

## About this class

Back to MDPs

What happens when we don't have complete knowledge of the environment?

Monte-Carlo Methods

Temporal Difference Methods

Function Approximation

## An Example

Blackjack: Goal is to obtain cards whose sum is as great as possible without exceeding 21. All face cards count as 10, and an Ace can be worth either 1 or 11.

Game proceeds as follows: two cards are dealt to both the dealer and the player. One of the dealer's cards is facedown and the other one is faceup. If the player immediately has 21, the game is over, with the player winning if the dealer has less than 21, and the game ending in a draw otherwise.

Otherwise, the player continues by choosing whether to *hit* (get another card) or *stick*. If the total exceeds 21 she goes bust and loses. Otherwise when she sticks, the dealer starts playing using a fixed strategy – she sticks on any sum of 17 or greater. IF the dealer goes

bust the player wins, otherwise the winner is determined by who has a sum closer to 21.

Assume cards are dealt from an infinite deck (i.e. with replacement)

Formulation as an MDP:

1. Episodic, undiscounted
2. Rewards of +1 (winning), 0 (draw), -1 (losing)
3. Actions: hit, stick
4. State space: determined by
  - (a) Player's current sum (12-21, because player always hits below 12)
  - (b) Presence of a *usable* ace (that doesn't have to be counted as 1)
  - (c) Dealer's faceup card

Total of 200 states

Problem: find the value function for a policy that always hits unless the current total is 20 or 21.

Suppose we wanted to apply a dynamic programming method. We would need to figure out all the transition and reward probabilities! This is not easy to do for a problem like this.

Monte Carlo methods can work with sample episodes alone!

It's easy to generate sample episodes for our Blackjack example.

## The Absence of a Transition Model

We now want to estimate *action* values rather than *state* values. So estimate  $Q^\pi(s, a)$

Problem? If  $\pi$  is deterministic, we'll never learn the values of taking different actions in particular states...

Must maintain exploration. This is sometimes dealt with through the concept of *exploring starts* – randomize over all actions at the first state in each episode.

Somewhat problematic assumption – nature won't always be so kind – but it should work OK for Blackjack

In first-visit MC, to evaluate a policy  $\pi$ , we repeatedly generate episodes using  $\pi$ , and then store the return achieved following the *first* occurrence of each state in the episode. Then, averaging these over many simulations gives us the expected value of each state under policy  $\pi$

## Coming Up With Better Policies

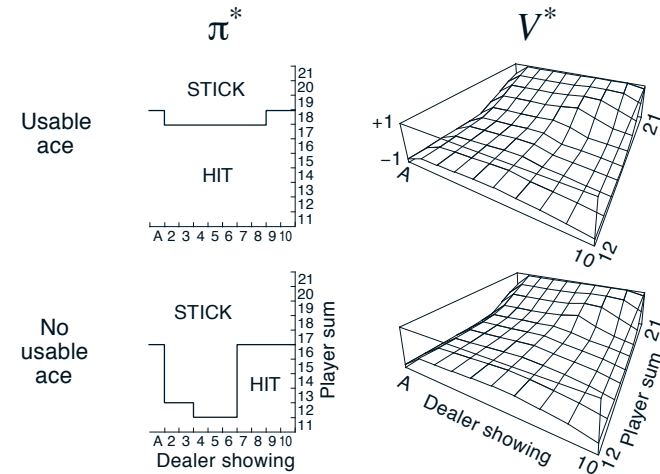
We can interleave policy evaluation with policy improvement as before.

$$\pi_0 \xrightarrow{E} Q^{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} \dots \xrightarrow{I} \pi^* \xrightarrow{E} Q^*$$

We've just figured out how to do policy evaluation.

Policy improvement is even easier because now we have the direct expected rewards for each action in each state  $Q(s, a)$  so just pick the best action among these

The optimal policy for Blackjack:



## On-Policy Learning

On-policy methods attempt to evaluate the same policy that is being used to make decisions

Get rid of the assumption of exploring starts. Now use an  $\epsilon$ -greedy method where some  $\epsilon$  proportion of the time you don't take the greedy action, but instead take a random action

Soft policies: all actions have non-zero probabilities of being selected in all states

For any  $\epsilon$ -soft policy  $\pi$ , any  $\epsilon$ -greedy strategy with respect to  $Q^\pi$  is guaranteed to be an improvement over  $\pi$ .

If we move the  $\epsilon$ -greedy requirement inside the environment, so that we say *nature* randomizes your action  $1 - \epsilon$  proportion of the time, then

the best one can do with general strategies in the new environment is the same as the best one could do with  $\epsilon$ -greedy strategies in the old environment.

## Adaptive Dynamic Programming

Simple idea – take actions in the environment (follow some strategy like  $\epsilon$ -greedy with respect to your current belief about what the value function is) and update your transition and reward models according to observations. Then update your value function by doing full dynamic programming on your current believed model.

In some sense this does as well as possible, subject to the agent's ability to learn the transition model. But it is highly impractical for anything with a big state space (Backgammon has  $10^{50}$  states)

## Temporal-Difference Learning

What is MC estimation doing?

$$V(s_t) \leftarrow (1 - \alpha_t)V(s_t) + \alpha_t R_t$$

where  $R_t$  is the return received following being in state  $s_t$ .

Suppose we switch to a constant step-size  $\alpha$  (this is a trick often used in nonstationary environments)

TD methods basically bootstrap off of existing estimates instead of waiting for the whole reward sequence  $R$  to materialize

$$V(s_t) \leftarrow (1 - \alpha)V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1})]$$

(based on actual observed reward and new state)

This target uses the current value as an estimate of  $V$  whereas the Monte Carlo target

uses the sample reward as an estimate of the expected reward

If we actually want to converge to the optimal policy, the decision-making policy must be GLIE (greedy in the limit of infinite exploration) – that is, it must become more and more likely to take the greedy action, so that we don't end up with faulty estimates (this problem can be exacerbated by the fact that we're bootstrapping)

## Q-Learning: A Model-Free Approach

Even without a model of the environment, you can learn effectively. Q-learning is conceptually similar to TD-learning, but uses the Q function instead of the value function

1. In state  $s$ , choose some action  $a$  using policy derived from current  $Q$  (for example,  $\epsilon$ -greedy), resulting in state  $s'$  with reward  $r$ .

2. Update:

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(r + \gamma \max_{a'} Q(s', a'))$$

You don't need a model for either learning or action selection!

As environments become more complex, using a model can help more (anecdotally)

## Generalization in Reinforcement Learning

So far, we've thought of Q functions and utility functions as being represented by tables

Question: can we parameterize the state space so that we can learn (for example) a linear function of the parameterization?

$$V_{\theta}(s) = \theta_1 f_1(s) + \theta_2 f_2(s) + \dots + \theta_n f_n(s)$$

Monte Carlo methods: We obtain sample of  $V(s)$  and then learn the  $\theta$ 's to minimize squared error.

In general, often makes more sense to use an online procedure, like the Widrow-Hoff rule:

Suppose our linear function predicts  $V_{\theta}(s)$  and we actually would "like" it to have predicted something else, say  $v$ . Define the error as  $E(s) = (V_{\theta}(s) - v)^2/2$ . Then the update rule is:

$$\begin{aligned} \theta_i &\leftarrow \theta_i - \alpha \frac{\partial E(s)}{\partial \theta_i} \\ &= \theta_i + \alpha(v - V_{\theta}(s)) \frac{\partial V_{\theta}(s)}{\partial \theta_i} \end{aligned}$$

If we look at the TD-learning updates in this framework, we see that we essentially replace what we'd "like" it to be with the learned backup (sum of the reward and the value function of the next state):

$$\theta_i \leftarrow \theta_i + \alpha [R(s) + \gamma V_{\theta}(s') - V_{\theta}(s)] \frac{\partial V_{\theta}(s)}{\partial \theta_i}$$

This can be shown to converge to the closest function to the true function when linear function approximators are used, but it's not clear



how good a linear function will be at approximating non-linear functions in general, and all bets on convergence are off when we move to non-linear spaces.

The power of function approximation: allows you to generalize to values of states you haven't yet seen!

In backgammon, Tesauro constructed a player as good as the best humans although it only examined one out of every  $10^{44}$  possible states. Caveat: this is one of the few successes that has been achieved with function approximation and RL. Most of the time it's hard to get a good parameterization and get it to work.