Unification of Cognitive Control, Emotion, and Learning
    - Robert P. Mariniere III (Dissertation)
- 2009
  - Learning to Play Mario
    - Shiwali Mohan (Technical Report)
  - Hierarchical Reinforcement Learning in the Taxicab Domain
    - Mitchell K. Bloch (Technical Report)
- 2010
  - Instance-Based Online Learning of Deterministic Relational Action Models
    - Joseph Xu, John E. Laird (AAAI)
  - Using Imagery to Simplify Perceptual Abstraction in Reinforcement Learning Agents
    - Samuel Westermute (AAAI)
- 2011
  - Learning to Use Episodic memory
    - Nicholas Gorski, John E. Laird (Cognitive Systems Research)
- 2012
  - The Soar Cognitive Architecture
    - Laird (MIT Press)

19 June 2012

RL Tutorial

44

## Extra Time

- Extend Water-Jug base agent to use RL
  - See Soar-RL Tutorial for guidance