

Sequential decision making

Control:

Q-Learning & Deep Q-Learning (DQN)  
( and Eligibility Trace)

# On-policy MC Control

---

Initialize, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$ :

$Q(s, a) \leftarrow$  arbitrary

$Returns(s, a) \leftarrow$  empty list

$\pi(a|s) \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

Repeat forever:

(a) Generate an episode using  $\pi$

(b) For each pair  $s, a$  appearing in the episode:

$G \leftarrow$  return following the first occurrence of  $s, a$

Append  $G$  to  $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

(c) For each  $s$  in the episode:

$A^* \leftarrow \arg \max_a Q(s, a)$

For all  $a \in \mathcal{A}(s)$ :

$$\pi(a|s) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(s)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(s)| & \text{if } a \neq A^* \end{cases}$$

# Sarsa: On-Policy TD Control

---

Turn this into a control method by always updating the policy to be greedy with respect to the current estimate:

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

Repeat (for each step of episode):

Take action  $A$ , observe  $R, S'$

Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

$$Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma Q(S', A') - Q(S, A)]$$

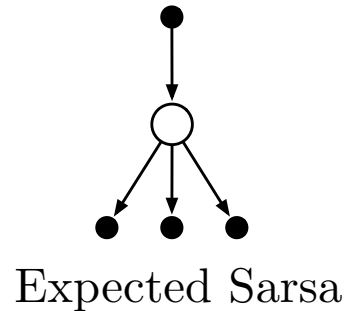
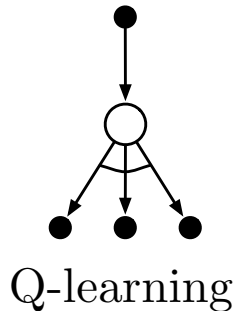
$S \leftarrow S'; A \leftarrow A';$

until  $S$  is terminal

# Expected Sarsa

- Instead of the *sample* value-of-next-state, use the expectation!

$$\begin{aligned} Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right] \\ &\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right] \end{aligned}$$

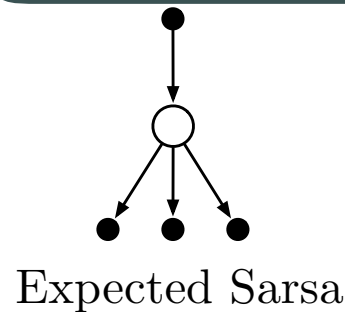
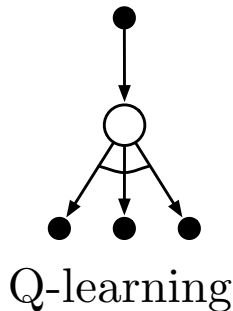


- Expected Sarsa's performs better than Sarsa (but costs more)

# Off-policy Expected Sarsa

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

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \mathbb{E}[Q(S_{t+1}, A_{t+1}) \mid S_{t+1}] - Q(S_t, A_t) \right]$$
$$\leftarrow Q(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q(S_t, A_t) \right]$$

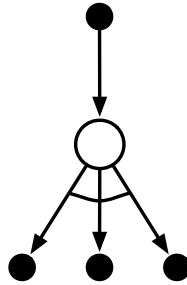


# Q-Learning: Off-Policy TD Control

---

One-step Q-learning:

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



Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

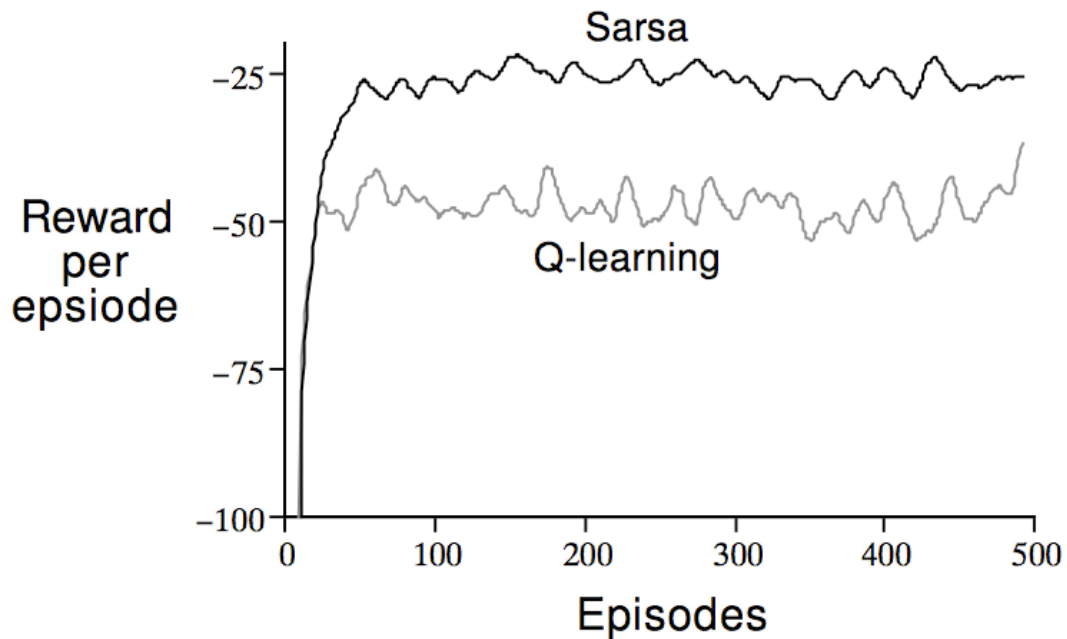
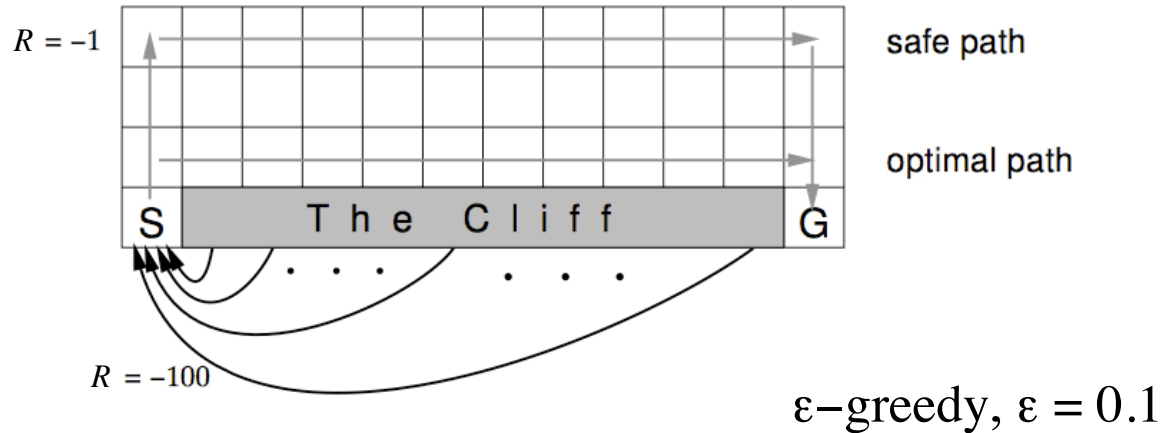
Take action  $A$ , observe  $R, S'$

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$$

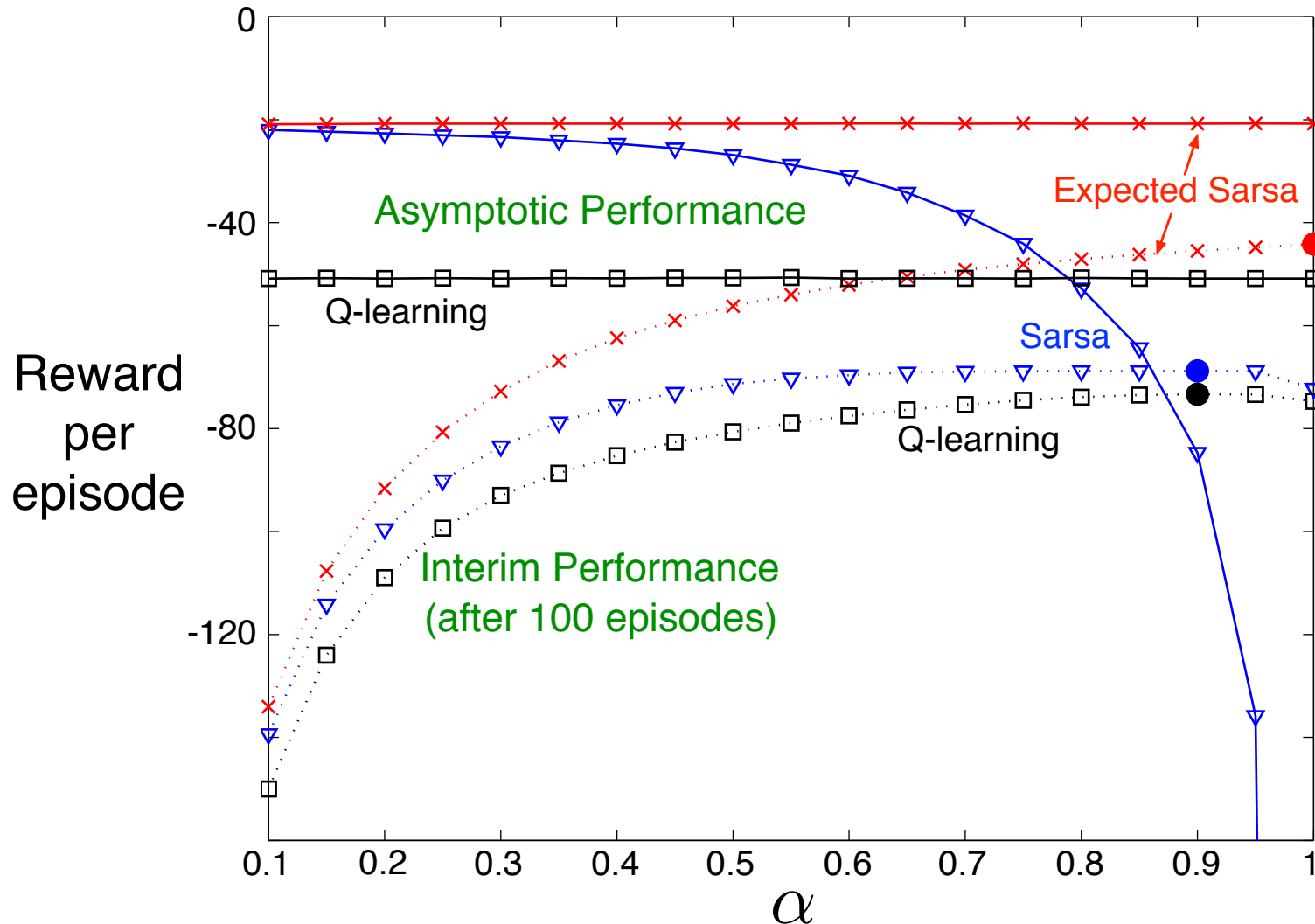
$S \leftarrow S'$ ;

until  $S$  is terminal

# Cliffwalking

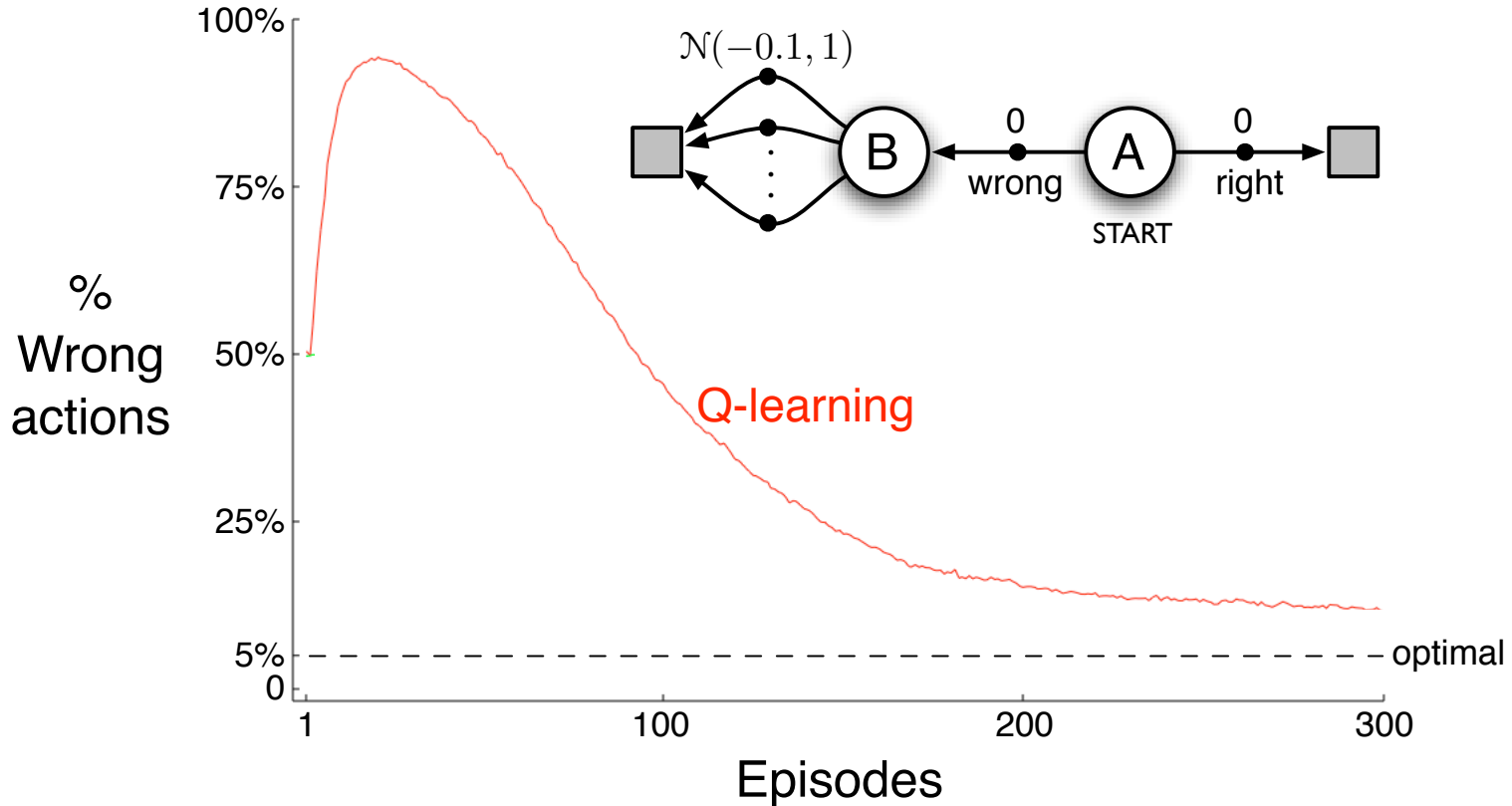


# Performance on the Cliff-walking Task





# Maximization Bias Example



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

# Double Q-Learning

- Train 2 action-value functions,  $Q_1$  and  $Q_2$
- Do Q-learning on both, but
  - never on the same time steps ( $Q_1$  and  $Q_2$  are indep.)
  - pick  $Q_1$  or  $Q_2$  at random to be updated on each step
- If updating  $Q_1$ , use  $Q_2$  for the value of the next state:

$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left( R_{t+1} + Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right)$$

- Action selections are (say)  $\varepsilon$ -greedy with respect to the sum of  $Q_1$  and  $Q_2$

# Double Q-Learning

Initialize  $Q_1(s, a)$  and  $Q_2(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily

Initialize  $Q_1(\text{terminal-state}, \cdot) = Q_2(\text{terminal-state}, \cdot) = 0$

Repeat (for each episode):

Initialize  $S$

Repeat (for each step of episode):

Choose  $A$  from  $S$  using policy derived from  $Q_1$  and  $Q_2$  (e.g.,  $\epsilon$ -greedy in  $Q_1 + Q_2$ )

Take action  $A$ , observe  $R, S'$

With 0.5 probability:

$$Q_1(S, A) \leftarrow Q_1(S, A) + \alpha \left( R + \gamma Q_2(S', \arg \max_a Q_1(S', a)) - Q_1(S, A) \right)$$

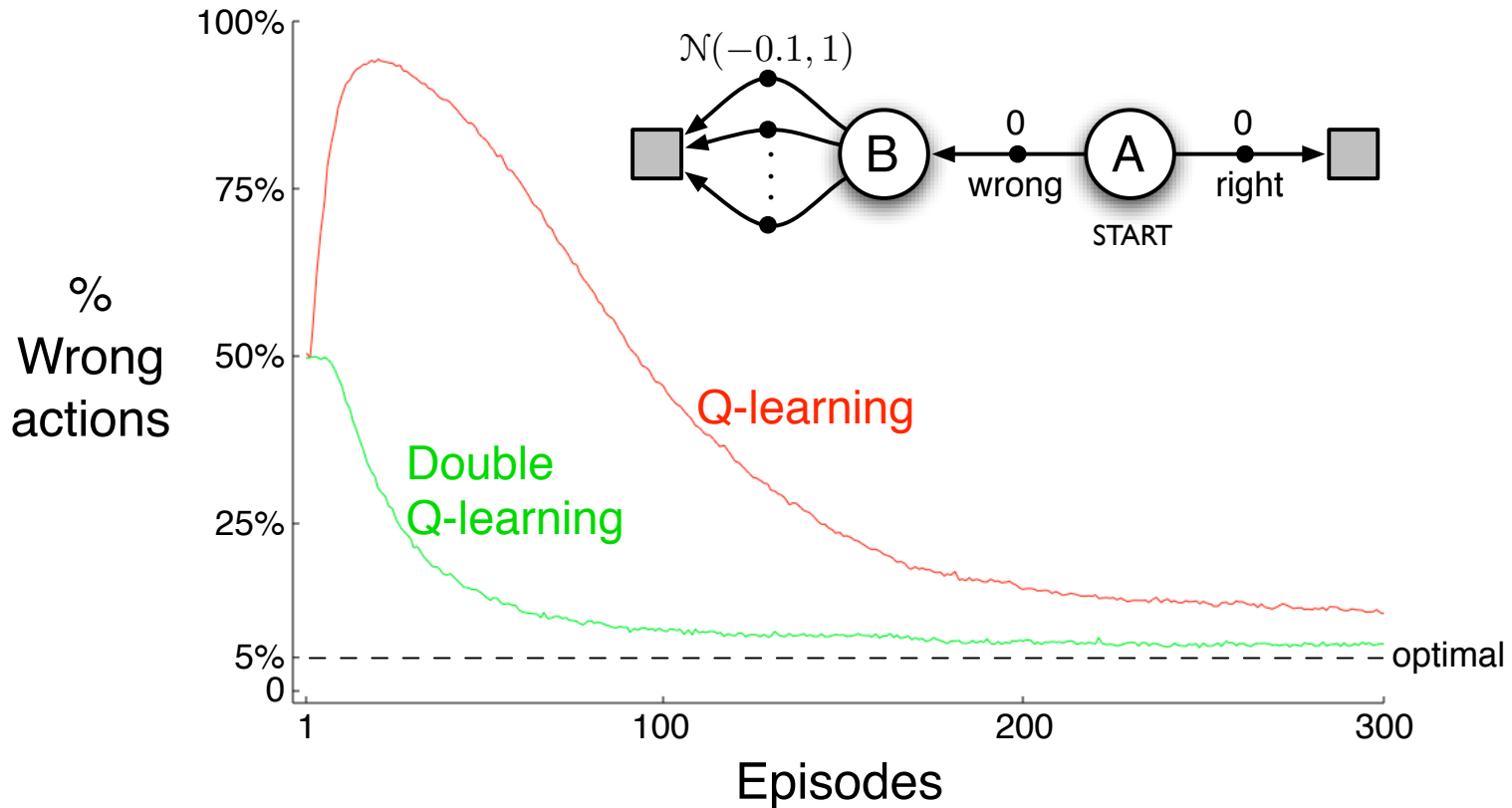
else:

$$Q_2(S, A) \leftarrow Q_2(S, A) + \alpha \left( R + \gamma Q_1(S', \arg \max_a Q_2(S', a)) - Q_2(S, A) \right)$$

$S \leftarrow S'$ ;

until  $S$  is terminal

# Example of Maximization Bias



Double Q-learning:

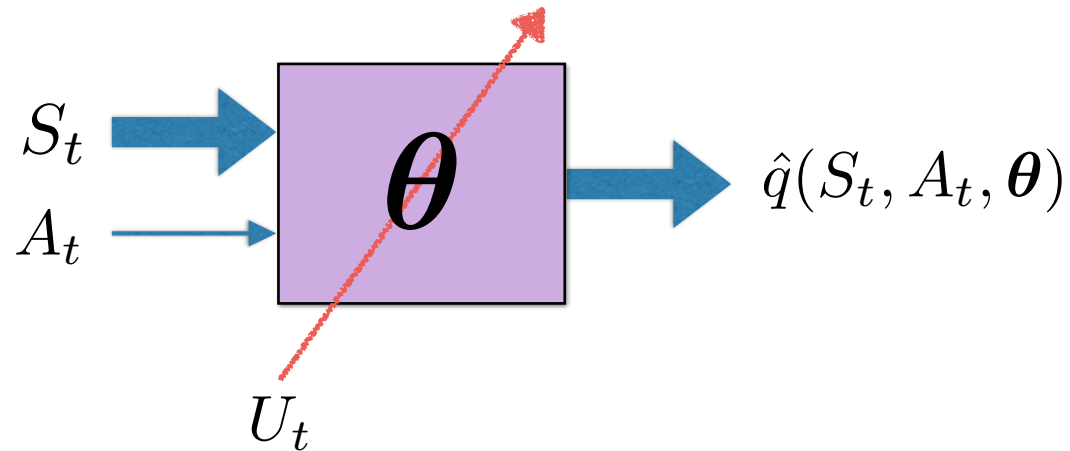
$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right]$$

# Summary

---

- Extend prediction to control by employing some form of GPI
  - On-policy control: **Sarsa, Expected Sarsa**
  - Off-policy control: **Q-learning, Expected Sarsa**
- Avoiding maximization bias with Double Q-learning

# Value function approximation (VFA) for control



# Recall: Different Targets

---

- **Monte Carlo:**  $G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$
- **TD:**  $G_t^{(1)} \doteq R_{t+1} + \gamma V_t(S_{t+1})$ 
  - Use  $V_t$  to estimate remaining return
- **$n$ -step TD:**
  - 2 step return:  $G_t^{(2)} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 V_t(S_{t+2})$
  - $n$ -step return:  $G_t^{(n)} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V_t(S_{t+n})$   
 $G_t^{(n)} \doteq G_t$  if  $t+n \geq T$

# Recall: Stochastic Gradient Descent (SGD)

General SGD:  $\theta \leftarrow \theta - \alpha \nabla_{\theta} Error_t^2$

For VFA:  $\leftarrow \theta - \alpha \nabla_{\theta} [Target_t - \hat{v}(S_t, \theta)]^2$

Chain rule:  $\leftarrow \theta - 2\alpha [Target_t - \hat{v}(S_t, \theta)] \nabla_{\theta} [Target_t - \hat{v}(S_t, \theta)]$

Semi-gradient:  $\leftarrow \theta + \alpha [Target_t - \hat{v}(S_t, \theta)] \nabla_{\theta} \hat{v}(S_t, \theta)$

Different RL algorithms provide different targets!

But share the “semi-gradient” aspect



# (Semi-)gradient methods carry over to control in the usual on-policy GPI way

- Always learn the action-value function of the current policy
- Always act near-greedily wrt the current action-value estimates
- The learning rule is:

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left[ U_t - \hat{q}(S_t, A_t, \boldsymbol{\theta}_t) \right] \nabla \hat{q}(S_t, A_t, \boldsymbol{\theta}_t)$$

update target, e.g.  $U_t = G_t$  (MC)

$U_t = R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \boldsymbol{\theta}_t)$  (Sarsa)

$$U_t = R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) \hat{q}(S_{t+1}, a, \boldsymbol{\theta}_t) \quad U_t = \sum_{s', r} p(s', r|S_t, A_t) \left[ r + \gamma \sum_{a'} \pi(a'|s') \hat{q}(s', a', \boldsymbol{\theta}_t) \right] \quad (\text{DP})$$

(Expected Sarsa)

# (Semi-)gradient methods carry over to control

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left[ U_t - \hat{q}(S_t, A_t, \boldsymbol{\theta}_t) \right] \nabla \hat{q}(S_t, A_t, \boldsymbol{\theta}_t)$$

## Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable function  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^n \rightarrow \mathbb{R}$

Initialize value-function weights  $\boldsymbol{\theta} \in \mathbb{R}^n$  arbitrarily (e.g.,  $\boldsymbol{\theta} = \mathbf{0}$ )

Repeat (for each episode):

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

    Repeat (for each step of episode):

        Take action  $A$ , observe  $R, S'$

        If  $S'$  is terminal:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha [R - \hat{q}(S, A, \boldsymbol{\theta})] \nabla \hat{q}(S, A, \boldsymbol{\theta})$$

        Go to next episode

        Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \boldsymbol{\theta})$  (e.g.,  $\varepsilon$ -greedy)

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha [R + \gamma \hat{q}(S', A', \boldsymbol{\theta}) - \hat{q}(S, A, \boldsymbol{\theta})] \nabla \hat{q}(S, A, \boldsymbol{\theta})$$

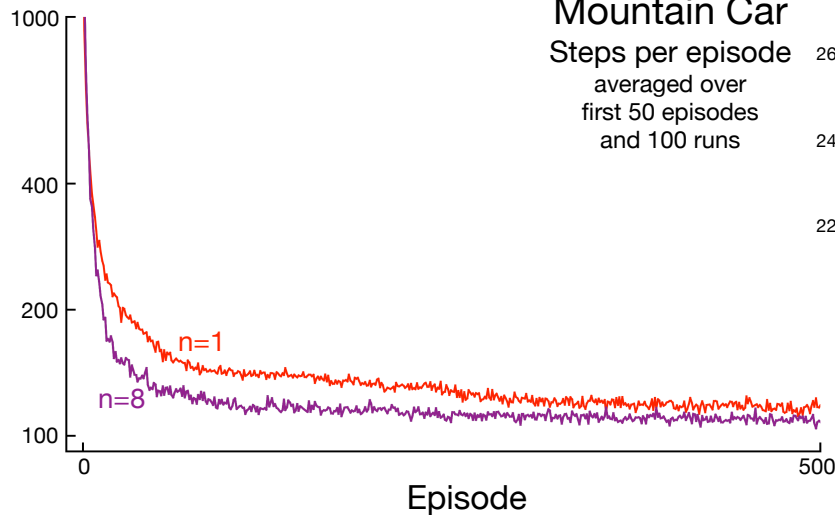
$S \leftarrow S'$

$A \leftarrow A'$

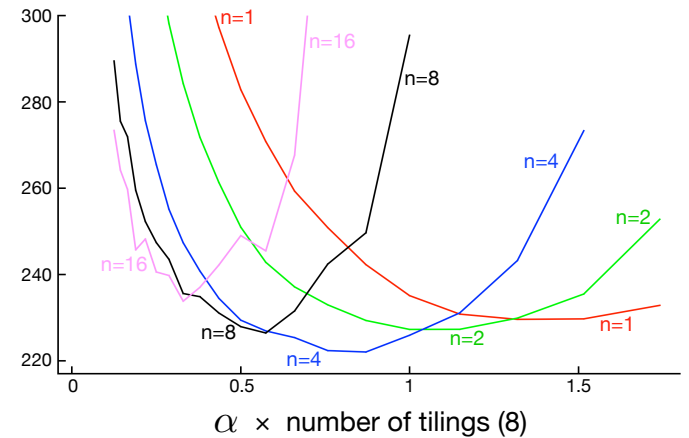
# $n$ -step semi-gradient Sarsa is better for $n > 1$

$$\theta_{t+n} \doteq \theta_{t+n-1} + \alpha \left[ G_t^{(n)} - \hat{q}(S_t, A_t, \theta_{t+n-1}) \right] \nabla \hat{q}(S_t, A_t, \theta_{t+n-1}), \quad 0 \leq t < T$$

Mountain Car  
Steps per episode  
log scale  
averaged over 100 runs



Mountain Car  
Steps per episode  
averaged over  
first 50 episodes  
and 100 runs



# Conclusions

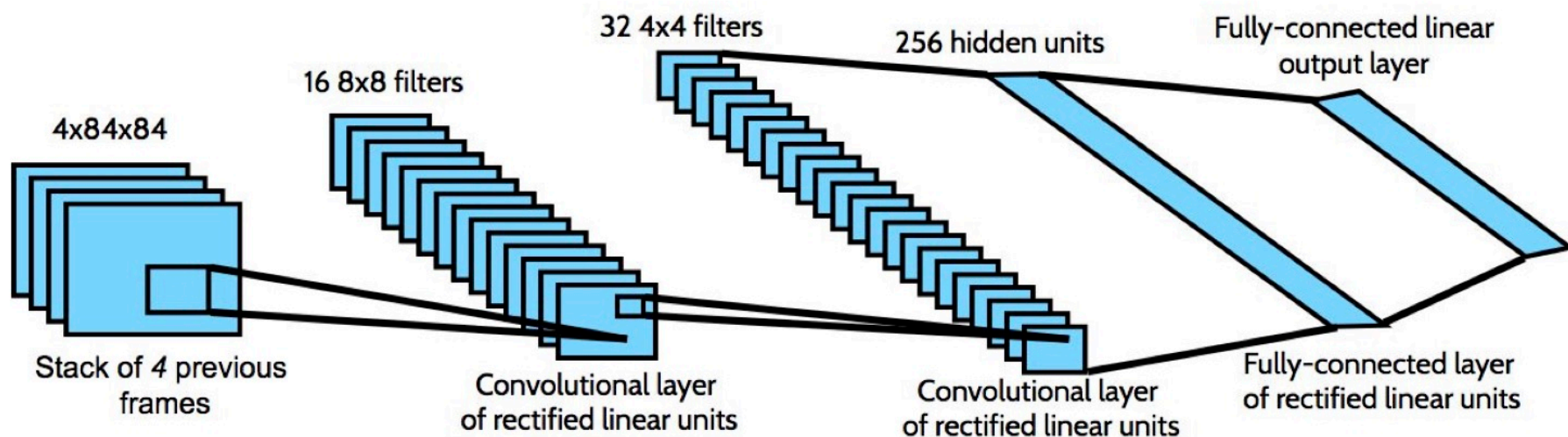
- **Control is straightforward** in the on-policy case
- **Formal results** (bounds) exist for the linear, on-policy case (eg. Gordon, 2000, Perkins & Precup, 2003 and follow-up work)
  - we get **chattering** near a good solution, **not convergence**

# DQN

(Mnih, Kavukcuoglu, Silver, et al., Nature 2015)

- Learns to play video games **from raw pixels**, simply by playing
- Can learn Q function by Q-learning

$$\Delta \mathbf{w} = \alpha \left( R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \mathbf{w}) - Q(S_t, A_t; \mathbf{w}) \right) \nabla_{\mathbf{w}} Q(S_t, A_t; \mathbf{w})$$



# DQN

(Mnih, Kavukcuoglu, Silver, et al., Nature 2015)

- Learns to play video games **from raw pixels**, simply by playing
- Can learn Q function by Q-learning

$$\Delta \mathbf{w} = \alpha \left( R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \mathbf{w}) - Q(S_t, A_t; \mathbf{w}) \right) \nabla_{\mathbf{w}} Q(S_t, A_t; \mathbf{w})$$

- Core components of DQN include:
  - Target networks (Mnih et al. 2015)

$$\Delta \mathbf{w} = \alpha \left( R_{t+1} + \gamma \max_a Q(S_{t+1}, a; \mathbf{w}^-) - Q(S_t, A_t; \mathbf{w}) \right) \nabla_{\mathbf{w}} Q(S_t, A_t; \mathbf{w})$$

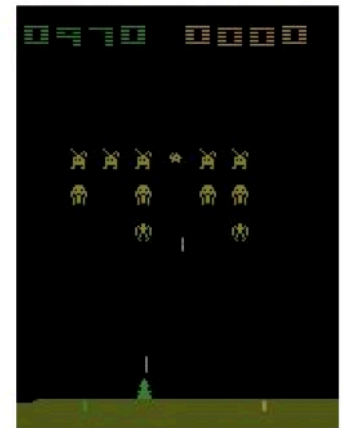
- Experience replay (Lin 1992): replay previous tuples (s, a, r, s')

# Target Network Intuition

(Slide credit: Vlad Mnih)

- Changing the value of one action will change the value of other actions and similar states.
- The network can end up chasing its own tail because of bootstrapping.
- Somewhat surprising fact - bigger networks are less prone to this because they alias less.

$$L_i(\theta_i) = \mathbb{E}_{s,a,s',r \sim D} \left( \underbrace{r + \gamma \max_{a'} Q(s', a'; \theta_i^-)}_{\text{target}} - Q(s, a; \theta_i) \right)^2$$



# DQN

(Mnih, Kavukcuoglu, Silver, et al., Nature 2015)

- Many later improvements to DQN
  - Double Q-learning (van Hasselt 2010, van Hasselt et al. 2015)
  - Prioritized replay (Schaul et al. 2016)
  - Dueling networks (Wang et al. 2016)
  - Asynchronous learning (Mnih et al. 2016)
  - Adaptive normalization of values (van Hasselt et al. 2016)
  - Better exploration (Bellemare et al. 2016, Ostrovski et al., 2017, Fortunato, Azar, Piot et al. 2017)
  - Distributional losses (Bellemare et al. 2017)
  - Multi-step returns (Mnih et al. 2016, Hessel et al. 2017)
  - ... many more ...

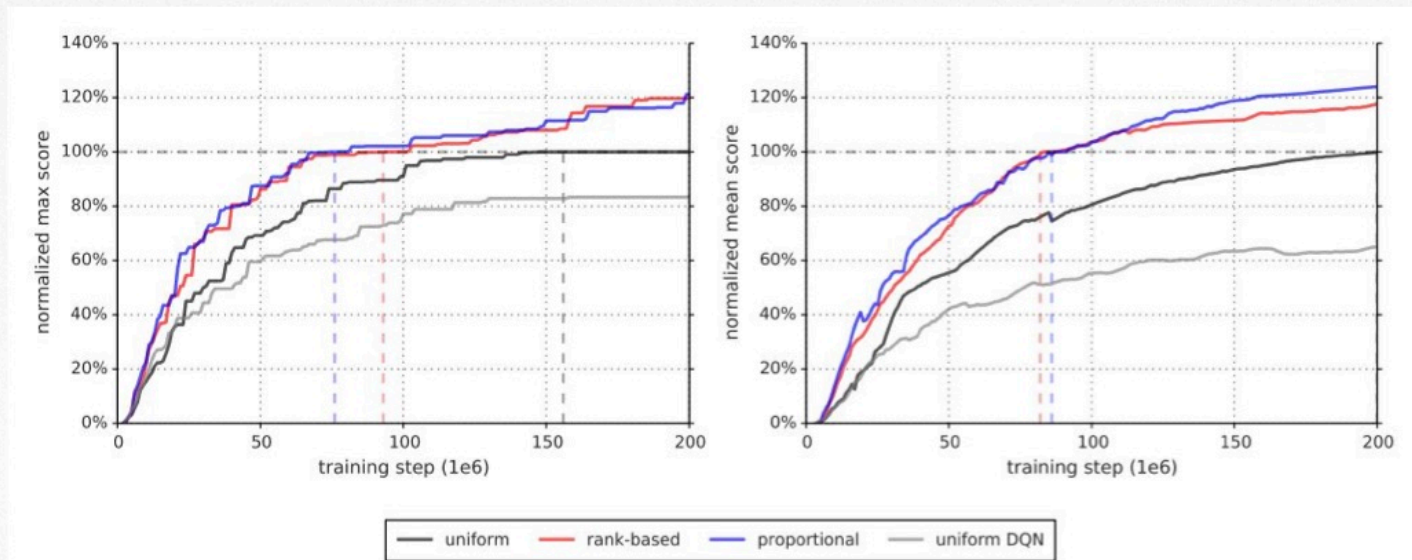


# Prioritized Experience Replay

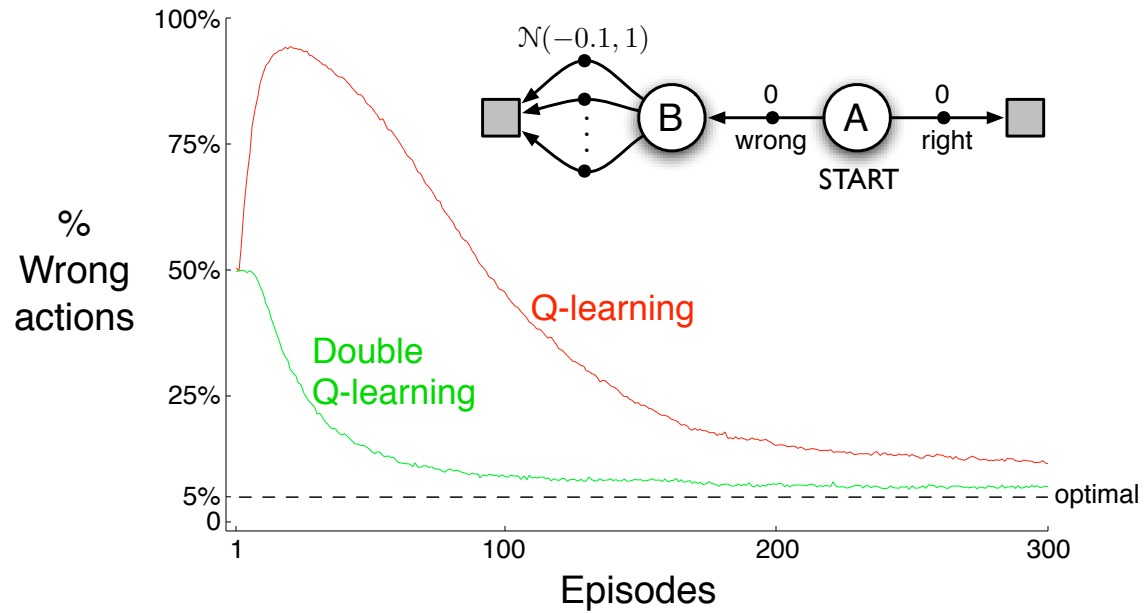
"Prioritized Experience Replay", Schaul et al. (2016)

- Idea: Replay transitions in proportion to TD error:

$$\left| r + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right|$$



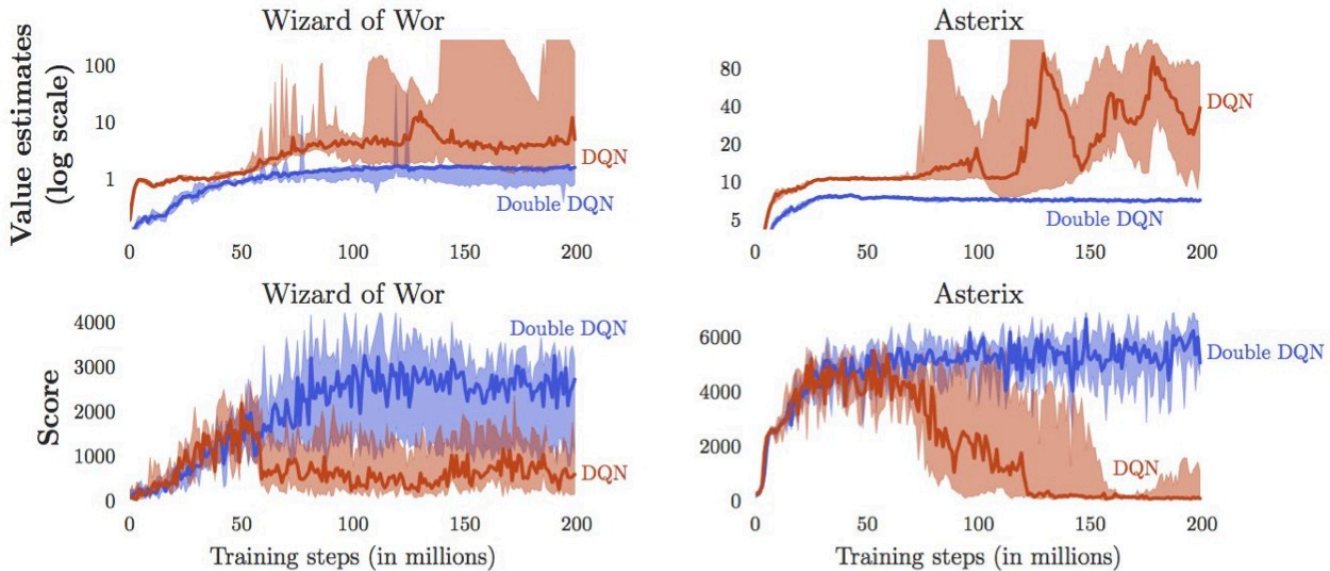
# Recall: Double DQN



## Double Q-learning:

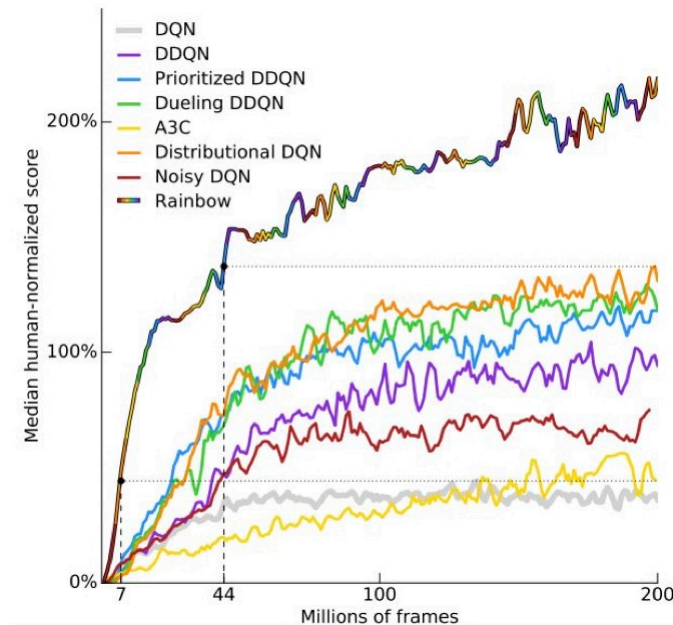
$$Q_1(S_t, A_t) \leftarrow Q_1(S_t, A_t) + \alpha \left[ R_{t+1} + \gamma Q_2(S_{t+1}, \arg \max_a Q_1(S_{t+1}, a)) - Q_1(S_t, A_t) \right]$$

# Double DQN



cf. van Hasselt et al, 2015)

# Which DQN improvements

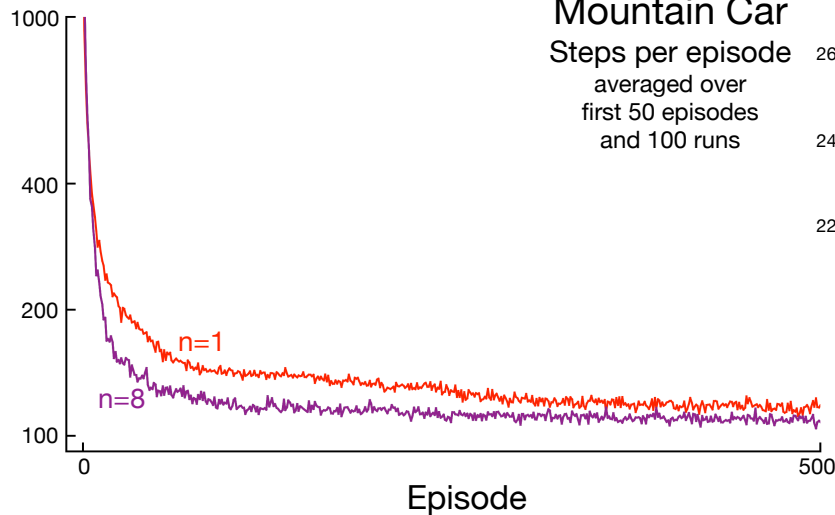


Rainbow model, (Hessel et al, 2017)

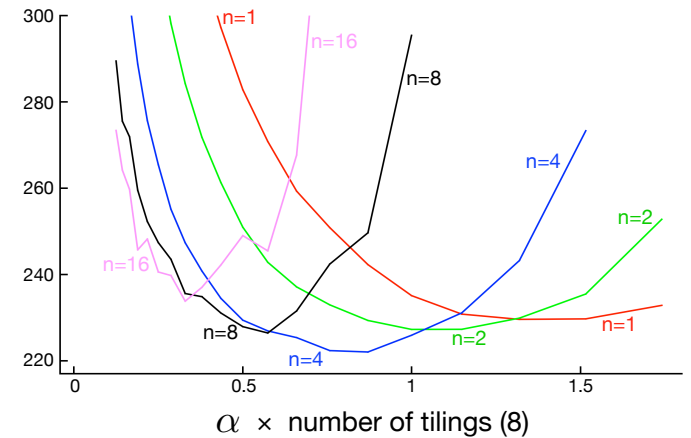
$n$ -step semi-gradient Sarsa is better for  $n > 1$

$$\theta_{t+n} \doteq \theta_{t+n-1} + \alpha \left[ G_t^{(n)} - \hat{q}(S_t, A_t, \theta_{t+n-1}) \right] \nabla \hat{q}(S_t, A_t, \theta_{t+n-1}), \quad 0 \leq t < T$$

Mountain Car  
Steps per episode  
log scale  
averaged over 100 runs



Mountain Car  
Steps per episode  
averaged over  
first 50 episodes  
and 100 runs



## Eligibility traces are

- Another way of interpolating between MC and TD methods
- A way of implementing *compound  $\lambda$ -return* targets
- A basic mechanistic idea — a short-term, fading memory
- A new style of algorithm development/analysis

## Recall $n$ -step targets

- For example, in the episodic case, with linear function approximation:

- 2-step target:

$$G_t^{(2)} \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 \boldsymbol{\theta}_{t+1}^\top \boldsymbol{\phi}_{t+2}$$

- $n$ -step target:  $G_t^{(n)} \doteq R_{t+1} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \boldsymbol{\theta}_{t+n-1}^\top \boldsymbol{\phi}_{t+n}$

with  $G_t^{(n)} \doteq G_t$  if  $t+n \geq T$

# Any set of update targets can be

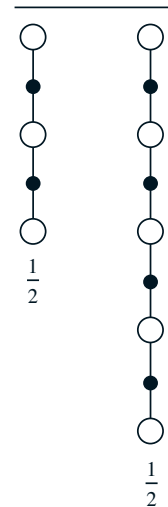
- For example, half a 2-step plus half a 4-

step

$$U_t = \frac{1}{2}G_t^{(2)} + \frac{1}{2}G_t^{(4)}$$

- Called a compound backup
  - Draw each component
  - Label with the weights for that

*A compound backup*



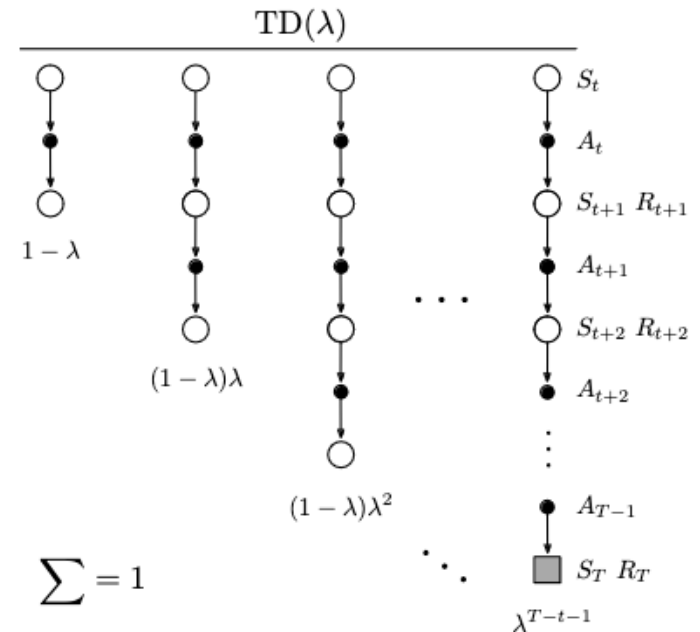


# The $\lambda$ -return is a compound update target

- The  $\lambda$ -return is a target that averages all  $n$ -step targets
- each weighted by  $\lambda^{n-1}$

$$G_t^\lambda \doteq (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n}$$

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{v}(S_{t+n}, \mathbf{w}_{t+n-1}), \quad 0 \leq t \leq T-n.$$



## Relation to TD(0) and MC

---

- The  $\lambda$ -return can be rewritten as:

$$G_t^\lambda = \underbrace{(1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_t^{(n)}}_{\text{Until termination}} + \underbrace{\lambda^{T-t-1} G_t}_{\text{After termination}}$$

- If  $\lambda = 1$ , you get the MC target:

$$G_t^\lambda = (1 - 1) \sum_{n=1}^{T-t-1} 1^{n-1} G_t^{(n)} + 1^{T-t-1} G_t = G_t$$

- If  $\lambda = 0$ , you get the TD(0) target:

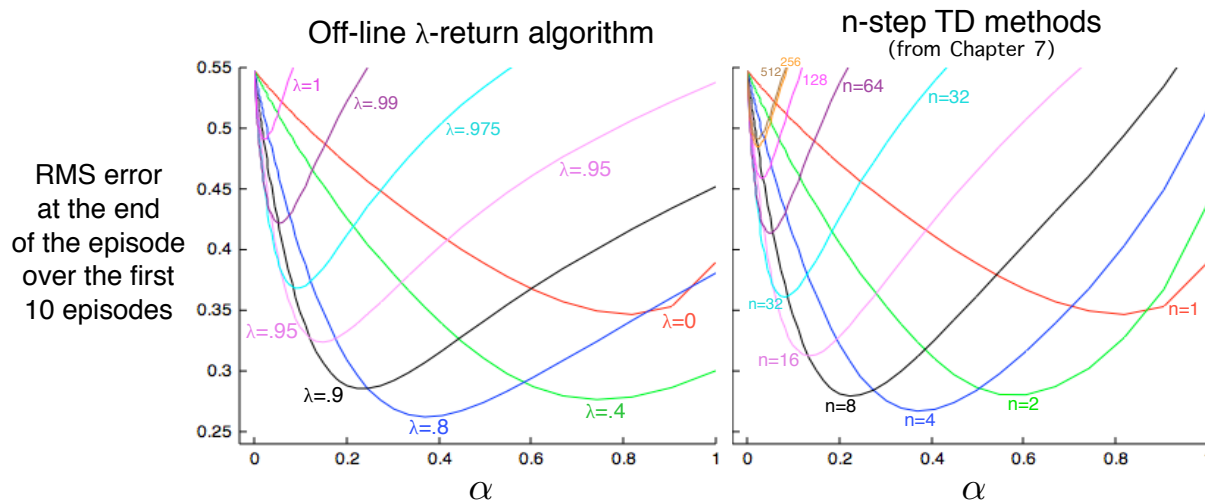
$$G_t^\lambda = (1 - 0) \sum_{n=1}^{T-t-1} 0^{n-1} G_t^{(n)} + 0^{T-t-1} G_t = G_t^{(1)} \quad 34$$

## The off-line $\lambda$ -return “algorithm”

- Wait until the end of the episode (offline)

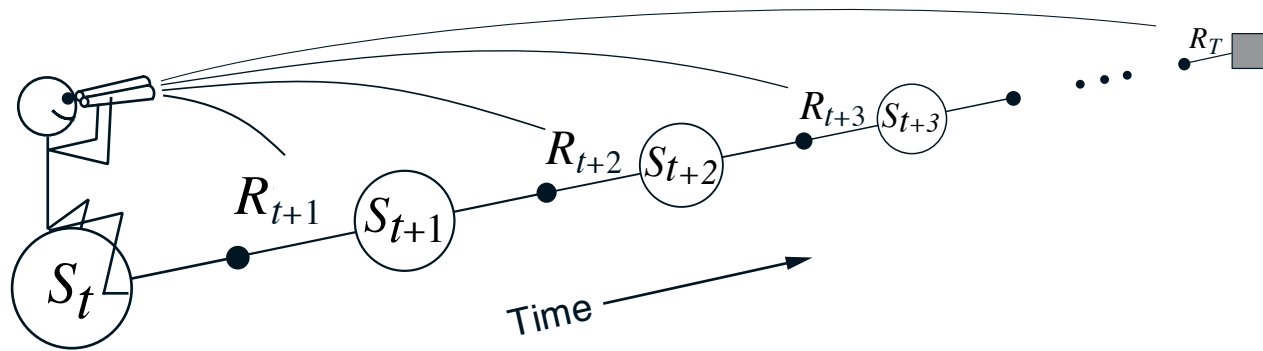
$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left[ G_t^\lambda - \hat{v}(S_t, \boldsymbol{\theta}_t) \right] \nabla \hat{v}(S_t, \boldsymbol{\theta}_t), \quad t = 0, \dots, T - 1$$

# The $\lambda$ -return alg performs similarly to

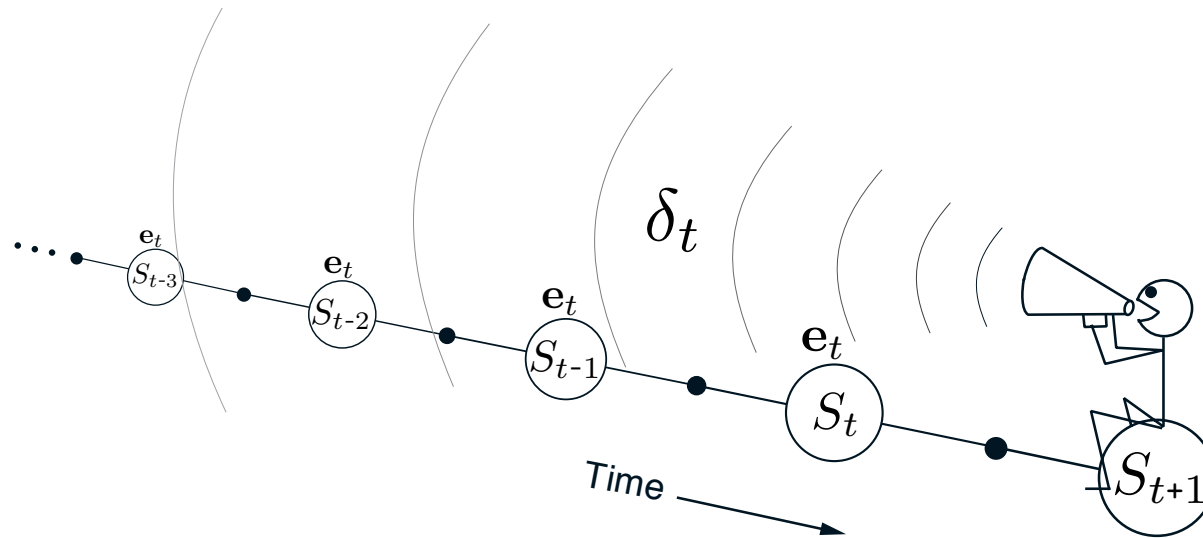


Intermediate  $\lambda$  is best (just like intermediate  $n$  is best)  
 $\lambda$ -return slightly better than  $n$ -step

The forward view looks forward from the state being updated to future states and rewards



The backward view looks back to the recently visited states (marked by eligibility traces)



- Shout the TD error backwards
- The traces fade with temporal distance by  $\gamma\lambda$

# Eligibility traces (mechanism)

---

- The forward view was for theory
- The backward view is for *mechanism* same shape as  $\theta$   
$$\mathbf{e}_t \in \mathbb{R}^n \geq \mathbf{0}$$
- New memory vector called *eligibility trace*
  - On each step, decay each component by  $\gamma\lambda$  and increment the trace for the current state by 1
  - *Accumulating trace*

$$\mathbf{e}_0 \doteq \mathbf{0},$$
$$\mathbf{e}_t \doteq \nabla \hat{v}(S_t, \boldsymbol{\theta}_t) + \gamma\lambda \mathbf{e}_{t-1}$$

# The Semi-gradient TD( $\lambda$ ) algorithm

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \delta_t \mathbf{e}_t$$

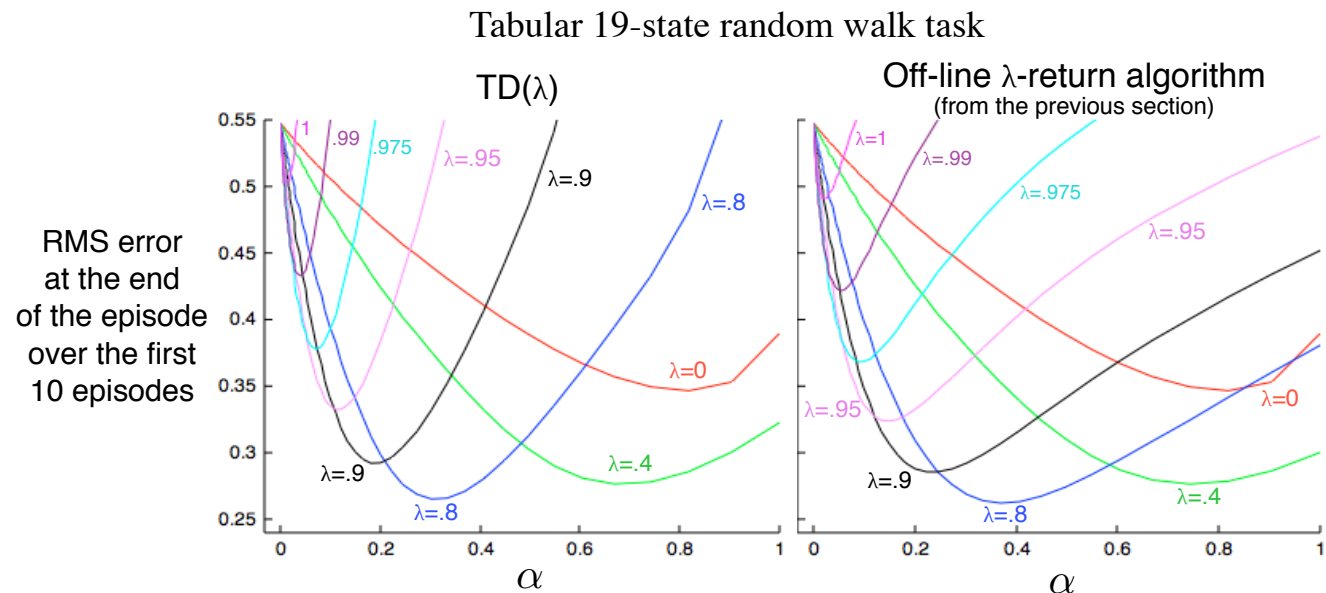
$$\delta_t \doteq R_{t+1} + \gamma \hat{v}(S_{t+1}, \boldsymbol{\theta}_t) - \hat{v}(S_t, \boldsymbol{\theta}_t)$$

$$\mathbf{e}_0 \doteq \mathbf{0},$$

$$\mathbf{e}_t \doteq \nabla \hat{v}(S_t, \boldsymbol{\theta}_t) + \gamma \lambda \mathbf{e}_{t-1}$$



# TD( $\lambda$ ) performs similarly to offline $\lambda$ -



Can we do better? Can we update online?

# Conclusions

- Value-function approximation by stochastic gradient descent enables RL to be applied to arbitrarily large state spaces
- Most algorithms just carry over the targets from the tabular case
- With bootstrapping (TD), we don't get true gradient descent methods
  - this complicates the analysis
  - but the linear, on-policy case is still guaranteed convergent
  - and learning is still *much faster*