# Q-Learning

## CSE-5364 - FALL 2017

**Learning to balance an inverted pendulum**

**Q-learning** is a model-free reinforcement technique. Specifically, Q-learning can be used to find an optimal action-selection policy for any given (finite) Markov Decision Process. It works by learning an action value function that ultimately gives the expected utility of taking a given action in each state and following the optimal policy thereafter. -wikipedia

**The following equation has been used for solving the problem -**

$$Q^*_{t+1}(s, a) = (1-\alpha) \times Q^*_t(s, a) + \alpha \times (r(s) + \gamma \times max_b Q^*_t(s', b))$$

where

$Q_{t+1}$: optimal Q value of current state

**s**: current state

**a**: current action

**α**: learning rate

$Q_t$: optimal Q value of previous state

**r(s)**: Reward associated with current state

**γ**: Discount factor

**s'**: next state

**b**: all possible actions

The code is attached alongside.


Below are the results of the learning curve obtained after implementing the above equation

Homework 3 Q-Learning

Learning Curve α: 0.1 γ: 0.97



Learning Curve α: 0.1 γ: 0.99



Learning Curve α: 0.3 γ: 0.99

Homework 3 Q-Learning