# Lecture 20: Off-Policy Learning

#### **Off-policy Methods**

- Learn the value of the *target policy*  $\pi$  from experience due to *behavior policy*  $\mu$
- □ For example,  $\pi$  is the greedy policy (and ultimately the optimal policy) while  $\mu$  is exploratory (e.g.,  $\varepsilon$ -soft)
- In general, we only require *coverage*, i.e., that  $\mu$  generates behavior that covers, or includes,  $\pi$

 $\pi(a|s) > 0$  implies  $\mu(a|s) > 0$ 

- **Idea:** *importance sampling* 
  - Weight each return by the *ratio of the probabilities* of the trajectory under the two policies

#### **Importance Sampling in General**

- Suppose we want to estimate the expected value of a function f depending on a random variable X drawn according to the *target* probability distribution P(X).
- If we had N samples  $x_i$  drawn from P(X), we could estimate the expectation using the empirical mean:

$$E_P[f] \approx \frac{1}{N} \sum_{i=1}^N f(x_i)$$

- But instead, we have only samples drawn according to a different *proposal* or *sampling* distribution Q(X).
- How can we do the estimation?

#### **Regular Importance Sampling**

• We do a simple trick:

$$E_P[f] = \sum_x f(x)P(X = x)$$
  
= 
$$\sum_x f(x)Q(X = x)\frac{P(X = x)}{Q(X = x)} = E_Q\left[f\frac{P}{Q}\right]$$

- Only requirement: if P(x) > 0 then Q(x) > 0
- So for an estimator, we should average each sample of the function,  $f(x_i)$  weighted by the ratio of its probability under the target and the sampling distribution:

$$E_p[f] \approx \frac{1}{N} \sum_{i=1}^N f(x_i) \frac{P(x_i)}{Q(x_i)}$$

### **Applying IS to Policy Evaluation**

- **I** Function for which we want the expectation is the return
- Target distribution P is the distribution of trajectories under *target policy*  $\pi$
- The Proposal distribution Q is distribution of trajectories under *behavior policy*  $\mu$
- Note that P and Q can be very different depending on the horizon!
- **But there is structure in P and Q that we can exploit**

#### **Importance Sampling Ratio**



In importance sampling, each return is weighted by the relative probability of the trajectory under the two policies

$$\rho_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) P(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} \mu(A_k | S_k) P(S_{k+1} | S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{\mu(A_k | S_k)}$$

- This is called the *importance sampling ratio*
- □ All importance sampling ratios have expected value 1

$$\mathsf{E}_{\mu}\left[\frac{\pi(A_{k} \mid S_{k})}{\mu(A_{k} \mid S_{k})}\right] = \sum_{a} \mu(a \mid S_{k}) \frac{\pi(a \mid S_{k})}{\mu(a \mid S_{k})} \sum_{a} \pi(a \mid S_{k}) = 1$$

#### **Per-reward Importance Sampling**

 $\square$  Another way of reducing variance, even if  $\gamma = 1$ 

Uses the fact that the return is a *sum of rewards* 

 $\rho_t^T G_t = \rho_t^T R_{t+1} + \gamma \rho_t^T R_{t+2} + \dots + \gamma^{k-1} \rho_t^T R_{t+k} + \dots + \gamma^{T-t-1} \rho_t^T R_T$ where  $\rho_t^T R_{t+k} = \frac{\pi(A_t | S_t)}{\mu(A_t | S_t)} \frac{\pi(A_{t+1} | S_{t+1})}{\mu(A_{t+1} | S_{t+1})} \dots \frac{\pi(A_{t+k} | S_{t+k})}{\mu(A_{t+k} | S_{t+k})} \dots \frac{\pi(A_{T-1} | S_{T-1})}{\mu(A_{T-1} | S_{T-1})} R_{t+k}$ 

#### **Per-reward Importance Sampling**

- $\square$  Another way of reducing variance, even if  $\gamma = 1$
- Uses the fact that the return is a *sum of rewards*

$$\rho_{t:T-1}G_t = \rho_{t:T-1}R_{t+1} + \dots + \gamma^{k-1}\rho_{t:T-1}R_{t+k} + \dots + \gamma^{T-t-1}\rho_{t:T-1}R_T$$

$$\rho_{t:T-1}R_{t+k} = \frac{\pi(A_t|S_t)}{b(A_t|S_t)}\frac{\pi(A_{t+1}|S_{t+1})}{b(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_{t+k}|S_{t+k})}{b(A_{t+k}|S_{t+k})} \cdots \frac{\pi(A_{T-1}|S_{T-1})}{b(A_{T-1}|S_{T-1})}R_{t+k}.$$

 $\therefore \mathbb{E}[\rho_{t:T-1}R_{t+k}] = \mathbb{E}[\rho_{t:t+k-1}R_{t+k}]$ 

$$\therefore \mathbb{E}[\rho_{t:T-1}G_t] = \mathbb{E}\left[\underbrace{\rho_{t:t}R_{t+1} + \dots + \gamma^{k-1}\rho_{t:t+k-1}R_{t+k} + \dots + \gamma^{T-t-1}\rho_{t:T-1}R_T}_{\tilde{G}_t}\right]$$

$$V(s) \doteq \frac{\sum_{t \in \mathfrak{T}(s)} \tilde{G}_t}{|\mathfrak{T}(s)|}$$

#### Implementation

- Importance sampling ratios fold into the eligibility trace
- Multiply at each step by an extra factor
- But on long trajectories traces will get cut a lot!





#### **Recall: Q-Learning is Off-Policy TD Control**

One-step Q-learning:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

Behavior is randomized, but we are evaluating the greedy policy

#### **Off-policy Expected Sarsa**

- Expected Sarsa generalizes to arbitrary behavior policies  $\mu$ 
  - in which case it includes Q-learning as the special case in which  $\pi$  is the greedy policy

$$Q(S_{t}, A_{t}) \leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

$$\leftarrow Q(S_{t}, A_{t}) + \alpha \Big[ R_{t+1} + \gamma \sum_{a} \pi(a \mid S_{t+1}) Q(S_{t+1}, a) - Q(S_{t}, A_{t}) \Big]$$

#### **Q-learning with Eligibility Traces**



#### **Tree Backup**



Reweight the traces by the product of target probabilities

#### **Blueprint Off-policy Q-Algorithms**

$$\Delta Q(x,a) = \sum_{t \ge 0} \gamma^t \Big(\prod_{1 \le s \le t} c_s\Big) \Big(\underbrace{r_t + \gamma \mathbb{E}_{\pi} Q(x_{t+1}, \cdot) - Q(x_t, a_t)}_{\delta_t}\Big)$$

Algorithm:	Trace coefficient:	Problem:
IS	$c_s = \frac{\pi(a_s   x_s)}{\mu(a_s   x_s)}$	high variance
$Q^{\pi}(\lambda)$	$c_s = \lambda$	not safe (off-policy)
$TB(\lambda)$	$c_s = \lambda \pi(a_s   x_s)$	not efficient (on-policy)

Use Retrace(
$$\lambda$$
) defined by  $c_s = \lambda \min\left(1, \frac{\pi(a_s|x_s)}{\mu(a_s|x_s)}\right)$ 

#### **Properties:**

- Low variance since  $c_s \leq 1$
- Safe (off policy): cut the traces when needed  $c_s \in \left[0, \frac{\pi(a_s|x_s)}{\mu(a_a|x_s)}\right]$
- Efficient (on policy): but only when needed. Note that  $c_s \ge \lambda \pi (a_s | x_s)$

#### **Retrace in Atari**



Games:

Asteroids, Defender, Demon Attack, Hero, Krull, River Raid, Space Invaders, Star Gunner, Wizard of Wor, Zaxxon

#### **Retrace vs Tree Backup**



#### **Off-policy is much harder with Function Approximation**

- Even linear FA
- **\Box** Even for prediction (two fixed policies  $\pi$  and  $\mu$ )
- **D** Even for Dynamic Programming
- □ The deadly triad: FA, TD, off-policy
  - Any two are OK, but not all three
  - With all three, we may get instability (elements of θ may increase to ±∞)

### **Two Off-Policy Learning Problems**

The easy problem is that of off-policy targets (future)

- Use importance sampling in the target
- The hard problem is that of the distribution of states to update (present): we are no longer updating according to the on-policy distribution

#### **Baird's counterexample**



 $\theta_7 + 2\theta_8$ 

#### TD(0) can diverge: A simple example



**TD fixpoint:**  $\theta^* = 0$ 

#### What causes the instability?

□ It has nothing to do with learning or sampling

- Even dynamic programming suffers from divergence with FA
- It has nothing to do with exploration, greedification, or control
  - Even prediction alone can diverge
- It has nothing to do with local minima or complex non-linear approximators
  - Even simple linear approximators can produce instability

### The deadly triad

- The risk of divergence arises whenever we combine three things:
  - Function approximation
    - significantly generalizing from large numbers of examples
  - Bootstrapping
    - learning value estimates from other value estimates, as in dynamic programming and temporal-difference learning
  - Off-policy learning
    - learning about a policy from data not due to that policy, as in Q-learning, where we learn about the greedy policy from data with a necessarily more exploratory policy

#### How to survive the deadly triad

- Least-squares methods like off-policy LSTD(λ) (Yu 2010, Mahmood et al. 2015, Bradtke & Barto 1996, Boyan 2000) computational costs scale with the *square* of the number of parameters
- True-gradient RL methods (Gradient-TD and proximalgradient-TD) (Maei et al, 2011, Mahadevan et al, 2015)
- Emphatic-TD methods (Sutton, White & Mahmood 2015, Yu 2015). These semi-gradient methods attain stability through an extension of the early on-policy theorems

#### **Linear Least-Squares**

• At minimum of  $LS(\mathbf{w})$ , the expected update must be zero

$$\mathbb{E}_{\mathcal{D}} \left[ \Delta \mathbf{w} \right] = 0$$

$$\alpha \sum_{t=1}^{T} \mathbf{x}(s_t) (v_t^{\pi} - \mathbf{x}(s_t)^{\top} \mathbf{w}) = 0$$

$$\sum_{t=1}^{T} \mathbf{x}(s_t) v_t^{\pi} = \sum_{t=1}^{T} \mathbf{x}(s_t) \mathbf{x}(s_t)^{\top} \mathbf{w}$$

$$\mathbf{w} = \left( \sum_{t=1}^{T} \mathbf{x}(s_t) \mathbf{x}(s_t)^{\top} \right)^{-1} \sum_{t=1}^{T} \mathbf{x}(s_t) v_t^{\pi}$$

For N features, direct solution time is  $O(N^3)$ 

Incremental solution time is  $O(N^2)$  using Shermann-Morrison

## LSTD

- We do not know true values  $v_t^{\pi}$
- In practice, our "training data" must use noisy or biased samples of  $v_t^{\pi}$ 
  - LSMC Least Squares Monte-Carlo uses return  $v_t^{\pi} \approx G_t$
  - LSTD Least Squares Temporal-Difference uses TD target  $v_t^{\pi} \approx R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$
- $\begin{array}{l} \mathsf{LSTD}(\lambda) \ \ \mathsf{Least} \ \mathsf{Squares} \ \mathsf{TD}(\lambda) \ \mathsf{uses} \ \lambda\text{-return} \\ v_t^\pi \approx \mathbf{G}_t^\lambda \end{array}$

In each case solve directly for fixed point of MC / TD / TD( $\lambda$ )

#### **Convergence Properties**

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
	MC	$\checkmark$	$\checkmark$	$\checkmark$
On-Policy	LSMC	$\checkmark$	$\checkmark$	-
	TD	$\checkmark$	$\checkmark$	×
	LSTD	$\checkmark$	$\checkmark$	-
	MC	$\checkmark$	$\checkmark$	$\checkmark$
Off-Policy	LSMC	✓	$\checkmark$	-
	TD	$\checkmark$	×	×
	LSTD	$\checkmark$	<b>√</b>	-

Algorithm	Table Lookup	Linear	Non-Linear	
Monte-Carlo Control	✓	()	×	
Sarsa	$\checkmark$	$(\checkmark)$	×	
Q-learning	$\checkmark$	X	×	
LSPI	$\checkmark$	$(\checkmark)$	-	

 $(\checkmark)$  = chatters around near-optimal value function

#### Value function geometry



Mean Square Projected Bellman Error (MSPBE)

#### **Gradient-Based TD**

- Bootstraps (genuine TD)
- Works with linear function approximation (stable, reliably convergent)
- □ Is simple, like linear TD O(n)
- Learns fast, like linear TD
- Can learn off-policy
- Learns from online causal trajectories (no repeat sampling from the same state)

#### TD is not the gradient of anything

#### TD(0) algorithm:

Assume there is a J such that:

$$\Delta \theta = \alpha \delta \phi$$
  
$$\delta = r + \gamma \theta^{\top} \phi' - \theta^{\top} \phi$$

$$\frac{\partial J}{\partial \theta_i} = \delta \phi_i$$

## Then look at the second derivative:

$$\frac{\partial^2 J}{\partial \theta_j \partial \theta_i} = \frac{\partial (\delta \phi_i)}{\partial \theta_j} = (\gamma \phi'_j - \phi_j) \phi_i$$

$$\frac{\partial^2 J}{\partial \theta_i \partial \theta_j} = \frac{\partial (\delta \phi_j)}{\partial \theta_i} = (\gamma \phi'_i - \phi_i) \phi_j$$

$$\frac{\partial^2 J}{\partial \theta_j \partial \theta_i} \neq \frac{\partial^2 J}{\partial \theta_i \partial \theta_j}$$
Contradiction:

#### Real 2<sup>nd</sup> derivatives must be symmetric

Etienne Barnard 199

#### The Gradient-TD Family of Algorithms

- True gradient-descent algorithms in the Projected Bellman Error
- **T** GTD( $\lambda$ ) and GQ( $\lambda$ ), for learning V and Q
- Solve two open problems:
  - convergent linear-complexity off-policy TD learning
  - convergent non-linear TD
- Extended to control variate, proximal forms by Mahadevan et al.

#### First relate the geometry to the iid statistics



#### **Derivation of the TDC algorithm**

#### TD with gradient correction (TDC) algorithm



#### **Convergence theorems**

□ All algorithms converge w.p.1 to the TD fix-point:

 $\square \text{ GTD, GTD-2 converges at one time scale}$ 

$$\alpha=\beta\longrightarrow 0$$

**TD-C** converges in a two-time-scale sense  $\alpha, \beta \longrightarrow 0 \qquad \frac{\alpha}{\beta} \longrightarrow 0$ 

#### **Off-policy result: Baird's counter-example**



Gradient algorithms converge. TD diverges.

#### A little more theory

$$\Delta\theta \propto \delta\phi = (r + \gamma\theta^{\mathsf{T}}\phi' - \theta^{\mathsf{T}}\phi)\phi$$

$$= \theta^{\mathsf{T}}(\gamma\phi' - \phi)\phi + r\phi$$

$$= \phi(\gamma\phi' - \phi)^{\mathsf{T}}\theta + r\phi$$

$$\mathbb{E}[\Delta\theta] \propto -\mathbb{E}\left[\phi(\phi - \gamma\phi')^{\mathsf{T}}\right]\theta + \mathbb{E}[r\phi]$$

$$\mathbb{E}[\Delta\theta] \propto -A\theta + b$$
convergent if  
 $A \text{ is pos. def.}$ 
therefore, at  
the TD
 $\theta^* = A^{-1}b$ 
LSTD computes this directly  
fixpoint:
 $-\frac{1}{2}\nabla_{\theta}\text{MSPBE} = -A^{\mathsf{T}}C^{-1}(A\theta - b)$ 
always pos. def.
 $C = \mathbb{E}\left[\phi\phi^{\mathsf{T}}\right]$ 
covariance  
matrix

## Example: Go

- Learn a linear value function (probability of winning) for 9x9 Go from self play
- One million features, each corresponding to a template on a part of the Go board



#### Summary

		ALGORITHM						
		$TD(\lambda),$ Sarsa $(\lambda)$	Approx. DP	$\begin{array}{l} LSTD(\lambda),\\ LSPE(\lambda) \end{array}$	Fitted-Q	Residual gradient	GDP	$\begin{array}{c} {\rm GTD}(\lambda),\\ {\rm GQ}(\lambda) \end{array}$
ISSUE	Linear computation	$\checkmark$	$\checkmark$	*	*	$\checkmark$	$\checkmark$	$\checkmark$
	Nonlinear convergent	*	*	*	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$
	Off-policy convergent	*	*	$\checkmark$	*	$\checkmark$	$\checkmark$	$\checkmark$
	Model-free, online	$\checkmark$	*	$\checkmark$	×	$\checkmark$	₩	$\checkmark$
	Converges to PBE = 0	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	*	$\checkmark$	$\checkmark$