

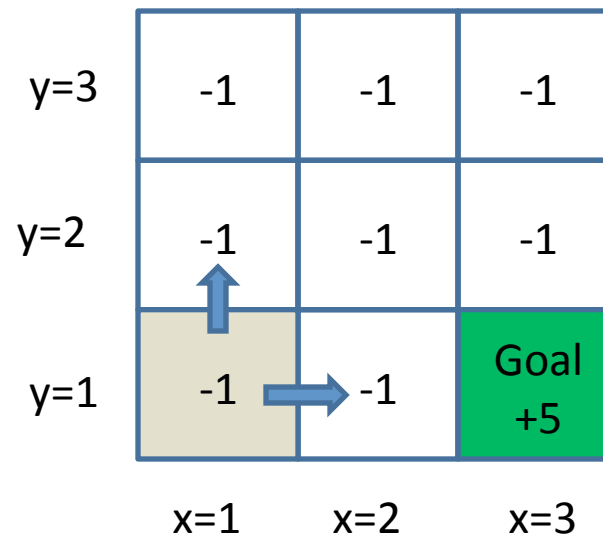
Value Function Approximation with Hierarchical Clustering in Soar-RL

Yongjia Wang
University of Michigan

Soar-RL Value Function

- Soar-RL relies on learning the Q value function
 - Q Value Function, $Q(\text{state}, \text{action})$

$(x=1, y=1)$, (move-up) $\rightarrow Q=1$
 $(x=1, y=1)$, (move-right) $\rightarrow Q=3$
 $(x=1, y=1)$, (move-down) $\rightarrow Q=2$
 $(x=1, y=1)$, (move-left) $\rightarrow Q=2$



Soar-RL Value Function

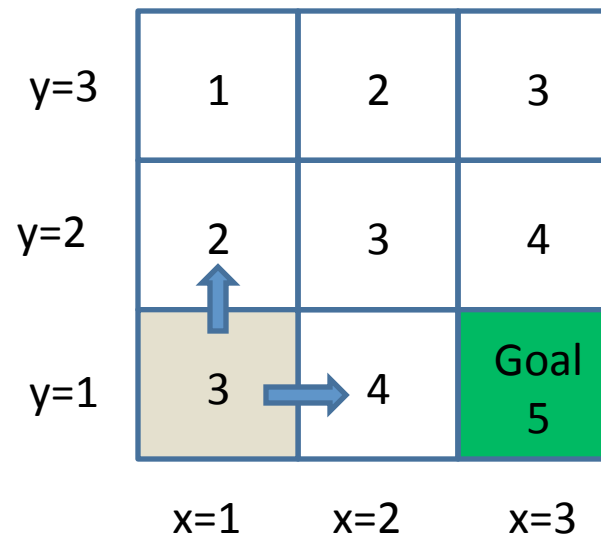
- Soar-RL relies on learning the Q value function
 - Q Value Function, $Q(\text{state}, \text{action})$
 - $\text{Value}(\text{state}) = \max(Q(\text{state}, \text{action}_i)) + r$

$(x=1, y=1), (\text{move-up}) \rightarrow Q=1$

$(x=1, y=1), (\text{move-right}) \rightarrow Q=3$

$(x=1, y=1), (\text{move-left}) \rightarrow Q=2$

$(x=1, y=1), (\text{move-down}) \rightarrow Q=2$

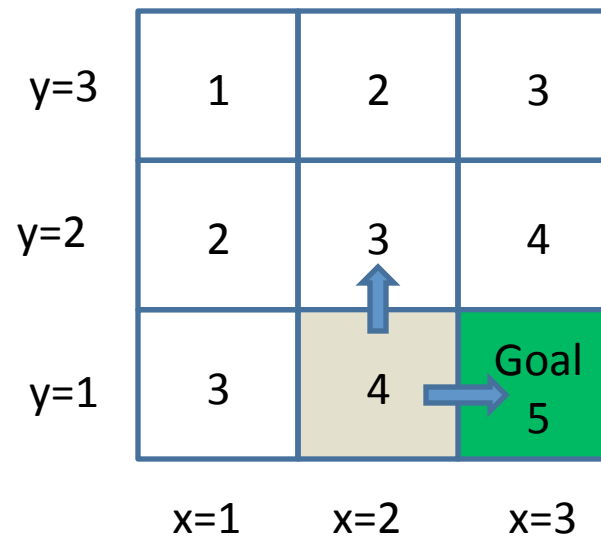


Soar-RL Value Function

- Soar-RL relies on learning the Q value function
 - Q Value Function, $Q(\text{state}, \text{action})$
 - $\text{Value}(\text{state}) = \max(Q(\text{state}, \text{action}_i))$

$(x=1, y=1), (\text{move-up}) \rightarrow Q=1$
 $(x=1, y=1), (\text{move-right}) \rightarrow Q=3$
 $(x=1, y=1), (\text{move-left}) \rightarrow Q=2$
 $(x=1, y=1), (\text{move-down}) \rightarrow Q=2$

$(x=2, y=1), (\text{move-up}) \rightarrow Q=2$
 $(x=2, y=1), (\text{move-down}) \rightarrow Q=3$
 $(x=2, y=1), (\text{move-left}) \rightarrow Q=2$
 $(x=2, y=1), (\text{move-right}) \rightarrow Q=4$



Soar-RL Value Function

- Soar-RL relies on learning the Q value function
 - Q Value Function, $Q(\text{state}, \text{action})$
 - $\text{Value}(\text{state}) = \max(Q(\text{state}, \text{action}_i))$
 - What about high dimensional continuous space?

(x=1, y=1, z=1,)

(x=1.1, y=1, z=1,)

(x=1.2, y=1, z=1,)

... ..

- Need Value Function Approximation

Motivations

- Improve Value Function Approximation in Soar-RL
 - High dimension continuous space
- Evaluate with challenging tasks
 - Characterize the tasks
 - Compare with baselines

Outline

- Background
- Evaluation Task
- Implementation
- Results
- Nuggets & Coal

Value Function Approximation -Parametric Form

Simple linear model:

$$Q(x,y,\text{up}) = w_1x + w_2y + w_3$$

$$Q(x,y,\text{down}) = w_1x + w_2y + w_4$$

$$Q(x,y,\text{left}) = w_1x + w_2y + w_5$$

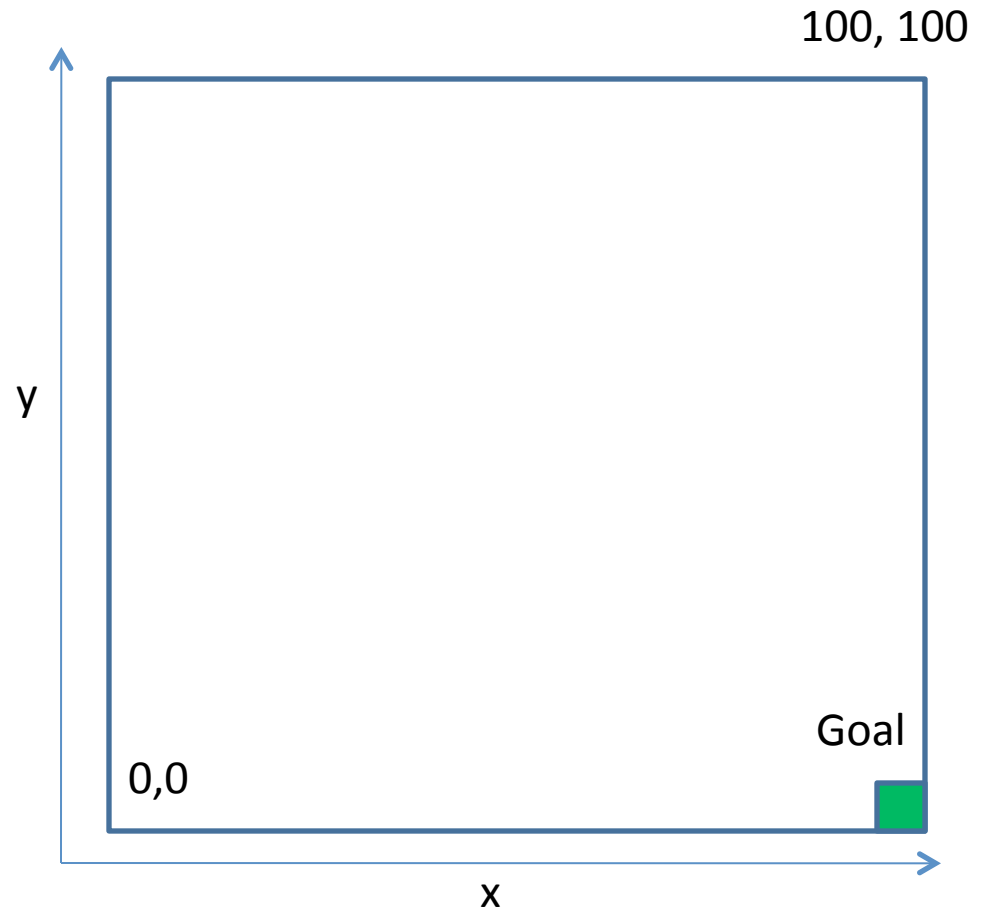
$$Q(x,y,\text{right}) = w_1x + w_2y + w_6$$

Solution:

$$w_1 > 0, w_2 < 0$$

$$w_3 < w_4$$

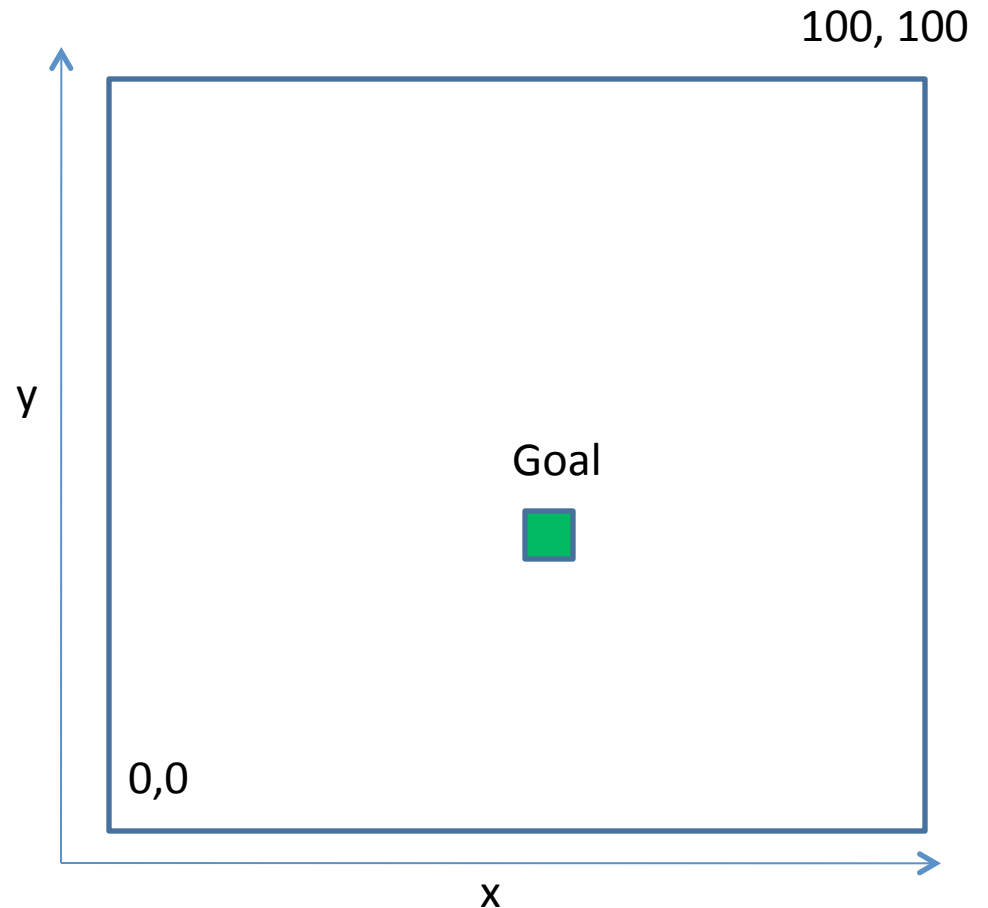
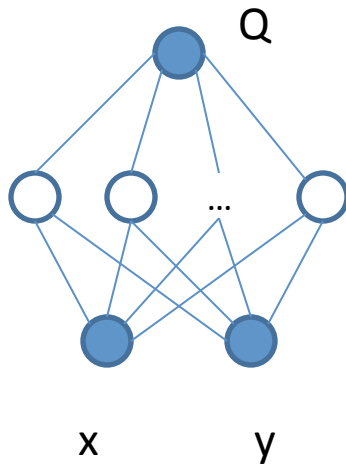
$$w_5 > w_6$$



Value Function Approximation -Parametric Form

Simple linear model won't work.
Need nonlinear basis function.

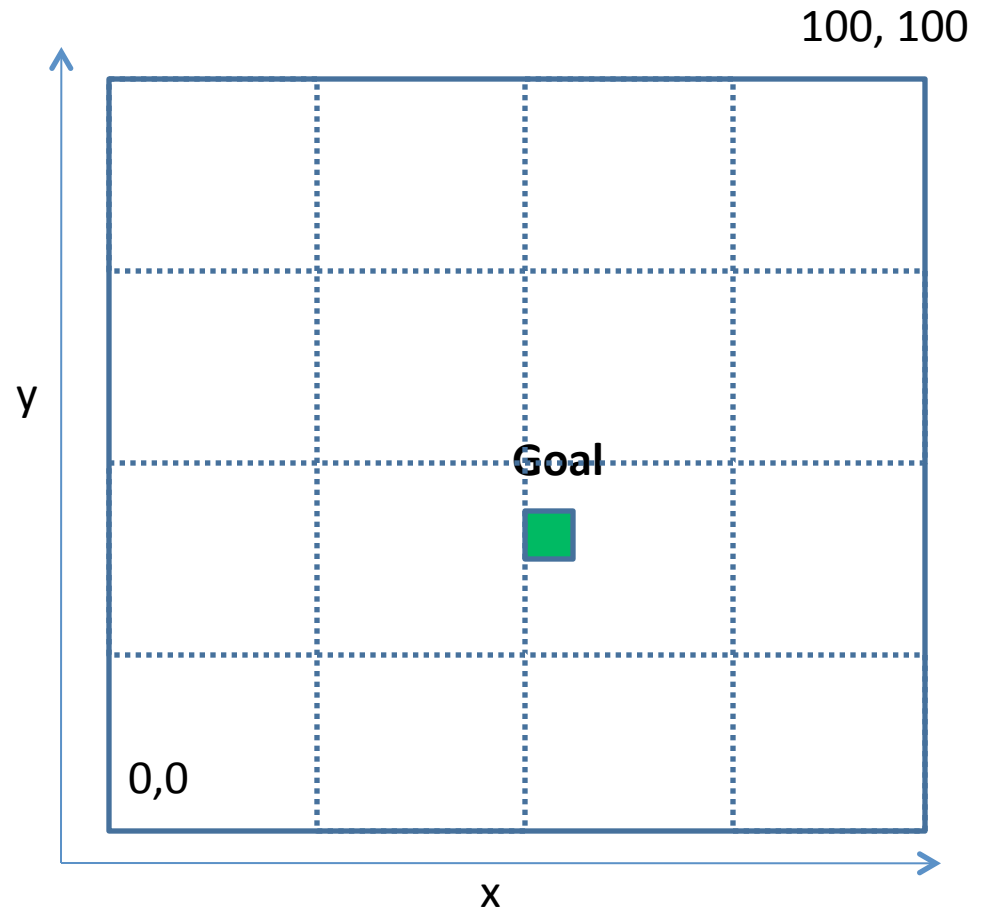
Neural network:



Global Approximators: updates in parameters during learning
affects the entire state space

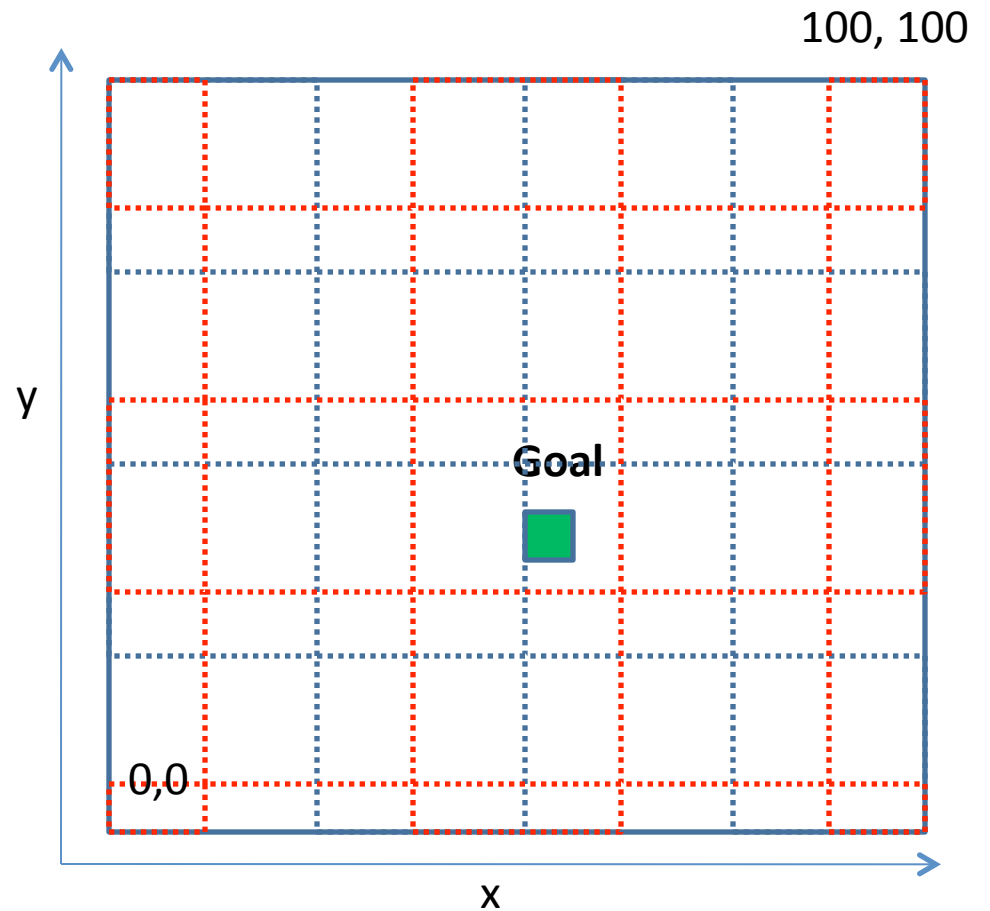
Value Function Approximation -Nonparametric methods

Tile coding



Value Function Approximation -Nonparametric methods

Coarse coding



Value Function Approximation -Nonparametric methods

Soar-RL supports arbitrary coarse coding
by using “coarse features” in RL rules

RL rule_1: $(x_0 < x < x_1) \rightarrow Q=q_1$

RL rule_2: $(y_0 < y < y_1) \rightarrow Q=q_2$

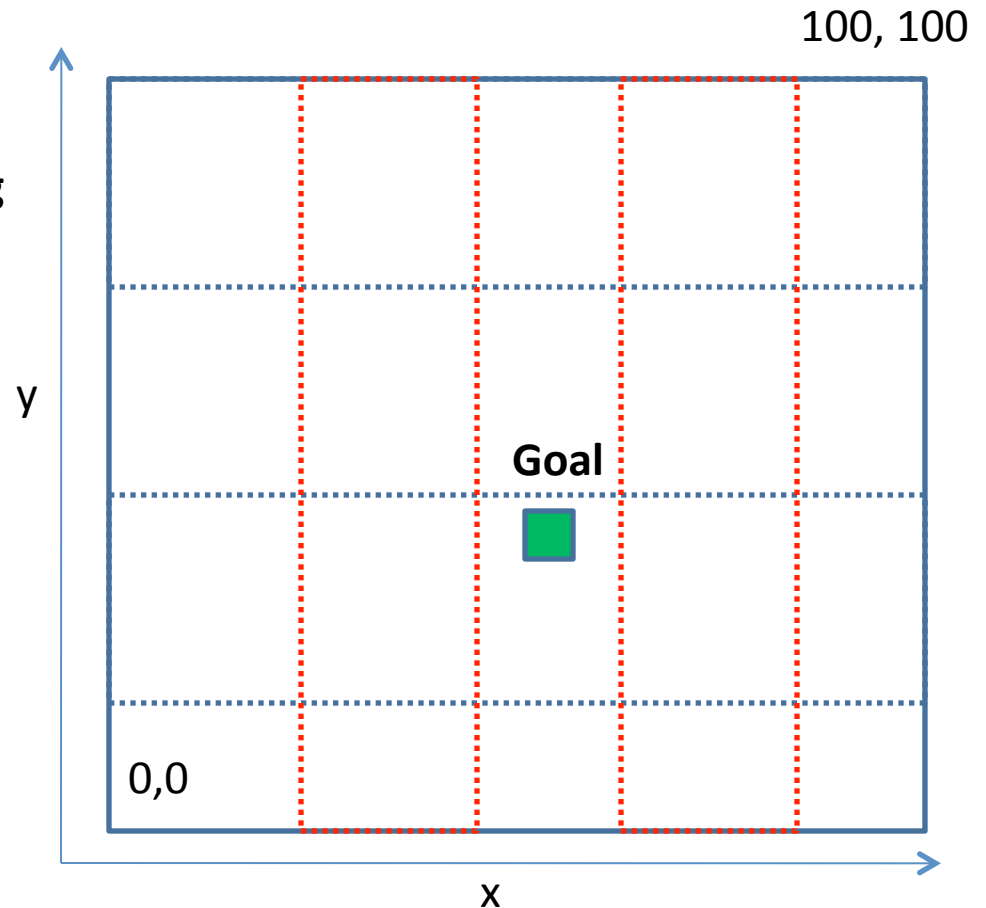
... ..

RL rule_n

$$Q(x_0 < x < x_1, y_0 < y < y_1) = q_1 + q_2$$

How to decide the boundaries?

How to encode rules for high dimensional space?

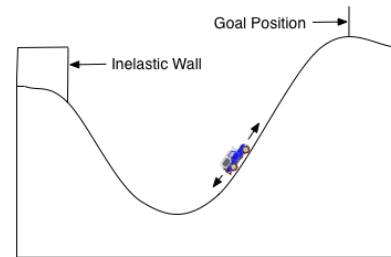
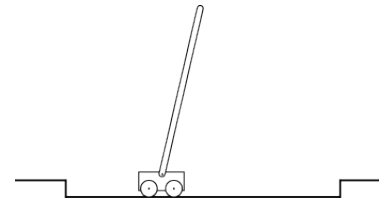


Motivations

- Improve Value Function Approximation in Soar-RL
 - Manually design coarse features can be challenging, even infeasible
 - Learn coarse features automatically
 - Compatible with existing VFA scheme in Soar-RL
- Evaluate with challenging tasks
 - Characterize the tasks
 - Compare with baselines

Traditional RL

- Dynamic control problem
 - Pole balancing
 - Mountain car
 - Robotic arm
 - Helicopter



Object Oriented RL

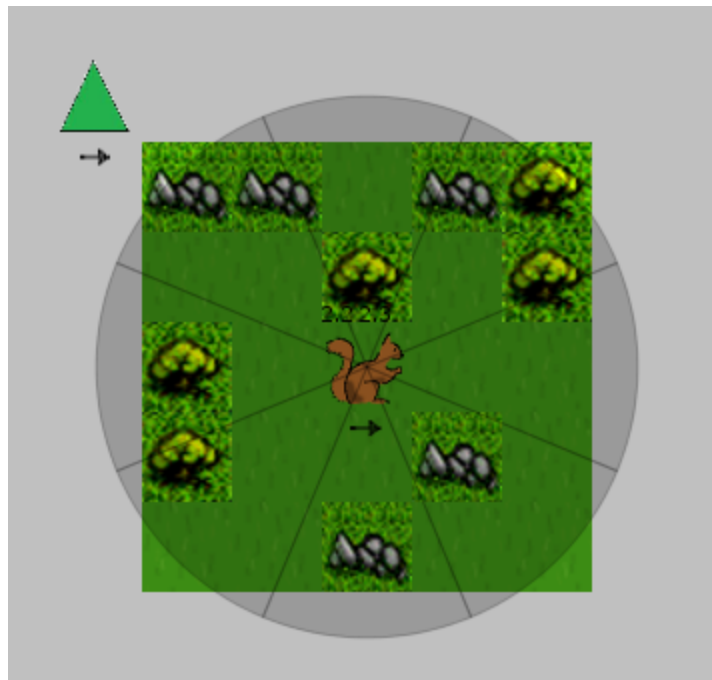
- Object Oriented MDP (OOMDP) framework
 - *Diuk et al. 2008*
 - Representation and learning of action models, not value function approximation
 - Limited object diversity



Motivations

- Improve Value Function Approximation in Soar-RL
 - Manually design coarse features can be challenging
 - Learn coarse features by hierarchical clustering
 - Compatible with existing VFA scheme in Soar-RL
- Evaluate with challenging tasks (Object oriented environments)
 - State representation consists of objects
 - Objects are functionally and perceptually diverse, represented by feature vectors
 - Complex interactions among objects

Evaluation Environment



Distance: 3

Obstacles: bush at 1, rock at 2

Prey:

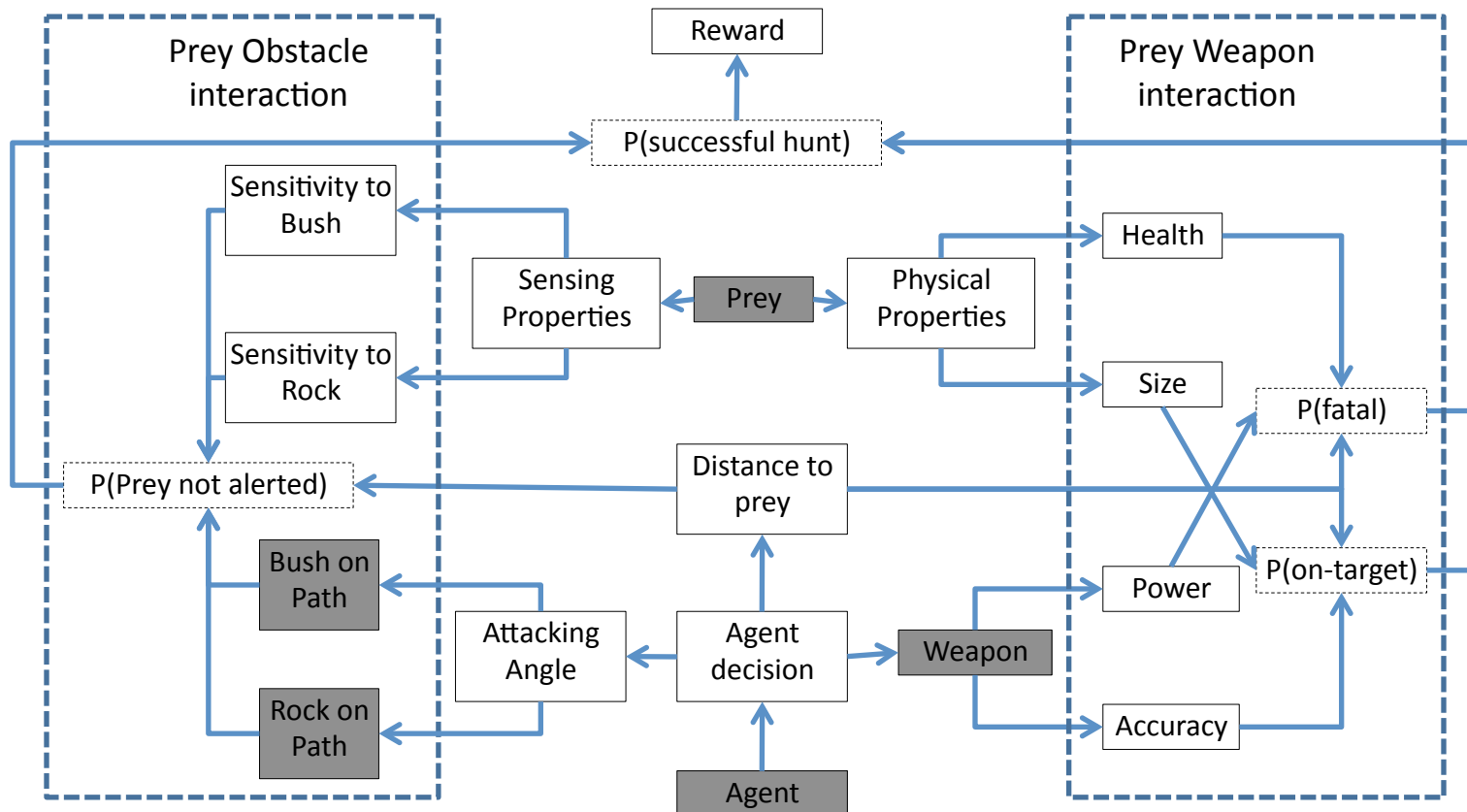
size	health	bush-sensitivity	rock-sensitivity
1	2	2	4

Weapon:

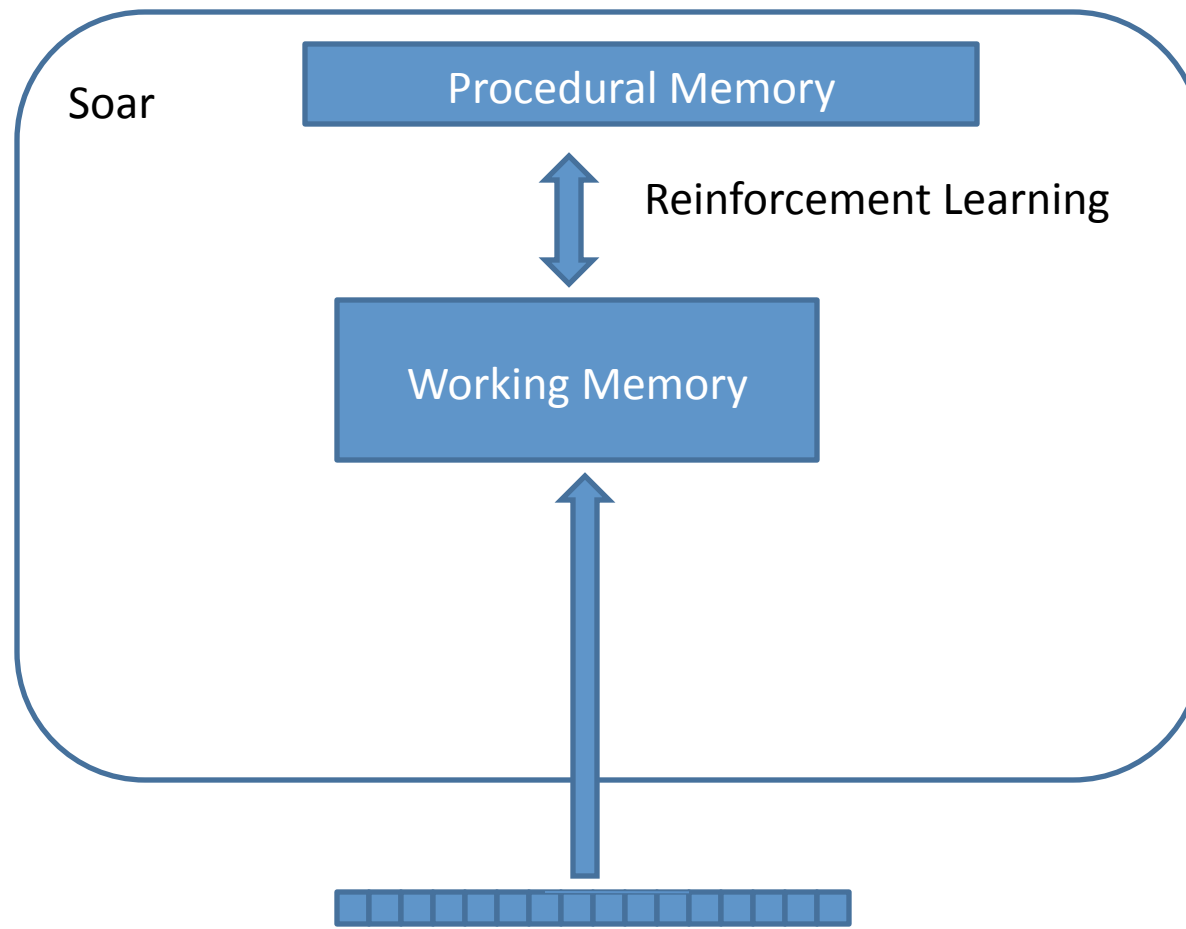
accuracy	power
5	2

1. Not alert the prey
2. Shot on target
3. Fatally wound the prey
4. Reward

Environment Model

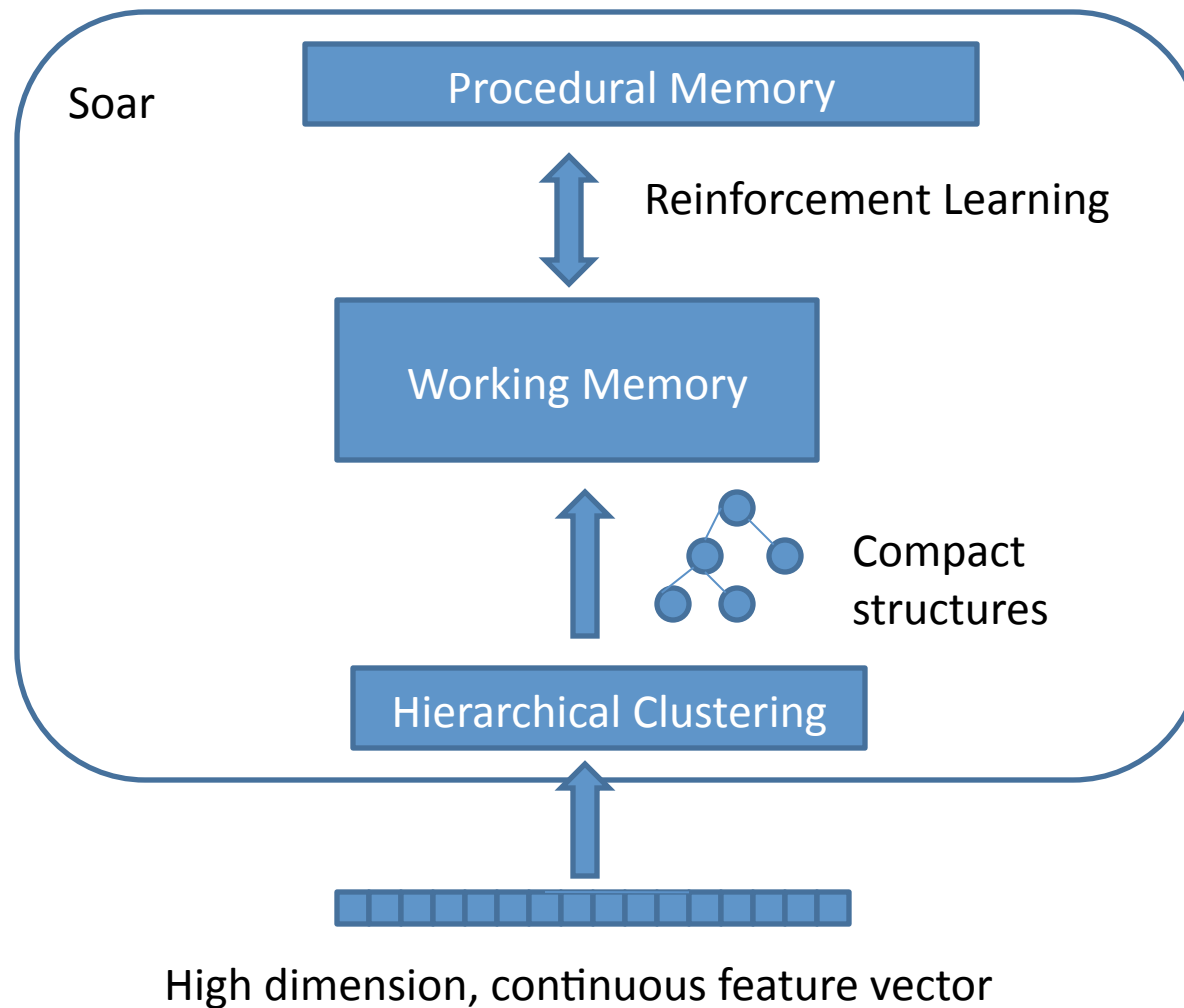


Integrate Hierarchical Clustering with Soar-RL



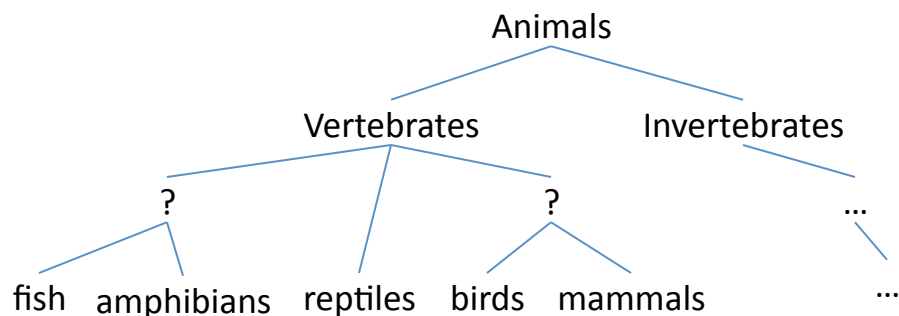
High dimension, continuous state space

Integrate Hierarchical Clustering with Soar-RL



Integrate Category Learning with RL

- COBWEB (D. Fisher 1987)
 - Incremental learning
 - process one instance a time
 - Hierarchical clustering
 - Take numeric feature vector as input
 - Automatically create category symbols
 - Hierarchical structure

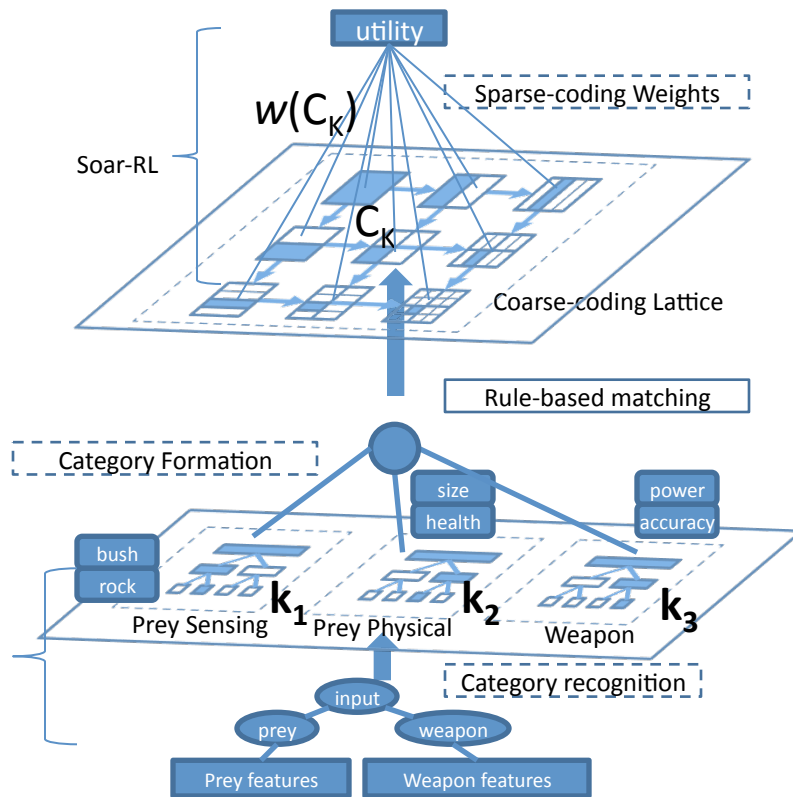


Details of the Algorithm

Procedural Memory

Working Memory

COBWEB



Vector $K = (k_1, k_2, k_3)$, represents the activated clusters from each of the three hierarchies.

A cell C_k , represents the rule matching the clusters K , and $a(C_k)$ is the activation:

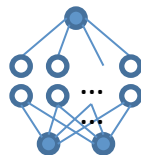
$$a(C_k) = \begin{cases} 1 & \text{if matches} \\ 0 & \text{if not matches} \end{cases}$$

Connection weight $w(C_k)$ is the Q value associated with the RL rule. The predicted value:

$$y = \sum_{C_k} w(C_k) a(C_k)$$

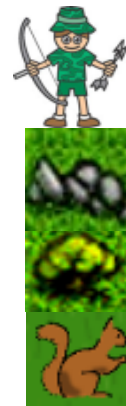
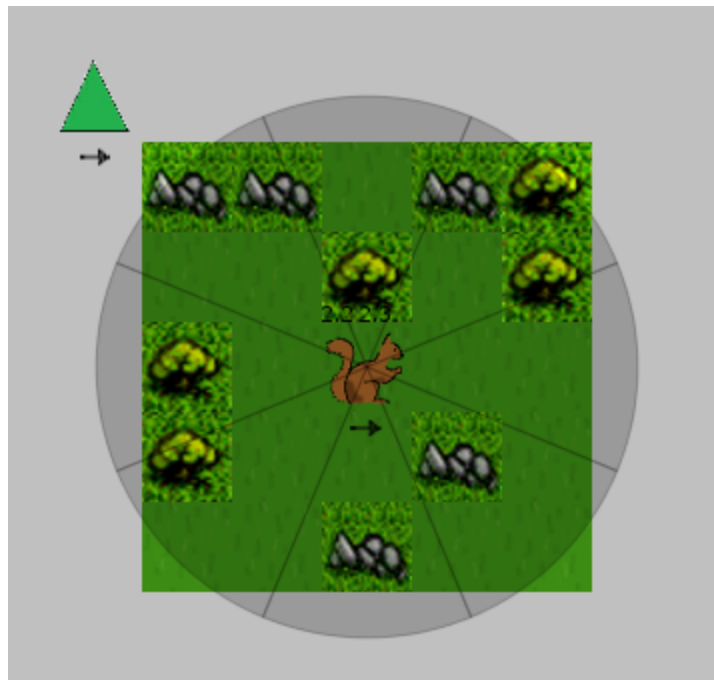
Update in Soar-RL is calculated in the same way as in stochastic gradient descent methods

$$\Delta w(C_k) = \frac{\alpha}{\sum_{C_k} a(C_k)} (t - y) a(C_k)$$



Like two-hidden-layer network, except for competitive learning and sparse activity in the "hidden layer".

Evaluation Environment



Distance: 3

Obstacles: bush at 1, rock at 2

Prey:

size	health	bush-sensitivity	rock-sensitivity
1	2	2	4

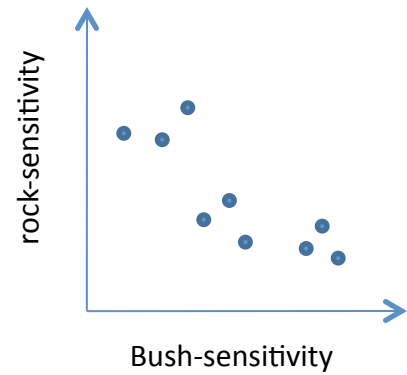
Weapon:

accuracy	power
5	2

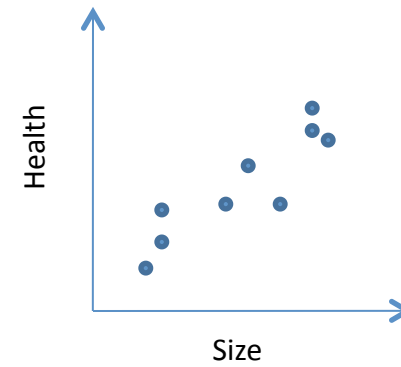
1. Not alert the prey
2. Shot on target
3. Fatally wound the prey
4. Reward

Data Used for Evaluation

Prey

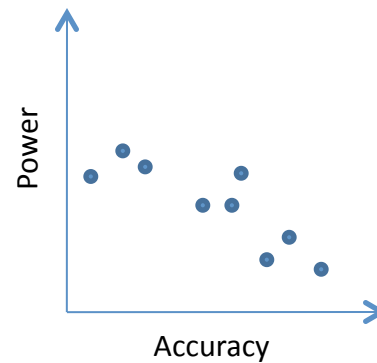


Prey adapted to specific terrain



Health and size are positively correlated

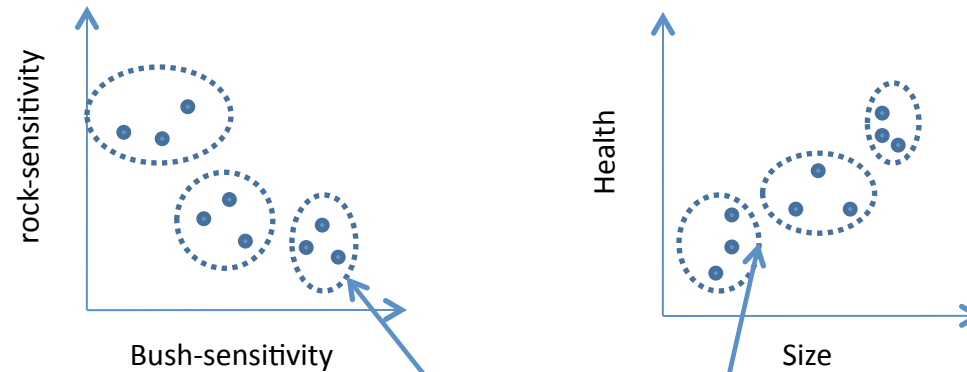
Weapon



Power and accuracy are negatively correlated

Data Used for Evaluation

Prey

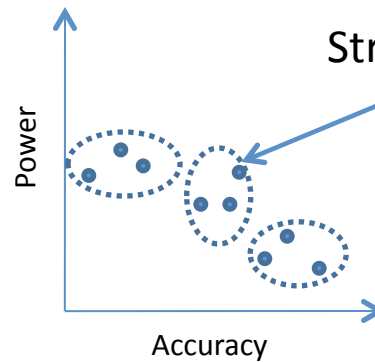


Prey adapted to specific terrain

Health and size are positively correlated

Structures in the data

Weapon

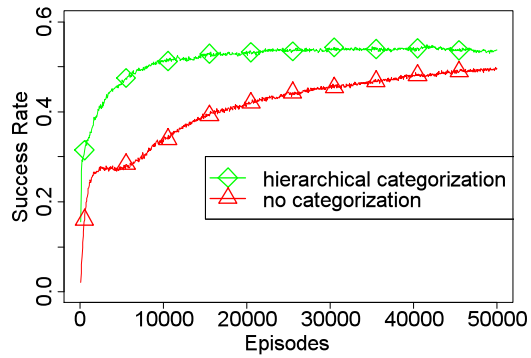


Power and accuracy are negatively correlated

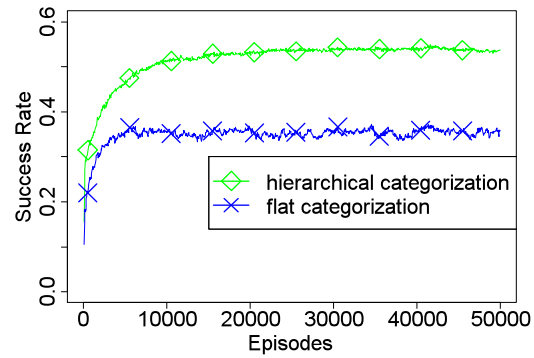
Evaluation

- Target system (Hierarchical clustering + RL)
 - Hierarchical Categorization
 - Using specific functional features
 - prey sensing, prey physical, weapon
- Baseline 1
 - No categorization
 - No generalization, slow learning
- Baseline 2
 - Flat categorization
 - Can be over-general, fail to discriminate more specific cases
- Baseline 3
 - A single hierarchy for prey including all perceptual features
 - Hierarchy with lower quality: some perceptual features are irrelevant to the task (such as color)

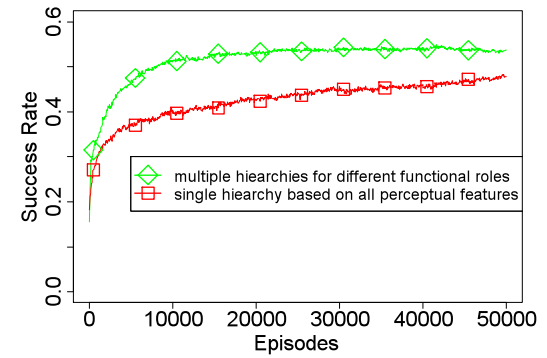
Results



Against baseline1
No generalization



Against baseline2
Over-generalization



Against baseline3
Hierarchy of lower quality

	Categorization	Hierarchical	Functional Hierarchies
Target	Yes	Yes	Yes
Baseline1	No	NA	Yes
Baseline2	Yes	No	Yes
Baseline3	Yes	Yes	No

Nuggets & Coal

- Nuggets
 - Integrated hierarchical category learning and Soar-RL for object-oriented environments with complex features
 - Provide means to take advantage of prior knowledge about functional features
- Coal
 - Single point evaluation
 - Limited complexity in object diversity
 - Only one way flow of information: from categorization to VFA