## Lecture 5: Value Function Approximation

Emma Brunskill

CS234 Reinforcement Learning.

Winter 2022

The value function approximation structure for today closely follows much of David Silver's Lecture 6.
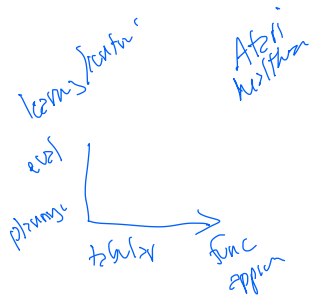
# L5 Refresh Your Knowledge

- The basic idea of TD methods are to make state-next state pairs fit the constraints of the Bellman equation on average (question by: Phil Thomas)
  1. True
  2. False
  3. Not sure

- In tabular MDPs, if using a decision policy that visits all states an infinite number of times, and in each state randomly selects an action, then (select all)
  1. Q-learning will converge to the optimal Q-values
  2. SARSA will converge to the optimal Q-values
  3. Q-learning is learning off-policy
  4. SARSA is learning off-policy
  5. Not sure

- A TD error $> 0$ can occur even if the current $V(s)$ is correct $\forall s$: [select all]
  1. False
  2. True if the MDP has stochastic state transitions
  3. True if the MDP has deterministic state transitions
  4. Not sure

# L5 Refresh Your Knowledge

- The basic idea of TD methods are to make state-next state pairs fit the constraints of the Bellman equation on average (question by: Phil Thomas)
  True

- In tabular MDPs, if using a decision poilcy that visits all states an infinite number of times, and in each state randomly selects an action, then (select all)

  1. Q-learning will converge to the optimal Q-values (True)
  2. SARSA will converge to the optimal Q-values (False)
  3. Q-learning is learning off-policy (True)
  4. SARSA is learning off-policy (False)

- A TD error $> 0$ can occur even if the current $V(s)$ is correct $\forall s$: [select all]

  1. False
  2. True if the MDP has stochastic state transitions (True)
  3. True if the MDP has deterministic state transitions (False)
  4. Not sure

- 

learn/eval/contr'

Afer
hu/lthan

eval

plann'

$\rightarrow$ tabular $\rightarrow$ func appln

# Maximization Bias[1] Proof

- Consider single-state MDP ($|S| = 1$) with 2 actions, and both actions have 0-mean random rewards, ($\mathbb{E}(r|a = a_1) = \mathbb{E}(r|a = a_2) = 0$).
- Then $Q(s, a_1) = Q(s, a_2) = 0 = V(s)$
- Assume there are prior samples of taking action $a_1$ and $a_2$
- Let $\hat{Q}(s, a_1), \hat{Q}(s, a_2)$ be the finite sample estimate of $Q$
- Use an unbiased estimator for $Q$: e.g. $\hat{Q}(s, a_1) = \frac{1}{n(s,a_1)} \sum_{i=1}^{n(s,a_1)} r_i(s, a_1)$
- Let $\hat{\pi} = \arg\max_a \hat{Q}(s, a)$ be the greedy policy w.r.t. the estimated $\hat{Q}$
- *Even though each estimate of the state-action values is unbiased*, the estimate of $\hat{\pi}$'s value $\hat{V}^{\hat{\pi}}$ can be biased:
$\hat{V}^{\hat{\pi}}(s) = \mathbb{E}[\max \hat{Q}(s, a_1), \hat{Q}(s, a_2)]$
$\geq \max[\mathbb{E}[\hat{Q}(s, a_1)], [\hat{Q}(s, a_2)]]$
$= max[0, 0] = V^{\pi}$,
where the inequality comes from Jensen's inequality.

---

[1]Example from Mannor, Simester, Sun and Tsitsiklis. Bias and Variance Approximation in Value Function Estimates. Management Science 2007

# Double Q-Learning

- The greedy policy w.r.t. estimated $Q$ values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of $Q_1(s_1, a_i)$ and $Q_2(s_1, a_i)$ $\forall a$.
  - Use one estimate to select max action: $a^* = \arg\max_a Q_1(s_1, a)$
  - Use other estimate to estimate value of $a^*$: $Q_2(s, a^*)$
  - Yields unbiased estimate: $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$

# Double Q-Learning

- The greedy policy w.r.t. estimated $Q$ values can yield a maximization bias during finite-sample learning
- Avoid using max of estimates as estimate of max of true values
- Instead split samples and use to create two independent unbiased estimates of $Q_1(s_1, a_i)$ and $Q_2(s_1, a_i)$ $\forall a$.
    - Use one estimate to select max action: $a^* = \arg\max_a Q_1(s_1, a)$
    - Use other estimate to estimate value of $a^*$: $Q_2(s, a^*)$
    - Yields unbiased estimate: $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$
- Why does this yield an unbiased estimate of the max state-action value?
  Using independent samples to estimate the value

- If acting online, can alternate samples used to update $Q_1$ and $Q_2$, using the other to select the action chosen
- Next slides extend to full MDP case (with more than 1 state)

# Double Q-Learning

---

1: Initialize $Q_1(s, a)$ and $Q_2(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
2: **loop**
3:   Select $a_t$ using $\epsilon$-greedy $\pi(s) = \arg\max_a Q_1(s_t, a) + Q_2(s_t, a)$
4:   Observe $(r_t, s_{t+1})$
5:   **if** (with 0.5 probability) **then**
6:     $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma Q_2(s_{t+1}, \arg\max_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t))$
7:   **else**
8:     $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma Q_1(s_{t+1}, \arg\max_a Q_2(s_{t+1}, a)) - Q_2(s_t, a_t))$
9:   **end if**
10:   $t = t + 1$
11: **end loop**

---

2x          same          unclear

Compared to Q-learning, how does this change the: memory requirements, computation requirements per step, amount of data required?
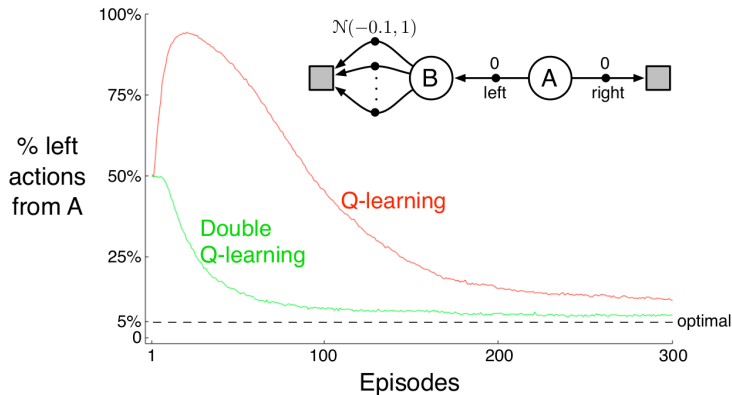
# Double Q-Learning

---

1: Initialize $Q_1(s, a)$ and $Q_2(s, a), \forall s \in S, a \in A$ $t = 0$, initial state $s_t = s_0$
2: **loop**
3:     Select $a_t$ using $\epsilon$-greedy $\pi(s) = \arg\max_a Q_1(s_t, a) + Q_2(s_t, a)$
4:     Observe $(r_t, s_{t+1})$
5:     **if** (with 0.5 probability) **then**
6:         $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma Q_2(s_{t+1}, \arg\max_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t))$
7:     **else**
8:         $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma Q_1(s_{t+1}, \arg\max_a Q_2(s_{t+1}, a)) - Q_2(s_t, a_t))$
9:     **end if**
10:     $t = t + 1$
11: **end loop**

---

Compared to Q-learning, how does this change the: memory requirements, computation requirements per step, amount of data required?

Doubles the memory, same computation requirements, data requirements are subtle– might reduce amount of exploration

Due to the maximization bias, in this case Q-learning spends much more time selecting suboptimal actions than double Q-learning.
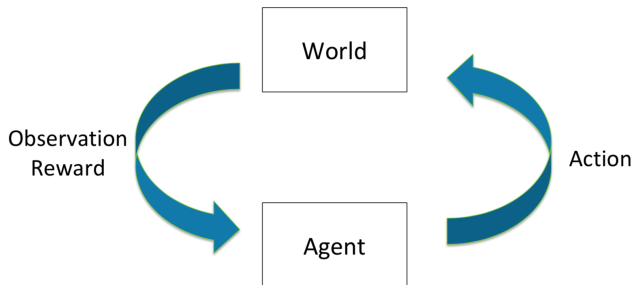
# Class Structure

- Last time: Control (making decisions) without a model of how the world works  *tabular setting*
- **This time: Linear value function approximation**
- Next time: Deep reinforcement learning

- Value function approximation
- Monte Carlo policy evaluation with linear function approximation
- TD policy evaluation with linear function approximation
- Control methods with linear value function approximation

# Reinforcement Learning



- Goal: Learn to select actions to maximize total expected future reward

# Last Time: Tabular Representations for Model-free Control

- Last time: how to learn a good policy from experience
- So far, have been assuming we can represent the value function or state-action value function as a vector/ matrix
    - Tabular representation
- Many real world problems have enormous state and/or action spaces
- Tabular representation is insufficient
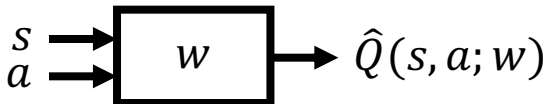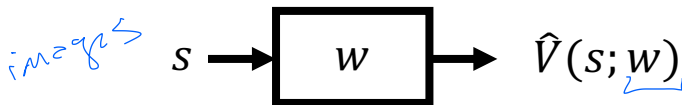
# Motivation for Function Approximation

- Don't want to have to explicitly store or learn for every single state a
  - Dynamics or reward model
  - Value
  - State-action value
  - Policy
- Want more compact representation that generalizes across state or states and actions

# Benefits of Function Approximation

- Reduce memory needed to store $(P, R)/V/Q/\pi$
- Reduce computation needed to compute $(P, R)/V/Q/\pi$
- Reduce experience needed to find a good $P, R/V/Q/\pi$

# Value Function Approximation (VFA)

- Represent a (state-action/state) value function with a parameterized function instead of a table

$$s \longrightarrow \boxed{w} \longrightarrow \hat{V}(s; w)$$

*images*

$$\begin{array}{c} s \\ a \end{array} \longrightarrow \boxed{w} \longrightarrow \hat{Q}(s, a; w)$$

- **Which function approximator?**

# Function Approximators

- Many possible function approximators including
  - Linear combinations of features
  - Neural networks
  - Decision trees
  - Nearest neighbors
  - Fourier/ wavelet bases
- In this class we will focus on function approximators that are differentiable (Why?)
- Two very popular classes of differentiable function approximators
  - Linear feature representations (Today)
  - Neural networks (Next lecture)

- Given known dynamics and reward models, and a tabular representation
  - Discussed how to do policy evaluation and then control (value iteration and policy iteration)
- Given no models, and a tabular representation
  - Discussed how to do policy evaluation (MC/TD) and then control (MC, SARSA, Q-learning)
- **Given no models, and function approximation**
  - **Today will discuss how to do policy evaluation and then control**

# Review: Gradient Descent

- Consider a function $J(\mathbf{w})$ that is a differentiable function of a parameter vector $\mathbf{w}$
- Goal is to find parameter $\mathbf{w}$ that minimizes $J$
- The gradient of $J(\mathbf{w})$ is

$$\nabla J(\omega) = \left[ \frac{\partial J}{\partial \omega_1} \quad \frac{\partial J}{\partial \omega_2} \quad \cdots \quad \frac{\partial J}{\partial \omega_n} \right]$$

$$\omega = \omega - \Delta \omega$$

# Value Function Approximation for Policy Evaluation with an Oracle

- First assume we could query any state $s$ and an oracle would return the true value for $V^\pi(s)$
- Similar to supervised learning: assume given $(s, V^\pi(s))$ pairs
- The objective is to find the best approximate representation of $V^\pi$ given a particular parameterized function $\hat{V}(s; w)$

# Stochastic Gradient Descent

- Goal: Find the parameter vector $\boldsymbol{w}$ that minimizes the loss between a true value function $V^\pi(s)$ and its approximation $\hat{V}(s; \boldsymbol{w})$ as represented with a particular function class parameterized by $\boldsymbol{w}$.

- Generally use mean squared error and define the loss as

$$J(\boldsymbol{w}) = \mathbb{E}_\pi[(V^\pi(s) - \hat{V}(s; \boldsymbol{w}))^2]$$

$$s \sim \pi$$
$$true$$

- Can use gradient descent to find a local minimum $\left(V^\pi - \hat{V}\right) \nabla_w V(s, w)$

$$\Delta \boldsymbol{w} = -\frac{1}{2}\alpha \nabla_{\boldsymbol{w}} J(\boldsymbol{w})$$

- Stochastic gradient descent (SGD) uses a finite number of (often one) samples to compute an approximate gradient:

$$\nabla_w J(w) = \nabla_w E_\pi \left[ V^\pi(s) - \hat{V}(s, w) \right]^2$$
$$= E_\pi \, 2 \left( V^\pi(s) - \hat{V}(s, w) \right] \nabla_w V(s, w)$$

- Expected SGD is the same as the full gradient update

# Stochastic Gradient Descent

- Goal: Find the parameter vector $\boldsymbol{w}$ that minimizes the loss between a true value function $V^\pi(s)$ and its approximation $\hat{V}(s; \boldsymbol{w})$ as represented with a particular function class parameterized by $\boldsymbol{w}$.

- Generally use mean squared error and define the loss as

$$J(\boldsymbol{w}) = \mathbb{E}_\pi[(V^\pi(s) - \hat{V}(s; \boldsymbol{w}))^2]$$

- Can use gradient descent to find a local minimum

$$\Delta \boldsymbol{w} = -\frac{1}{2}\alpha\nabla_{\boldsymbol{w}} J(\boldsymbol{w})$$

- Stochastic gradient descent (SGD) uses a finite number of (often one) samples to compute an approximate gradient:

$$
\begin{aligned}
\Delta_w J(w) &= \Delta_w E_\pi[V^\pi(s) - \hat{V}(s; w)]^2 \\
&= E_\pi[2(V^\pi(s) - \hat{V}(s; w))\Delta_w \hat{V}(s, w)]
\end{aligned}
$$

- Expected SGD is the same as the full gradient update

# Model Free VFA Policy Evaluation

- Don't actually have access to an oracle to tell true $V^\pi(s)$ for any state $s$
- Now consider how to do model-free value function approximation for prediction / evaluation / policy evaluation without a model

# Model Free VFA Prediction / Policy Evaluation

- Recall model-free policy evaluation (Lecture 3)
  - Following a fixed policy $\pi$ (or had access to prior data)
  - Goal is to estimate $V^\pi$ and/or $Q^\pi$
- Maintained a lookup table to store estimates $V^\pi$ and/or $Q^\pi$
- Updated these estimates after each episode (Monte Carlo methods) or after each step (TD methods)
- **Now: in value function approximation, change the estimate update step to include fitting the function approximator**
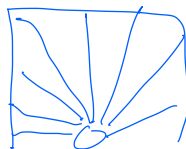
# Break

## Outline for Today

- Value function approximation
- **Monte Carlo policy evaluation with linear function approximation**
- TD policy evaluation with linear function approximation
- Control methods with linear value function approximation

# Feature Vectors

- Use a feature vector to represent a state $s$



$$\mathbf{x}(s) = \begin{pmatrix} x_1(s) \\ x_2(s) \\ \ldots \\ x_n(s) \end{pmatrix}$$

# Linear Value Function Approximation for Prediction With An Oracle

- Represent a value function (or state-action value function) for a particular policy with a weighted linear combination of features

$$\hat{V}(s; \boldsymbol{w}) = \sum_{j=1}^{n} x_j(s) w_j = \boldsymbol{x}(s)^T \boldsymbol{w}$$

$$\nabla_w x(s)^T w = x(s)$$

- Objective function is

$$J(\boldsymbol{w}) = \mathbb{E}_\pi[(V^\pi(s) - \hat{V}(s; \boldsymbol{w}))^2]$$

- Recall weight update is

$$\Delta \boldsymbol{w} = -\frac{1}{2} \alpha \nabla_{\boldsymbol{w}} J(\boldsymbol{w})$$

prediction error

feature value

- Update is: $\Delta w = -\frac{1}{2} \alpha \left( V^\pi(s) - x(s)^T w \right) x(s)$

# Feature Vectors

- Use a feature vector to represent a state $s$

$$\mathbf{x}(s) = \begin{pmatrix} x_1(s) \\ x_2(s) \\ \ldots \\ x_n(s) \end{pmatrix}$$

# Linear Value Function Approximation for Prediction With An Oracle

- Represent a value function (or state-action value function) for a particular policy with a weighted linear combination of features

$$\hat{V}(s; \boldsymbol{w}) = \sum_{j=1}^{n} x_j(s) w_j = \boldsymbol{x}(s)^T \boldsymbol{w}$$

- Objective function is

$$J(\boldsymbol{w}) = \mathbb{E}_\pi [(V^\pi(s) - \hat{V}(s; \boldsymbol{w}))^2]$$

- Recall weight update is

$$\Delta \boldsymbol{w} = -\frac{1}{2} \alpha \nabla_{\boldsymbol{w}} J(\boldsymbol{w})$$

- Update is:

- Update = step-size × prediction error × feature value

- Return $G_t$ is an unbiased but noisy sample of the true expected return $V^\pi(s_t)$
- Therefore can reduce MC VFA to doing supervised learning on a set of (state,return) pairs: $\langle s_1, G_1 \rangle, \langle s_2, G_2 \rangle, \dots, \langle s_T, G_T \rangle$
  - Substitute $G_t$ for the true $V^\pi(s_t)$ when fit function approximator

$$\triangle w = \alpha \left( V^\pi(s) - \hat{V}(s, w) \right) x(s_t)$$
$$= \alpha \left( G_t(s) - x(s_t)w \right) x(s_t)$$

# Monte Carlo Value Function Approximation

- Return $G_t$ is an unbiased but noisy sample of the true expected return $V^\pi(s_t)$
- Therefore can reduce MC VFA to doing supervised learning on a set of (state,return) pairs: $\langle s_1, G_1 \rangle, \langle s_2, G_2 \rangle, \ldots, \langle s_T, G_T \rangle$
  - Substitute $G_t$ for the true $V^\pi(s_t)$ when fit function approximator
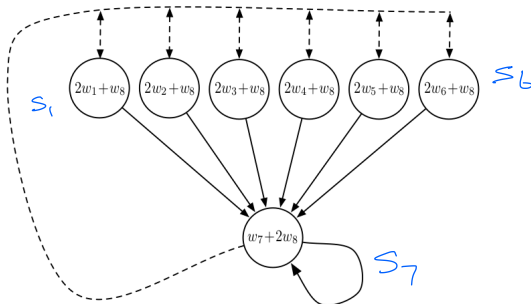- Concretely when using linear VFA for policy evaluation

$$\begin{aligned}
\Delta \boldsymbol{w} &= \alpha(G_t - \hat{V}(s_t; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{V}(s_t; \boldsymbol{w}) \\
&= \alpha(G_t - \hat{V}(s_t; \boldsymbol{w}))\boldsymbol{x}(s_t) \\
&= \alpha(G_t - \boldsymbol{x}(s_t)^T\boldsymbol{w})\boldsymbol{x}(s_t)
\end{aligned}$$

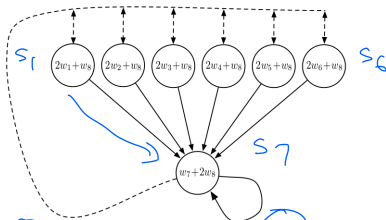- Note: $G_t$ may be a very noisy estimate of true return

1: Initialize w $= 0$, $k = 1$
2: **loop**
3:    Sample $k$-th episode $(s_{k,1}, a_{k,1}, r_{k,1}, s_{k,2}, \ldots, s_{k,L_k})$ given $\pi$
4:    **for** $t = 1, \ldots, L_k$ **do**
5:       **if** First visit to $(s)$ in episode $k$ **then**
6:          $G_t(s) = \sum_{j=t}^{L_k} r_{k,j}$   *noh   could  weight  by $\gamma$*
7:          Update weights:
                $w = \alpha \left( G_t(s) - X(s)^T w \right) X(s)$

8:       **end if**
9:    **end for**
10:    $k = k + 1$
11: **end loop**

- $x(s_1) = [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]\ x(s_2) = [0\ 2\ 0\ 0\ 0\ 0\ 0\ 1]\ \ldots x(s_6) = [0\ 0\ 0\ 0\ 0\ 2\ 0\ 1]$
  $x(s_7) = [0\ 0\ 0\ 0\ 0\ 0\ 1\ 2]$     $r(s) = 0\ \forall s$     2 actions $a_1$ solid line, $a_2$ dotted
- Small prob $s_7$ goes to terminal state $s_T$

- $x(s_1) = [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$ $x(s_2) = [0\ 2\ 0\ 0\ 0\ 0\ 0\ 1]$ ... $x(s_6) = [0\ 0\ 0\ 0\ 0\ 2\ 0\ 1]$
  $x(s_7) = [0\ 0\ 0\ 0\ 0\ 0\ 1\ 2]$    $r(s) = 0\ \forall s$       2 actions $a_1$ solid line, $a_2$ dotted
- Small prob $s_7$ goes to terminal state $s_T$
- Consider trajectory $(s_1, a_1, 0, s_7, a_1, 0, s_7, a_1, 0, s_T)$. $G(s_1) = 0$
- Let $w_0 = [1\ 1\ 1\ 1\ 1\ 1\ 1\ 1]$. MC update: $\Delta w = \alpha(G_t - x(s_t)^T w)x(s_t)$

$$x(s_t)^T w = 3$$

$$w = w + \alpha\ (0 - 3)\ [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$$

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value as

$$MSVE_{\mu}(\boldsymbol{w}) = \sum_{s \in S} \mu(s)(V^{\pi}(s) - \hat{V}^{\pi}(s; \boldsymbol{w}))^2$$

- where
  - $\mu(s)$: probability of visiting state $s$ under policy $\pi$. Note $\sum_s \mu(s) = 1$
  - $\hat{V}^{\pi}(s; \boldsymbol{w}) = \boldsymbol{x}(s)^T \boldsymbol{w}$, a linear value function approximation

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value as

$$MSVE_\mu(\mathbf{w}) = \sum_{s \in S} \mu(s)(V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$

- where
  - $\mu(s)$: probability of visiting state $s$ under policy $\pi$. Note $\sum_s \mu(s) = 1$
  - $\hat{V}^\pi(s; \mathbf{w}) = \mathbf{x}(s)^T \mathbf{w}$, a linear value function approximation
- Monte Carlo policy evaluation with VFA converges to the weights $\mathbf{w}_{MC}$ which has the minimum mean squared error possible with respect to the distribution $\mu$:

$$MSVE_\mu(\mathbf{w}_{MC}) = \min_{\mathbf{w}} \sum_{s \in S} \mu(s)(V^\pi(s) - \hat{V}^\pi(s; \mathbf{w}))^2$$
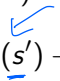
# Break

# Today: Focus on Generalization using Linear Value Function

- Preliminaries
- Monte Carlo policy evaluation with linear function approximation
- **TD policy evaluation with linear function approximation**
- Control methods with linear value function approximation

- Uses bootstrapping and sampling to approximate $V^\pi$
- Updates $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) \;=\; V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$, a biased estimate of the true value $V^\pi(s)$
- Represent value for each state with a separate table entry

# Temporal Difference (TD(0)) Learning with Value Function Approximation

- Uses bootstrapping and sampling to approximate true $V^\pi$
- Updates estimate $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$, a biased estimate of the true value $V^\pi(s)$
- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \boldsymbol{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- 3 forms of approximation:
  - sampling (vs expectation over s')
  - bootstrapping
  - VFA

# Temporal Difference (TD(0)) Learning with Value Function Approximation

- Uses bootstrapping and sampling to approximate true $V^\pi$
- Updates estimate $V^\pi(s)$ after each transition $(s, a, r, s')$:

$$V^\pi(s) = V^\pi(s) + \alpha(r + \gamma V^\pi(s') - V^\pi(s))$$

- Target is $r + \gamma V^\pi(s')$, a biased estimate of the true value $V^\pi(s)$
- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \boldsymbol{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- 3 forms of approximation:
    1. Sampling
    2. Bootstrapping
    3. Value function approximation

# Temporal Difference (TD(0)) Learning with Value Function Approximation

- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \boldsymbol{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- Can reduce doing TD(0) learning with value function approximation to supervised learning on a set of data pairs:
  - $\langle s_1, \underbrace{r_1 + \gamma \hat{V}^\pi(s_2; \boldsymbol{w})} \rangle, \langle s_2, r_2 + \gamma \underbrace{\hat{V}(s_3; \boldsymbol{w})} \rangle, \dots$   *TD targets*
- Find weights to minimize mean squared error

$$J(\boldsymbol{w}) = \mathbb{E}_\pi[(r_j + \gamma \hat{V}^\pi(s_{j+1}, \boldsymbol{w}) - \hat{V}(s_j; \boldsymbol{w}))^2]$$

# Temporal Difference (TD(0)) Learning with Value Function Approximation

- In value function approximation, target is $r + \gamma \hat{V}^\pi(s'; \mathbf{w})$, a biased and approximated estimate of the true value $V^\pi(s)$
- Supervised learning on a different set of data pairs:
  $\langle s_1, r_1 + \gamma \hat{V}^\pi(s_2; \mathbf{w})\rangle, \langle s_2, r_2 + \gamma \hat{V}(s_3; \mathbf{w})\rangle, \ldots$
- In linear TD(0)

$$
\begin{aligned}
\Delta \mathbf{w} &= \alpha(r + \gamma \hat{V}^\pi(s'; \mathbf{w}) - \hat{V}^\pi(s; \mathbf{w}))\nabla_{\mathbf{w}} \hat{V}^\pi(s; \mathbf{w}) \\
&= \alpha(r + \gamma \hat{V}^\pi(s'; \mathbf{w}) - \hat{V}^\pi(s; \mathbf{w}))\mathbf{x}(s) \\
&= \alpha(r + \gamma \mathbf{x}(s')^T \mathbf{w} - \mathbf{x}(s)^T \mathbf{w})\mathbf{x}(s)
\end{aligned}
$$

# TD(0) Linear Value Function Approximation for Policy Evaluation

1: Initialize $w = 0$, $k = 1$
2: **loop**
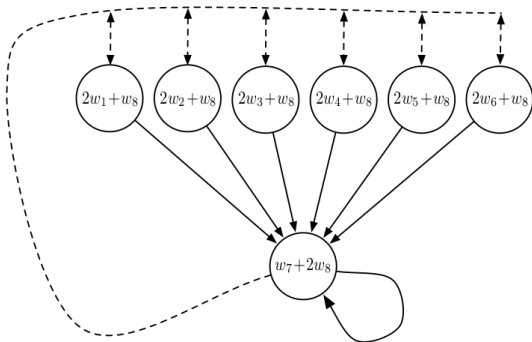3:     Sample tuple $(s_k, a_k, r_k, s_{k+1})$ given $\pi$
4:     Update weights:

$$w = w + \alpha(r + \gamma \underbrace{x(s')^T w}_{V(s')} - x(s)^T w)x(s)$$

5:     $k = k + 1$
6: **end loop**

- TD update: $\Delta \boldsymbol{w} = \alpha(r + \gamma \boldsymbol{x}(s')^T \boldsymbol{w} - \boldsymbol{x}(s)^T \boldsymbol{w}) \boldsymbol{x}(s)$

---

[1]Figure from Sutton and Barto 2018

- $x(s_1) = [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$ $x(s_2) = [0\ 2\ 0\ 0\ 0\ 0\ 0\ 1]$ ... $x(s_6) = [0\ 0\ 0\ 0\ 0\ 2\ 0\ 1]$
  $x(s_7) = [0\ 0\ 0\ 0\ 0\ 0\ 1\ 2]$   $r(s) = 0\ \forall s$    2 actions $a_1$ solid line, $a_2$ dotted
- Small prob $s_7$ goes to terminal state $s_T$
- Consider tuple $(s_1, a_1, 0, s_7)$.
- Let $w_0 = [1\ 1\ 1\ 1\ 1\ 1\ 1\ 1]$. TD update: $\Delta w = \alpha(r + \gamma x(s')^T w - x(s)^T w)x(s)$
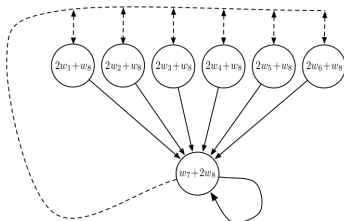
$V^\pi(s_1) = x(s_1)w = 3$        $V^\pi(s_7) = 3$

$\Delta \omega = \alpha(\gamma \cdot 3 - 3)\ [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$

vs MC        $(0 - 3)$

# Baird Example with TD(0) On Policy Evaluation [1]



- $x(s_1) = [2\ 0\ 0\ 0\ 0\ 0\ 0\ 1]$ $x(s_2) = [0\ 2\ 0\ 0\ 0\ 0\ 0\ 1] \ldots x(s_6) = [0\ 0\ 0\ 0\ 0\ 2\ 0\ 1]$
  $x(s_7) = [0\ 0\ 0\ 0\ 0\ 0\ 1\ 2]$     $r(s) = 0\ \forall s$     2 actions $a_1$ solid line, $a_2$ dotted

- Small prob $s_7$ goes to terminal state $s_T$

- Consider tuple $(s_1, a_1, 0, s_7)$.

- Let $w_0 = [1\ 1\ 1\ 1\ 1\ 1\ 1\ 1]$. TD update: $\Delta w = \alpha(r + \gamma x(s')^T w - x(s)^T w)x(s)$

- TD target is $r + \gamma x(s')^T w$. $r = 0$     $x(s')^T w = 3$.

- $x(s)^T w = 3$

- $\Delta w = \alpha(3\gamma - 3)x(s_1)$

[1]Figure from Sutton and Barto 2018

# Convergence Guarantees for TD Linear VFA for Policy Evaluation: Preliminaries

- For infinite horizon, the Markov Chain defined by a MDP with a particular policy will eventually converge to a probability distribution over states $d(s)$
- $d(s)$ is called the stationary distribution over states of $\pi$
- $\sum_s d(s) = 1$
- $d(s)$ satisfies the following balance equation:

*Markov systems*

$$d(s') = \sum_s \sum_a \pi(a|s)p(s'|s, a)d(s)$$

# Convergence Guarantees for Linear Value Function Approximation for Policy Evaluation

- Define the mean squared error of a linear value function approximation for a particular policy $\pi$ relative to the true value given the distribution $d$ as

$$MSVE_d(\boldsymbol{w}) = \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- where
  - $d(s)$: stationary distribution of $\pi$ in the true decision process
  - $\hat{V}^\pi(s; \boldsymbol{w}) = \boldsymbol{x}(s)^T \boldsymbol{w}$, a linear value function approximation

- TD(0) policy evaluation with VFA converges to weights $\boldsymbol{w}_{TD}$ which is within a constant factor of the min mean squared error possible given distribution $d$:

$$\gamma = .9$$

$$MSVE_d(\boldsymbol{w}_{TD}) \leq \frac{1}{1 - \gamma} \min_{\boldsymbol{w}} \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

# Check Your Understanding L5N1: Poll

- TD(0) policy evaluation with VFA converges to weights $\boldsymbol{w}_{TD}$ which is within a constant factor of the min mean squared error possible for distribution $d$:

$$MSVE_d(\boldsymbol{w}_{TD}) \leq \frac{1}{1-\gamma} \min_{\boldsymbol{w}} \sum_{s \in S} d(s)(V^\pi(s) - \hat{V}^\pi(s; \boldsymbol{w}))^2$$

- If the VFA is a tabular representation (one feature for each state), what is the $MSVE_d$ for TD?

1. Depends on the problem
2. MSVE $= 0$ for TD
3. Not sure

- TD(0) policy evaluation with VFA converges to weights $\boldsymbol{w}_{TD}$ which is within a constant factor of the min mean squared error possible for distribution $d$:

$$MSVE_d(\boldsymbol{w}_{TD}) \leq \frac{1}{1-\gamma} \min_{\boldsymbol{w}} \sum_{s \in S} d(s)(V^{\pi}(s) - \hat{V}^{\pi}(s; \boldsymbol{w}))^2$$

- If the VFA is a tabular representation (one feature for each state), what is the $MSVE_d$ for TD?

$MSVE = 0$ for TD

# Break

- Value function approximation
- Monte Carlo policy evaluation with linear function approximation
- TD policy evaluation with linear function approximation
- Control methods with linear value function approximation

# Table of Contents

# Control using Value Function Approximation

- Use value function approximation to represent state-action values $\hat{Q}^\pi(s, a; \boldsymbol{w}) \approx Q^\pi$
- Interleave
  - Approximate policy evaluation using value function approximation
  - Perform $\epsilon$-greedy policy improvement
- Can be unstable. Generally involves intersection of the following:
  - Function approximation
  - Bootstrapping
  - **Off-policy learning**

# Action-Value Function Approximation with an Oracle

- $\hat{Q}^\pi(s, a; \boldsymbol{w}) \approx Q^\pi$
- Minimize the mean-squared error between the true action-value function $Q^\pi(s, a)$ and the approximate action-value function:

$$J(\boldsymbol{w}) = \mathbb{E}_\pi[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \boldsymbol{w}))^2]$$

- Use stochastic gradient descent to find a local minimum

$$
\begin{aligned}
-\frac{1}{2}\nabla_{\boldsymbol{W}} J(\boldsymbol{w}) &= \mathbb{E}\left[(Q^\pi(s, a) - \hat{Q}^\pi(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}^\pi(s, a; \boldsymbol{w})\right] \\
\Delta(\boldsymbol{w}) &= -\frac{1}{2}\alpha\nabla_{\boldsymbol{w}} J(\boldsymbol{w})
\end{aligned}
$$

- Stochastic gradient descent (SGD) samples the gradient

- The weight update for control for MC and TD-style methods will be near identical to the policy evaluation steps. Try to see if you can predict which are the right weight update equations for the different methods (select all that are true)
- (1) is the SARSA control update
- (2) is the MC control update
- (3) is the Q-learning control update
- (4) is the MC control update
- (5) is the Q-learning control update

$(s\ a\ r\ s'\ a')$

$$\Delta \boldsymbol{w} = \alpha(r + \gamma \hat{Q}(\underline{s', a'}; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w}) \quad (1) \quad \text{true}$$

$$\Delta \boldsymbol{w} = \alpha(G_t + \gamma \hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w}) \quad (2)$$

$$\Delta \boldsymbol{w} = \alpha(r + \gamma \max_{\underline{a'}} \hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w}) \quad (3) \quad \text{true}$$

$$\Delta \boldsymbol{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s_t, a_t; \boldsymbol{w}) \quad (4) \quad \text{true}$$

$$\Delta \boldsymbol{w} = \alpha(r + \gamma \max_{s'} \hat{Q}(s', a; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w}) \quad (5) \quad \times$$

# Check Your Understanding L5N2: Answers

- The weight update for control for MC and TD-style methods will be near identical to the policy evaluation steps. Try to see if you can predict which are the right weight update equations for the different methods.

- (1) is the SARSA control update
  $\Delta \mathbf{w} = \alpha(r + \gamma \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})$

- (3) is the Q-learning control update
  $\Delta \mathbf{w} = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; \mathbf{w}) - \hat{Q}(s, a; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s, a; \mathbf{w})(3)$

- (4) is the MC control update
  $\Delta \mathbf{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \mathbf{w})) \nabla_{\mathbf{w}} \hat{Q}(s_t, a_t; \mathbf{w})$

# Linear State Action Value Function Approximation with an Oracle

- Use features to represent both the state and action

$$\boldsymbol{x}(s,a) = \begin{pmatrix} x_1(s,a) \\ x_2(s,a) \\ \dots \\ x_n(s,a) \end{pmatrix}$$

- Represent state-action value function with a weighted linear combination of features

$$\hat{Q}(s,a;\boldsymbol{w}) = \boldsymbol{x}(s,a)^T \boldsymbol{w} = \sum_{j=1}^{n} x_j(s,a) w_j$$

- Stochastic gradient descent update:

$$\nabla_{\boldsymbol{w}} J(\boldsymbol{w}) = \nabla_{\boldsymbol{w}} \mathbb{E}_{\pi}[(Q^\pi(s,a) - \hat{Q}^\pi(s,a;\boldsymbol{w}))^2]$$

## Incremental Model-Free Control Approaches

- Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value

- In Monte Carlo methods, use a return $G_t$ as a substitute target

$$\Delta \boldsymbol{w} = \alpha(G_t - \hat{Q}(s_t, a_t; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s_t, a_t; \boldsymbol{w})$$

- For SARSA instead use a TD target $r + \gamma\hat{Q}(s', a'; \boldsymbol{w})$ which leverages the current function approximation value

$$\Delta \boldsymbol{w} = \alpha(r + \gamma\hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w})$$

- For Q-learning instead use a TD target $r + \gamma\max_{a'}\hat{Q}(s', a'; \boldsymbol{w})$ which leverages the max of the current function approximation value

$$\Delta \boldsymbol{w} = \alpha(r + \gamma\max_{a'}\hat{Q}(s', a'; \boldsymbol{w}) - \hat{Q}(s, a; \boldsymbol{w}))\nabla_{\boldsymbol{w}}\hat{Q}(s, a; \boldsymbol{w})$$

$U_t$ is 1 value fn

- Informally, updates involve doing an (approximate) Bellman backup followed by best trying to fit underlying value function to a particular feature representation

- Bellman operators are contractions, but value function approximation fitting can be an expansion
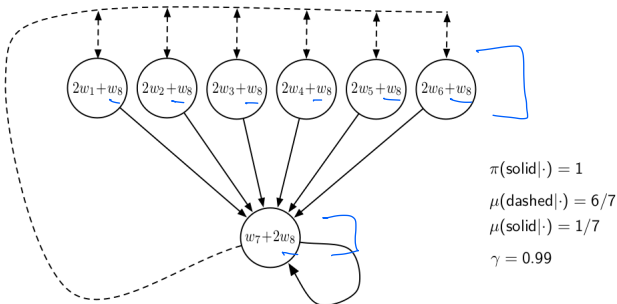
$$|BV_1 - BV_2|_\infty \le \gamma |V_1 - V_2|_\infty \quad \text{tabular}$$

$$|OV_1 - OV_2|_\infty \ge |V_1 - V_2|_\infty$$

↑
VFA

Geoff Gordon 1995

$\pi(\text{solid}|\cdot) = 1$

$\mu(\text{dashed}|\cdot) = 6/7$

$\mu(\text{solid}|\cdot) = 1/7$

$\gamma = 0.99$

- Behavior policy and target policy are not identical
- Value can diverge     *w can go to ∞*

# Convergence of Control Methods with VFA

| Algorithm | Tabular | Linear VFA |
|:---:|:---:|:---:|
| Monte-Carlo Control | ✓ | ✓ |
| Sarsa | ✓ | ✓ |
| Q-learning | ✓ | ✗ |

# Important Open Area: Off Policy Learning with Function Approximation

- Extensive work in better TD-style algorithms with value function approximation, some with convergence guarantees: see Chp 11 SB
- Will come up further later in this course

The space of all value functions

$B_\pi v_\mathbf{w}$

Bellman error (BE)

Value error (VE)

$v_\pi$

PBE

$\Pi B_\pi v_\mathbf{w}$

$v_\mathbf{w}$

$\Pi v_\pi \equiv \min \|\text{VE}\|$

$\text{PBE} = \vec{0}$

$\min \|\text{BE}\|$

$w_1$

$w_2$

The subspace of all value functions representable as $v_\mathbf{w}$

---

# What You Should Understand

- Be able to implement TD(0) and MC on policy evaluation with linear value function approximation
- Be able to define what TD(0) and MC on policy evaluation with linear VFA are converging to and when this solution has 0 error and non-zero error.
- Be able to implement Q-learning and SARSA and MC control algorithms
- List the 3 issues that can cause instability and describe the problems qualitatively: function approximation, bootstrapping and off policy learning

# Class Structure

- Last time: Control (making decisions) without a model of how the world works
- This time: Value function approximation
- **Next time**: Deep reinforcement learning

# Batch Monte Carlo Value Function Approximation

- May have a set of episodes from a policy $\pi$
- Can analytically solve for the best linear approximation that minimizes mean squared error on this data set
- Let $G(s_i)$ be an unbiased sample of the true expected return $V^\pi(s_i)$

$$\arg \min_{\textbf{w}} \sum_{i=1}^{N} (G(s_i) - \textbf{x}(s_i)^T \textbf{w})^2$$

- Take the derivative and set to 0

$$\textbf{w} = (X^T X)^{-1} X^T \textbf{G}$$

- where $\textbf{G}$ is a vector of all $N$ returns, and $X$ is a matrix of the features of each of the $N$ states $\textbf{x}(s_i)$
- Note: not making any Markov assumptions

# Value Function Approximation (VFA)

- Represent a (state-action/state) value function with a parameterized function instead of a table

$$s \longrightarrow \boxed{w} \longrightarrow \hat{V}(s; w)$$

$$\begin{matrix} s \\ a \end{matrix} \longrightarrow \boxed{w} \longrightarrow \hat{Q}(s, a; w)$$