

# RL Reinforcement control

Many RL algorithms:

- Q-Learning } Value Iteration
- TD( $\lambda$ ) } Value Iteration
- Programs } Policy search

## Policy search:

Idea: optimize value function by searching in the space of policies

Traditionally:

$$V(s) = \max_a R(s, a) + \gamma \sum_{s'} P(s' | s, a) V(s')$$

Alternatively:

$$V(\pi) = R(s, \pi(s)) + \gamma \sum_{s'} P(s' | s, \pi(s)) [R(s', \pi(s')) + \gamma \sum_{s''} P(s'' | s', \pi(s')) \dots]$$

Optimize  $V(\pi)$  by gradient ascent on the parameters;  $\pi(s)$

In RL  $P(s'_i | s_i, a_i)$  is unknown but we get to observe trajectories  $\langle s_1, a_1, s'_1, a'_1, s''_1, a''_1, \dots \rangle$

$$V(\pi) = \sum_{i=1}^M R(s_i, \pi(s_i)) + \gamma R(s'_i, \pi(s'_i)) + \gamma^2 R(s''_i, \pi(s''_i)) + \dots$$

Problem: we can't optimize this expression since it changes for every policy  $\pi$ . In other words the sequence  $\rho, \rho', \rho''$  depends on the actual policy,

Solution: express  $V$  in terms of  $\pi$  & a sequence of random numbers. Instead of recording state sequences, record random number sequences.