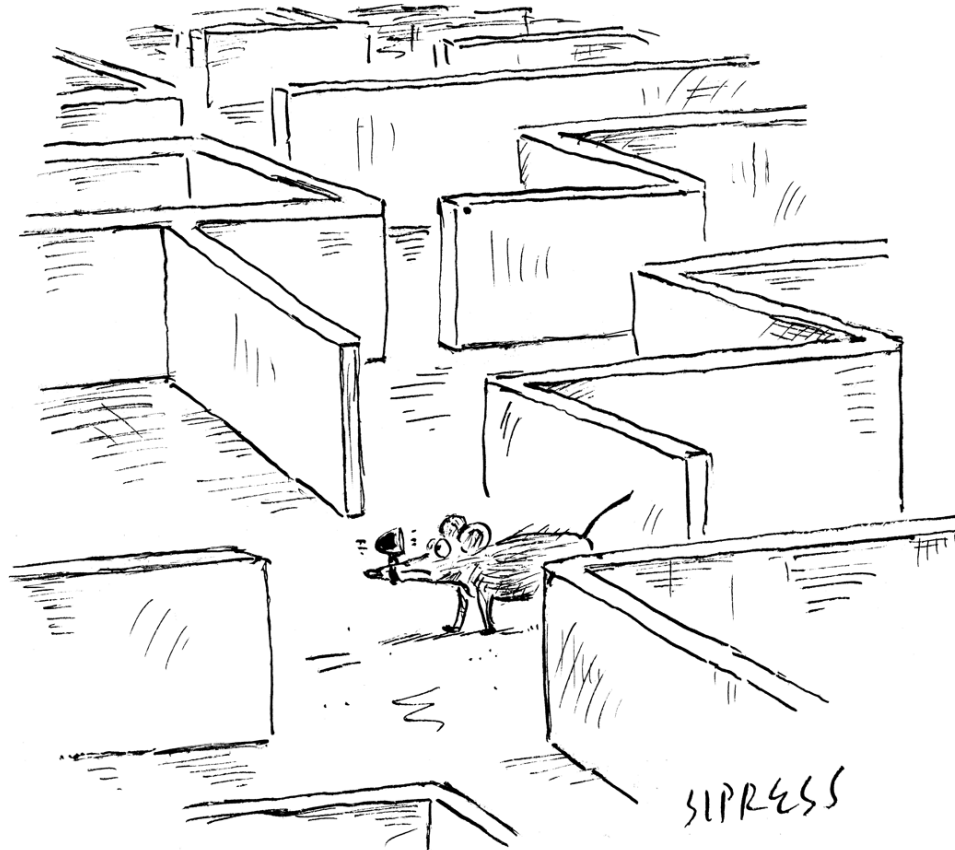


# Q-learning



*“Recalculating ... recalculating ...”*

CSC411/2515: Machine Learning and Data Mining, Winter 2018

Michael Guerzhoy and Lisa Zhang

# Recall Terminology

- State:  $S$ , Action:  $A$ , Reward:  $r$
- Policy:  $\pi_{\theta}(s, a)$
- Return:  $G_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$
- Value Function:

$$V^{\pi_{\theta}}(s) = E[G_t | S_t = s]$$

- Action-Value Function:

$$q^{\pi_{\theta}}(a, s) = E[G_t | S_t = s, A_t = a]$$

- Relationship:

$$V^{\pi_{\theta}}(s) = \sum_a \pi_{\theta}(a|s) q^{\pi_{\theta}}(a, s)$$

# Q-Learning

- Learn the Action-Value Function:

$$q^{\pi_{\theta}}(a, s) = E[G_t | S_t = s, A_t = a]$$

- Q-Learning: Learn the q function!

- Then, define policy to be

$$\pi_{\theta}(s, a) = \operatorname{argmax}_a q^{\pi_{\theta}}(a, s) \text{ or}$$
$$\pi_{\theta}(s, a) \propto q^{\pi_{\theta}}(a, s)$$

- Or use an epsilon-greedy policy:

- Choose  $\pi_{\theta}(s, a) = \operatorname{argmax}_a q^{\pi_{\theta}}(a, s)$  most of the time
- Choose a random action some of the time

# Bellman Equation (1)

- The Value Function can be decomposed:

$$V^{\pi_{\theta}}(s) = E[G_t | S_t = s]$$

$$V^{\pi_{\theta}}(s_t)$$

$$= E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | S_t = s]$$

$$= E[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) | S_t = s]$$

$$= E[r_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= E[r_{t+1} + \gamma E[V^{\pi_{\theta}}(S_{t+1})] | S_t = s]$$

# Bellman Equation (2)

- The Action-Value Function can also be decomposed:

$$q^{\pi\theta}(a, s) = E[G_t | S_t = s, A_t = a]$$

$$q^{\pi\theta}(a, s)$$

$$= E[r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | S_t = s, A_t = a]$$

$$= E[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \dots) | S_t = s, A_t = a],$$

$$= E[r_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$= E[r_{t+1} + \gamma E[V^{\pi\theta}(S_{t+1})] | S_t = s, A_t = a]$$

$$= E[r_{t+1} + \gamma E[q^{\pi\theta}(A_{t+1}, S_{t+1})] | S_t = s, A_t = a]$$

# Optimal Policy

- Suppose we have an optimal policy  $\pi^*$ , then we should have the following Bellman Equations

$$V^*(s) = \max_a q^*(a, s)$$

$$q^*(a, s) = E[r_{t+1} | S_t = s, A_t = a] + \gamma E \left[ \max_{a'} q^*(a', s) | S_t = s, A_t = a \right]$$

For small problems where:

- There are a small number of discrete states
- We know the state transition probabilities  $P(S_{t+1} | s_t, a)$

We can solve this Bellman Equation directly.

# Q-Learning Intuition

- Simple algorithm to find the optimal policy without knowing the state transition probabilities (known as the **model**)
- Idea: Learn a  $q$  function by training the function to satisfy the Bellman Equation

$$q^*(a, s) = E[r_{t+1} | S_t = s, A_t = a] + \gamma E \left[ \max_{a'} q^*(a', s_{t+1}) | S_t = s, A_t = a \right]$$

For a sample  $s_t, a_t, r_{t+1}, s_{t+1}$  from the environment. For the optimal  $q$  function we should have:

$$q^*(a_t, s_t) = r_{t+1} + \gamma \max_{a'} q^*(a', s_{t+1})$$

# Q-Learning Intuition

For a sample  $s_t, a_t, r_{t+1}, s_{t+1}$  from the environment, and a q-function:

$$Loss = L\left(r_{t+1} + \gamma \max_{a'} q(a', s_{t+1}) - q(a_t, s_t)\right)$$

So, train a q-function with gradient descent!

L = usually L2 loss

$$Loss = \left(r_{t+1} + \gamma \max_{a'} q(a', s_{t+1}) - q(a_t, s_t)\right)^2$$



# Q-Learning Algorithm

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$ ;

until  $S$  is terminal

# Benefits

- **Off-policy method:** samples don't have to be from the current policy (though it helps)
- Don't need to wait until episode is finished to train!
- “Learn a guess from a guess”: Q-learning is just one algorithm in a family of algorithms that use this idea

# Q-Learning in practise

- Sample need to be diverse enough to see everything
- **Replay buffer:** sample  $(s_t, a_t, r_{t+1}, s_{t+1})$  put in replay buffer, take a batch from replay buffer to train (**Priority replay:** re-sample  $(s_t, a_t, r_{t+1}, s_{t+1})$  with a large error)
- Takes longer to converge than policy gradient
- Even when the policy converged, the q-function might not

# Representing $q(s, a)$

- Lookup table
  - Learn the value of  $q(s, a)$  directly
  - Works if the states and actions are discrete, and there are few of them
- $q_{\theta}(s, a)$ 
  - $q_{\theta}$  could be a deep (or shallow) neural network
  - Want to minimize
$$Loss = L(r_{t+1} + \gamma \max'_a q_{\theta}(a', s_{t+1}) - q_{\theta}(a_t, s_t))$$
  - Tricky to do directly with gradient descent, but can try to approximate

# Where to go from here?

- Reinforcement Learning (Sutton & Barto)
- David Silver's Video Lectures:  
<http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching.html>