

# RL-General Formulation

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in RL-General Formulation

## 1. RL and Uncertainty

- Examples

## 2. Elements of RL

- Policy, reward, value, model, state

## 3. Tabular vs Approximate Solutions

## RL and Uncertainty

- RL considers a goal-directed *agent* interacting with an *environment*
  - By contrast, *supervised learning* does not explicitly specify how such an ability would finally be useful
- Agent has to operate despite significant uncertainty about the environment

# Examples of RL

- A master chess player makes a move
  - Choice is informed both by planning/anticipating possible replies and counter-replies and by immediate, intuitive judgments of positions and moves
- Roomba decides to enter a new room for more trash or return to its battery recharging station
  - Decision based on the current charge and how easily it has been able to find recharger in the past.

# Elements of RL

- Four main elements of RL
  1. *Policy*
  2. *Reward* signal
  3. *Value* function
  4. optionally, a *Model* of the environment
- RL methods heavily rely on concept of *state*

# Policy

- A mapping from perceived states to actions to be taken when in those states
- Policy is the core of an RL agent
  - It alone is sufficient to determine behavior
- Can be a lookup table, or involve extensive computation, such as a search
- Policies may be stochastic

# Reward Signal

- On each time step, the environment sends to the agent a single number called the reward
- Agent's sole objective is to maximize the total reward it receives over the long run
- Reward signals may be stochastic functions of the state of the environment and actions taken

# Value Function

- Specifies what is good in the long run
- Rewards are primary, whereas values, as predictions of rewards, are secondary
- Without rewards there could be no values, and the only purpose of estimating values is to achieve more reward



# Model of the Environment

- Model might predict the resultant next state and next reward
- Models are used for planning
  - Consider future situations before they are experienced
- Model-free methods are trial and error learners

# State

- RL methods heavily rely on concept of state
- As input to the policy and value function
- Conveys how the environment is to the agent
- State signal is produced by a preprocessing system that is nominally part of the environment
- RL does not focus on determining the state

# Example of Tic-Tac-Toe

- Two players take turns on a 3x3 board
  - One player plays Xs and the other Os
  - Player wins by placing three marks in a row, horizontally, vertically, or diagonally, as the X player has in the game

X	O	O
O	X	X
		X

- If neither has three in a row, game is drawn
- Skilled player can never lose
  - Assume playing against imperfect player
  - Assume draws and losses to be equally bad
- How to find imperfections in opponent and win?

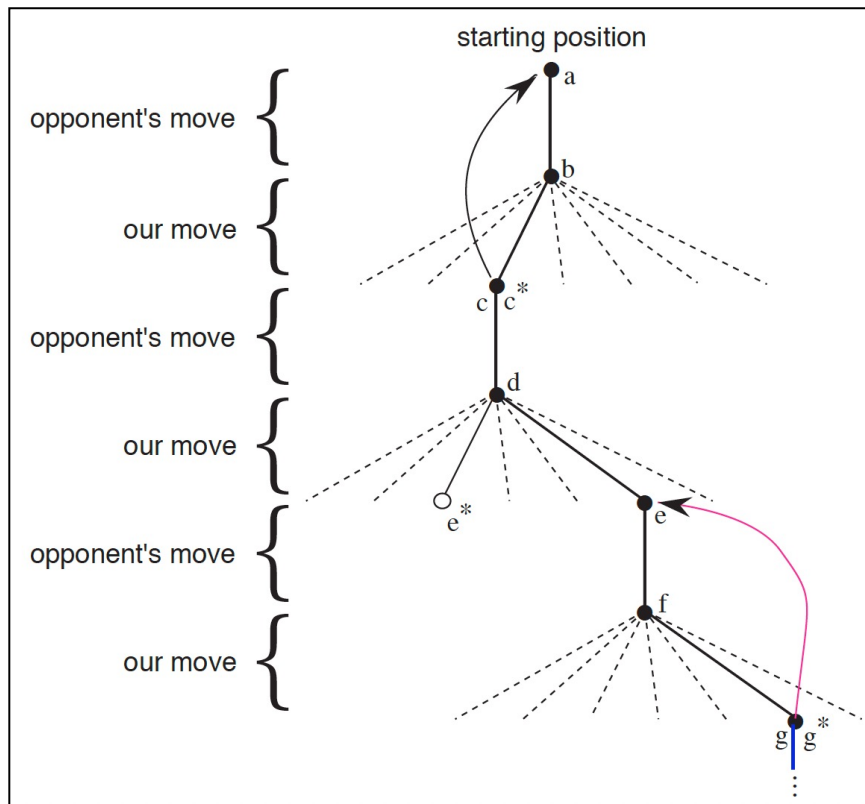
# Classical solutions fail

- Minimax assumes a particular way of playing by opponent
- Dynamic programming requires probabilities of moves opponent makes
- Evolutionary methods searches space of possible policies

# RL solution

- Use a value function
  - A table of numbers, one for each possible state
  - Each no (value) is latest estimate of probability of win from that state
    - Value for Xs in a row is 1, three Os in a row is 0, rest of states is 0.5
- Play many games
  - Most of the time select greedily moving to state with greatest value, select randomly from among other states

# Sequence of tic-tac-toe moves



Solid lines: moves taken during a game

Dashed lines: moves considered but didn't make

Our 2<sup>nd</sup> move was exploratory, even though one leading to  $e^*$  was ranked higher

Curved lines represents estimated values moving up the tree from later nodes (more accurate estimates of the probabilities of winning)

# Updating Value of state

- Move earlier state's value a fraction of the way toward the value of the later state
- Let  $s$  denote state before greedy move, and  $s'$  state after the move, then update to value of  $s$ ,
$$V(s) \leftarrow V(s) + \alpha[V(s') - V(s)]$$
  - $\alpha$  is step-size parameter for rate of learning
- An example of temporal difference learning
  - Changes are based on difference,  $V(s') - V(s)$ , between estimates at two different times

# Tabular vs Approximate Solutions

- Tabular
  - state and action spaces are small enough for value functions to be represented as tables
  - Exact solutions can be found
- Approximate solutions
  - Can be applied to larger class of problems



# Tabular Solutions

- *A single state problem*: bandit problems
- Those based Finite MDPs
  1. Dynamic Programming
    - Require a complete model of the environment
  2. Monte Carlo Methods
    - No model required
    - Unsited to incremental computation
  3. Temporal Difference Methods
    - Require no model, are incremental
    - More difficult to analyze
- Combining Models