# COMP 652: Machine Learning

Lecture 19

# Today

- Estimating value functions for large state spaces (Via Monte Carlo and supervised learning)
- □ Approximate policy iteration
- □ Relating state values Bellman equations, Bellman optimality equations
- Model-based reinforcement learning
- □ Model-free reinforcement learning



An MDP is defined by:

 $\Box$  Set of <u>states</u> S

 $\Box$  Set of <u>actions</u> A(s) available in each state s

□ <u>*Rewards*</u>:

$$r_{ss'}^a = E\left\{r_{t+1}|s_t = s, a_t = a, s_{t+1} = s'\right\}$$

Transition probabilities

$$p_{ss'}^{a} = P\left(s_{t+1} = s' | s_t = s, a_t = a\right)$$

- $\Box$  A deterministic policy maps every state to an action:  $\pi: S \mapsto A$ .
- The value function of a policy is the expected return the agent receives if following  $\pi$ , conditioned on the state in which the environment starts:

$$V^{\pi}(s) = E\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s, a_{t'} = \pi(s_{t'}) \text{ for all } t' \ge t\}$$
  
=  $E_{\pi}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s\}$ 

The action-value function of a policy the expected return of choosing an action and following the policy afterwards, conditioned on the state in which the environment starts:

$$Q^{\pi}(s,a) = E_{\pi}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s, a_t = a\}$$

(These are for  $\gamma \in (0, 1]$ , though finiteness for  $\gamma = 1$  requires extra conditions.)

# **Policy iteration**

- □ Last time, we describe Monte Carlo methods for estimating  $V^{\pi}$  and  $Q^{\pi}$ □ We showed the policy iteration algorithm, which alternates:
  - Estimating  $Q^{\pi}$
  - Updating  $\pi$  to be "greedy" with respect to  $Q^\pi$
- $\hfill\square$  Policy iteration terminates at a globally optimal policy,  $\pi^*$ 
  - For all s and all  $\pi'$ ,  $V^{\pi^*}(s) \ge V^{\pi'}(s)$
  - For all s, all a, and all  $\pi'$ ,  $Q^{\pi^*}(s,a) \ge Q^{\pi'}(s,a)$
- □ The Monte Carlo approach uses data quite inefficiently...a problem to which we return shortly.
- $\Box$  What if S is too large, even infinite, so that we cannot represent  $V^{\pi}$  explicitly?

### Dealing with large/continuous state spaces

- Consider the problem of  $V^{\pi}$  estimation when |S| is too large / continuous, so that explicit, tabular representation is not possible.
- $\Box$  What can we do?

- $\Box$  Consider a sample trajectory  $s_0, a_0, r_1, s_1, a_1, r_2, \ldots$  taken under policy  $\pi$
- $\Box$  Suppose we represent states by some feature mapping  $\phi(s)$
- □ We can formulate the supervised learning (regression) problem:

input	output
$\phi(s_0)$	$r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$
$\phi(s_1)$	$r_2 + \gamma r_3 + \gamma^2 r_4 + \dots$
$\phi(s_2)$	$r_3 + \gamma r_4 + \gamma^2 r_5 + \dots$
$\phi(s_3)$	$r_4 + \gamma r_5 + \gamma^2 r_6 + \dots$
÷	÷

- We can fit this data with any regressor: linear/polynomial regression, neural network, nearest-neighbor, regression tree, ....
- □ If we have multiple trajectories, each one contributes data similarly.

 $\Box$  Estimating  $Q^{\pi}$  is slightly trickier. How can we do it?

## Estimating $Q^{\pi}$ by supervised learning

 $\Box$  Can define a feature mapping  $\phi(s, a)$ 

□ One way (Monte Carlo exploring starts):

- Generate multiple trajectories,  $au_1$ ,  $au_2$ , ...,  $au_m$
- Trajectory  $\tau_i = \langle s_0^i, a_0^i, r_1^i, s_1^i, a_1^i, r_2^i, \ldots \rangle$  generated by choosing random initial state  $s_0^i$ , random initial action  $a_0^i$ , and following  $\pi$  afterwards
- Supervised learning data is:

input	output
$ \begin{array}{c} \phi(s_0^1,a_0^1) \\ \phi(s_0^2,a_0^2) \\ \phi(s_0^3,a_0^3) \\ \vdots \end{array} $	$ \begin{array}{c} r_1^1 + \gamma r_2^1 + \gamma^2 r_3^1 + \dots \\ r_1^2 + \gamma r_2^2 + \gamma^2 r_3^2 + \dots \\ r_1^3 + \gamma r_2^3 + \gamma^2 r_3^3 + \dots \\ \vdots \end{array} $

Alternatively, one can define a state-feature mapping,  $\phi(s)$ , and learn a separate approximator for each action's value (as a function of state).

□ This enables approximate policy iteration, one version of which is:

- Start with an arbitrary initial policy  $\pi_0$
- Repeat the following steps:
  - 1. Approximate  $Q^{\pi_i}$  for the current policy,  $\pi_i$ , by generated Monte Carlo Exploring Starts (MCES) data, and applying some supervised learning method.
  - 2. Define  $\pi_{i+1}(s)$  implicitly as  $\arg \max_a Q^{\pi_i}(s, a)$
- □ Does it converge?
- □ Does it converge to an optimal policy?

□ This enables approximate policy iteration, one version of which is:

- Start with an arbitrary initial policy  $\pi_0$
- Repeat the following steps:
  - 1. Approximate  $Q^{\pi_i}$  for the current policy,  $\pi_i$ , by generated Monte Carlo Exploring Starts (MCES) data, and applying some supervised learning method.
  - 2. Define  $\pi_{i+1}(s)$  implicitly as  $\arg \max_a Q^{\pi_i}(s, a)$
- $\Box$  Does it converge?
- ⇒ In general, no. For certain kinds of MDPs and/or function approximators, though, it does.
- $\Box$  Does it converge to an optimal policy?
- ⇒ In general, no. Though one can bound degree of suboptimality in terms of error in  $Q^{\pi}$  approximation.

- $\Box = V^{\pi}(s)$  an expected sum of future rewards
- But intuitively, a state's value must by related to the states to which it leads



#### **Relationships between state values (II)**

$$V^{\pi}(s) = E_{\pi}\{r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \dots | s_{t} = s\}$$
  
=  $\sum_{s'} P(s_{t+1} = s' | s_{t} = s) E_{\pi}\{r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \dots | s_{t} = s, s_{t+1} = s\}$   
=  $\sum_{s'} p_{ss'}^{\pi(s)} [r_{ss'}^{\pi(s)} + \gamma E_{\pi}\{r_{t+2} + \gamma r_{t+3} + \dots | s_{t} = s, s_{t+1} = s'\}]$   
=  $\sum_{s'} p_{ss'}^{\pi(s)} [r_{ss'}^{\pi(s)} + \gamma V^{\pi}(s')]$ 

- □ These are called the *Bellman equations* for policy value
- Treating the  $V^{\pi}(s)$  as a set of |S| variables, the above gives a <u>linear</u> system of |S| equations with |S| unknowns.
- $\Rightarrow We can use linear system solvers to find <math>V^{\pi}(s)$ , if we know  $p^{a}_{ss'}$  and  $r^{a}_{ss'}$ . (There's other ways to do it too, though.)
- $\Box$  Are solutions unique?

□ Let  $r^{\pi}$  be a the expected immediate reward upon following policy  $\pi$  in state s:

$$r^{\pi}(s) = \sum_{s'} p_{ss'}^{\pi(s)} r_{ss'}^{\pi(s)}$$

We will treat  $r^{\pi}$  as a length |S| column vector

Let  $P^{\pi}(s, s') = p_{ss'}^{\pi(s)}$  be the probability that state s' follows state s, when agent acts according to  $\pi$ .

We will treat this as an  $|S| \times |S|$  matrix.

 $\Box$  Then considering  $V^{\pi}$  as a length |S| column vector, the equation of the previous slide says that:

$$V^{\pi} = r^{\pi} + \gamma P^{\pi} V^{\pi}$$

Rearranging, we find that:

 $\square$ 

$$(I - \gamma P^{\pi})V^{\pi} = r^{\pi}$$

- $\Box$   $P^{\pi}$  is a Markov matrix. If it describes a non-terminating chain, its largest eigenvalue is 1. If it describe a terminating chain, its largest eigenvalue is < 1.
- $\Box \quad \text{If } P^{\pi} \text{ is non-terminating and } \gamma < 1 \text{, or if } P^{\pi} \text{ is terminating, then the largest eigenvalue of } (I \gamma P^{\pi}) \text{ is } < 0.$
- $\Box$  Thus, it is invertible, and  $V^{\pi}$  is uniquely determined as:

$$V^{\pi} = (I - \gamma P^{\pi})^{-1} V^{\pi}$$

- Matrix inversion is computationally expensive
   Iterative approaches are more commonly used for computing value functions, e.g.:
  - Initialize  $V^{\pi}(s)$  arbitrarily
  - Repeat, until changes are sufficiently small:

$$V^{\pi}(s) \leftarrow \sum_{s'} p_{ss'}^{\pi(s)}(r_{ss'}^{\pi(s)} + \gamma V^{\pi}(s')) \text{ for all } s$$

This can be shown to converge exponentially quickly to the correct values.
 (Aside: What does "exponentially quickly" mean?)

# Example

- □ Imagine the simple deterministic "maze" below
- □ Goal is to get into the dotted square in the lower right room
- $\Box$  Reward of +1 upon arrival to goal
- $\hfill\square$  Discount factor  $\gamma=0.9$
- $\square$   $\pi$  is correct, optimal policy
- $\square$  Initialize  $V^{\pi}(s) = 0$  for all s, except goal, which is set to +1
- □ Figure shows values by dots of radius proportional to value



#### **Action-value functions**

□ Similar reasoning holds for action-value functions. □ First, note that  $V^{\pi}(s) = Q^{\pi}(s, \pi(s))$ . Then:

$$Q^{\pi}(s,a) = E_{\pi}\{r_{t+1} + \gamma r_{t+2} + \gamma^{2} r_{t+3} + \dots | s_{t} = s, a_{t} = a\}$$
  

$$= \sum_{s'} p_{ss'}^{a} [r_{ss'}^{a} + \gamma E_{\pi}\{r_{t+2} + \gamma r_{t+3} + \dots | s_{t} = s, a_{t} = a, s_{t+1} = s]$$
  

$$= \sum_{s'} p_{ss'}^{a} [r_{ss'}^{a} + \gamma V^{\pi}(s')]$$
  

$$= \sum_{s'} p_{ss'}^{a} [r_{ss'}^{a} + \gamma Q^{\pi}(s', \pi(s'))]$$

U We have  $\sum_{s} |A(s)|$  linear equations in the same number of unknowns. A similar matrix argument shows the solution,  $Q^{\pi}$ , is unique.

□ An iterative-style algorithm can be used to compute it:

$$Q^{\pi}(s,a) \leftarrow \sum_{s'} p^a_{ss'}(r^a_{ss'} + \gamma Q^{\pi}(s',\pi(s')))$$

# Finding optimal policies

- $\hfill\square$  The previous techniques can be used to implement exact policy iteration
- However, there is another approach based on *Bellman optimality* equations (see below)
- □ Recall that there is at least one optimal policy  $\pi^*$ . (It satisfies  $V^{\pi^*}(s) \ge V^{\pi}(s)$  and  $Q^{\pi^*}(s, a) \ge Q^{\pi}(s, a)$  for all  $\pi, s, a$ .)
- $\Box$  For an optimal policy  $\pi^*$ , we must have:

$$V^{\pi^{*}}(s) = \max_{a} Q^{\pi^{*}}(s, a)$$
  
= 
$$\max_{a} \left[ \sum_{s'} p^{a}_{ss'}(r^{a}_{ss'} + \gamma V^{\pi^{*}}(s')) \right]$$
$$Q^{\pi^{*}}(s, a) = \sum_{s'} p^{a}_{ss'}(r^{a}_{ss'} + \gamma V^{\pi^{*}}(s'))$$
  
= 
$$\sum_{s'} p^{a}_{ss'}(r^{a}_{ss'} + \gamma \max_{a'} Q^{\pi^{*}}(s', a'))$$

These are *nonlinear* systems of equations

 $\square$ 

## Solving the Bellman optimality equations

Solutions can be found by linear programming
 More common, however, is *value iteration*:

More common, however, is *value iteration*:

$$V(s) \leftarrow \max_{a} \left[ \sum_{s'} p^a_{ss'} (r^a_{ss'} + \gamma V^{\pi^*}(s')) \right]$$

Or action-value iteration:

$$Q(s, a) \leftarrow \sum_{s'} p^a_{ss'}(r^a_{ss'} + \gamma \max_{a'} Q^{\pi^*}(s', a'))$$

□ These approaches converge exponentially quickly to the optimal value function  $V^*$  or action-value function  $Q^*$  (under the same conditions needed for well-definedness)

#### Back to reinforcement learning

 $\begin{tabular}{ll} \hline & \end{tabular} \begin{tabular}{ll} \hline & \end{t$ 

- □ The previous discussion assumes that we know  $p^a_{ss'}$  and  $r^a_{ss'}$ □ What if we don't?
  - Model-based value function learning
  - Model-free value function learning
  - Value function-free learning

- Model-based learning algorithms use experience from the environment to build an <u>approximate model</u>  $\hat{r}^a_{ss'}$ ,  $\hat{p}^a_{ss'}$
- Then we pretend the approximate model is correct and use it to compute the optimal value function/policy as above
- □ How do we estimate the rewards and transition probabilities?
  - If |S| is small enough for tabular representation,  $\hat{r}^a_{ss'}$  can be estimated as the empirical mean reward from every s - a - r - s' quadruple in the data.  $\hat{p}^a_{ss'}$  can be taken to be the empirical probability that s'follows s - a in the data.
  - Otherwise, solutions are ad hoc. For low-dimensional continuous state spaces, state aggregation or interpolators are often used.

# Example



- □ Followed Rand policy (equal chance of left or right action)
- $\Box$  10 trajectories
- □ Estimated rewards and transition probabilities

#### Results



 $\square$ 

#### Reward function exactly correct

Transition probabilities somewhat correct:

state $s$	2	3	4	5	6	7	8	9	10
$p_{s,s-1}^{left}$	0.6	0.85714	0.64706	0.95238	0.72	0.78571	0.875	0.8	
$p_{s,s+1}^{left}$	0.4	0.14286	0.35294	0.047619	0.28	0.21429	0.125	0.2	
$p_{s,s-1}^{right}$	0.25	0.2	0.27273	0.26087	0.23529	0.21429	0.083333	0.13333	
$p_{s,s+1}^{right}$	0.75	0.8	0.72727	0.73913	0.76471	0.78571	0.91667	0.86667	

 $Q^*$  estimate a bit off, but  $\pi^*$  estimate correct:

state s	2	3	4	5	6	7	8	9	10
$Q^*(s, left)$	10	9.9999	9.9998	9.9994	9.9926	9.975	9.9104	9.4586	
$Q^*(s, right)$	10	9.9998	9.9996	9.9944	9.9807	9.928	9.5016	7.9525	
$\pi^*(s)$	$\leftarrow$								

- + Tend to be very data-efficient
- + The models may be of independent interest
- + If reward function changes, or dynamics of only small part of the environment changes, lots of information can be reused
- If state set S or action set A are very large or infinite, it will be very hard to estimate the model from data—especially for the transition probabilities; this can lead to poor performance.
- Even if model is accurate, solving for  $\pi^*$  can be nontrivial. (It's polynomial in |S| and |A|, but if large, can be problematic.)

- $\hfill\square$  Suppose that we represented every state s with a  $\underline{\textit{feature vector}} \ \phi_s,$  of size  $k \leq |S|$
- $\Box$  We can represent all the feature vectors, for all the state, in a <u>feature matrix</u>  $\Phi$ , of size  $k \times |S|$ , where the sth column is  $\phi_s$
- Important special case: if each column has exactly one element equal to 1 and all the others are 0, the matrix represents a <u>state partition</u>, where the state space has been partitioned in k disjoint subsets.
- In general, the features (also called <u>basis functions</u>) can be anything (Gaussian, sine-cosine, etc)

□ Consider the Bellman equations:

$$\mathbf{V}^{\pi} = \mathbf{r}^{\pi} + \gamma \mathbf{P}^{\pi} \mathbf{V}^{\pi}$$

where  $\mathbf{V}^{\pi}$  is a column vector representing  $V^{\pi}(s)$ ,  $\mathbf{r}^{\pi}$  is the expected rewards following each state, and  $\mathbf{P}^{\pi}$  is the matrix containing  $p_{ss'}^{\pi}(s)$ .

 $\square$  We multiply at the left by  $\Phi$ :

$$\mathbf{\Phi}\mathbf{V}^{\pi} = \mathbf{\Phi}\mathbf{r}^{\pi} + \gamma\mathbf{\Phi}\mathbf{P}^{\pi}\mathbf{V}^{\pi}$$

 $\Box$  We make  $\mathbf{\Phi}\mathbf{V}^{\pi}$  appear on the right hand side as well:

$$\mathbf{P}^{\pi}\mathbf{V}^{\pi} = \mathbf{P}^{\pi}\mathbf{I}\mathbf{V}^{\pi} = \mathbf{P}^{\pi}\mathbf{\Phi}^{T}\mathbf{\Phi}\mathbf{V}^{\pi}$$

□ Now we can re-write the Bellman equations:

$$\mathbf{\Phi}\mathbf{V}^{\pi} = \mathbf{\Phi}\mathbf{r}^{\pi} + \gamma\mathbf{\Phi}\mathbf{P}^{\pi}\mathbf{\Phi}^{T}\mathbf{\Phi}\mathbf{V}^{\pi} \Rightarrow (\mathbf{I} - \gamma\mathbf{\Phi}\mathbf{P}^{\pi}\mathbf{\Phi}^{T})\mathbf{\Phi}\mathbf{V}^{\pi} = \mathbf{\Phi}\mathbf{r}^{\pi}$$

 $\Box$  We re-write the above equation as:

$$\mathbf{\Phi}\mathbf{V}^{\pi} = (\mathbf{I} - \gamma\mathbf{\Phi}\mathbf{P}^{\pi}\mathbf{\Phi}^{T})^{-1}\mathbf{\Phi}\mathbf{r}^{\pi}$$

(I'm glossing over assumption needed to ensure the inverse exists.)

- $\Box$  Let  $\hat{\mathbf{r}}^{\pi} = \mathbf{\Phi}\mathbf{r}^{\pi}$ ; this is a vector of size k, representing the reward for every feature
- E.g., in the special case of state partitioning, the reward associated with a partition will be the sum of the rewards for the states in that partition (why?)
- Let  $\hat{\mathbf{P}}^{\pi} = \mathbf{\Phi} \mathbf{P}^{\pi} \mathbf{\Phi}^{T}$ ; this is a  $k \times k$  matrix showing transitions between features
- Since this is typically much smaller than the original matrix, it can be estimated more accurately with less data

$$\mathbf{\Phi}\mathbf{V}^{\pi} = (\mathbf{I} - \gamma\mathbf{\Phi}\mathbf{P}^{\pi}\mathbf{\Phi}^{T})^{-1}\mathbf{\Phi}\mathbf{r}^{\pi}$$

- Let  $\hat{\mathbf{V}}^{\pi} = \mathbf{\Phi} \mathbf{V}^{\pi}$ ; this is the approximation of the value function using the features
- E.g., in the case of a state partition, each partition will have a value associated with it, and all states in the partition share the same value
   Obviously, not all value functions can be represented correctly anymore.
   The Bellman equations for approximate values become:

$$\hat{\mathbf{V}}^{\pi} = (\mathbf{I} - \gamma \hat{\mathbf{P}}^{\pi}) \hat{\mathbf{r}}^{\pi}$$

### Trade-off

- $\hfill\square$  The above systems has k equations with k unknowns
- $\square$  Model-based approximate methods will estimate  $\hat{\mathbf{r}}^{\pi}$  and  $\hat{\mathbf{P}}^{\pi}$  from data
- $\Box$  The smaller k is, the less data we need to do this estimation, and the easier it is to solve the system
- $\Box$  But the smaller k is, the less accurate will the value function be
- Instead of estimating  $V^{\pi}$ ,  $Q^{\pi}$  can be estimated, leading to an alternative approximate policy iteration algorithm

## Modeling value, not dynamics

- $\hfill\square$  In model-free value function-based RL, we directly estimate value function, but not  $r^a_{ss'}$  or  $p^a_{ss'}$
- The Monte Carlo methods describe last lecture and at the start of this lecture are one way to do so
- However, the iterative (dynamic programming) approaches to value function computation provide inspiration for another class of approaches

 $\Box \quad \text{Consider a trajectory, with actions selected according to policy } \pi:$  $\cdots \quad \underbrace{s_{t}}_{\mathscr{A}_{t}} \bullet \underbrace{f_{t+1}}_{\mathscr{A}_{t+1}} \underbrace{s_{t+2}}_{\mathscr{A}_{t+2}} \underbrace{s_{t+3}}_{\mathscr{A}_{t+2}} \underbrace{s_{t+3}}_{\mathscr{A}_{t+3}} \cdots$ 

 $\Box$  The Bellman equation is:

$$V^{\pi}(s_t) = E_{\pi} \left[ r_{t+1} + \gamma V^{\pi}(s_{t+1}) | s_t \right]$$

which suggested the dynamic programming update:

$$V(s_t) \leftarrow E_{\pi} \left[ r_{t+1} + \gamma V(s_{t+1}) | s_t \right]$$

- In general, we do not know this expected value, but we do have an possibly-biased sample estimate of it,  $r_{t+1} + \gamma V(s_{t+1})$
- $\Box$  We can make an update <u>towards</u> the sample value, with step size  $\alpha$ :

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

We can rewrite the previous as:

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

- □ The term after the  $\alpha$  is called the "temporal difference" it is the difference between what our estimate  $V(s_t)$  suggested we would see, and  $r_{t+1} + \gamma V(s_{t+1})$
- $\Box$  Does it converge to  $V^{\pi}$ ?

 $\square$ 

We can rewrite the previous as:

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

- The term after the  $\alpha$  is called the "temporal difference" it is the difference between what our estimate  $V(s_t)$  suggested we would see, and  $r_{t+1} + \gamma V(s_{t+1})$
- $\Box$  Does it converge to  $V^{\pi}$ ? Yes! If:
  - We have infinitely much data collected under policy  $\pi$
  - If all states s appear infinitely often in the data
  - If learning rate(s)  $\alpha$  decrease towards zero at an appropriate rate (Robbins-Monroe conditions)

 $\square$ 

- $\Box$  That's fine for learning  $V^{\pi}.$  A similar procedure can be designed to learn  $Q^{\pi}.$
- □ What about learning optimal policies?
- Suppose we generate experience  $(s_0, a_0, r_1, s_1, a_1, r_2, ...)$  from the environment, and update an action-value function as:

$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha(r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a'))$$

This is motivated by the Bellman optimality equation:

$$Q^*(s,a) = \sum_{s'} p^a_{ss'}(r^a_{ss'} + \gamma \max_{a}' Q^*(s',a'))$$

 $\Box$  This is called *Q*-Learning

- □ Does Q-Learning converge? Yes!
- The experience can be generated under any policy  $\pi$  at all or not even according to a policy, strictly speaking
- $\Box$  We need infinitely much data
- □ Every possible state-action pair must occur infinitely many times
- □ Learning rates need to be scheduled appropriately