

Reinforcement Learning

Reinforcement Learning

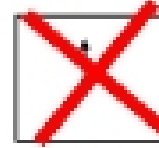
- R&N Chapter 21

- Demos and Data Contributions from
Vivek Mehta (vivekm@cs.cmu.edu)
Rohit Kelkar (ryk@cs.cmu.edu)

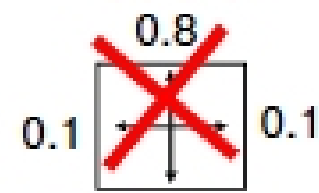
Reinforcement Learning

	1	2	3	4
3				✗
2		✗		✗
1				

Intended action a:



$T(s,a,s')$



- Same (fully observable) MDP as before except:
 - We don't know the model of the environment
 - We don't know $T(.,.,.)$
 - We don't know $R(.)$
- Task is still the same:
 - Find an optimal policy

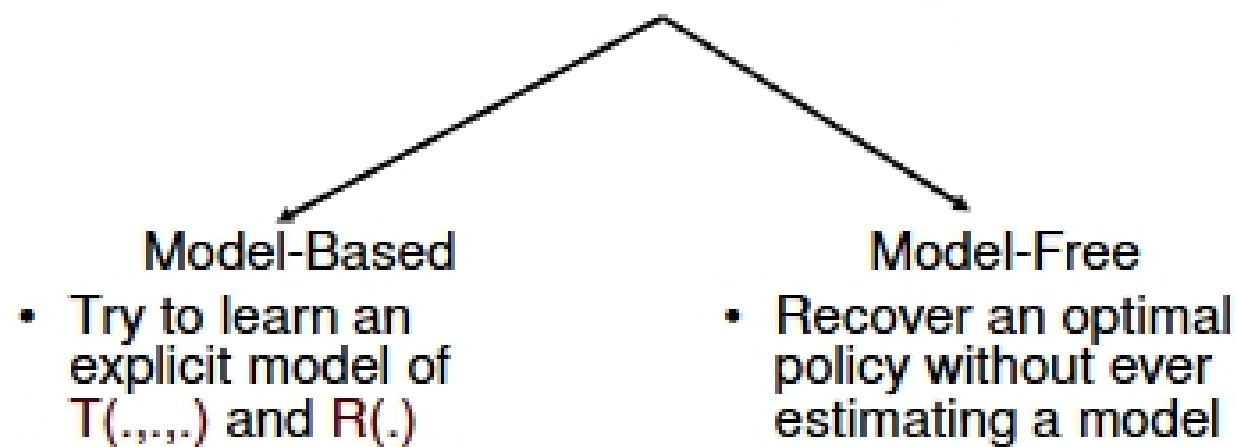
General Problem

- All we can do is try to execute actions and record the resulting rewards
 - World: You are in state 102, you have a choice of 4 actions
 - Robot: I'll take action 2
 - World: You get a reward of 1 and you are now in state 63, you have a choice of 3 actions
 - Robot: I'll take action 3
 - World: You get a reward of -10 and you are now in state 12, you have a choice of 4 actions
 -

Learning from experience

Classes of Techniques

Reinforcement Learning



Model-Based

- If we knew a good estimate $T^{\text{est}}(.,.,.,.)$ of $T(.,.,.,.)$ and $R(.)$, we could evaluate the optimal policy by solving the fundamental MDP relations:

$$U^{\text{est}}(s) = R(s) + \gamma \max_a \left(\sum_{s'} T^{\text{est}}(s, a, s') U^{\text{est}}(s') \right)$$

$$\pi^*(s) = \operatorname{argmax}_a \left(\sum_{s'} T^{\text{est}}(s, a, s') U^{\text{est}}(s') \right)$$