

# 10703 Deep Reinforcement Learning and Control

Russ Salakhutdinov

Machine Learning Department  
rsalakhu@cs.cmu.edu

## Function Approximation

# Used Materials

- **Disclaimer:** Much of the material and slides for this lecture were borrowed from Rich Sutton's class and David Silver's class on Reinforcement Learning.

# Large-Scale Reinforcement Learning

- ▶ Reinforcement learning can be used to solve large problems, e.g.
  - Backgammon:  $10^{20}$  states
  - Computer Go:  $10^{170}$  states
  - Helicopter: continuous state space
- ▶ How can we scale up the **model-free methods** for prediction and control?

# Value Function Approximation (VFA)

- ▶ So far we have represented value function by a **lookup table**
  - Every **state**  $s$  has an entry  $V(s)$ , or
  - Every **state-action** pair  $(s,a)$  has an entry  $Q(s,a)$
- ▶ Problem with large MDPs:
  - There are too many states and/or actions to store in memory
  - It is too slow to learn the value of each state individually
- ▶ Solution for large MDPs:
  - Estimate value function with **function approximation**

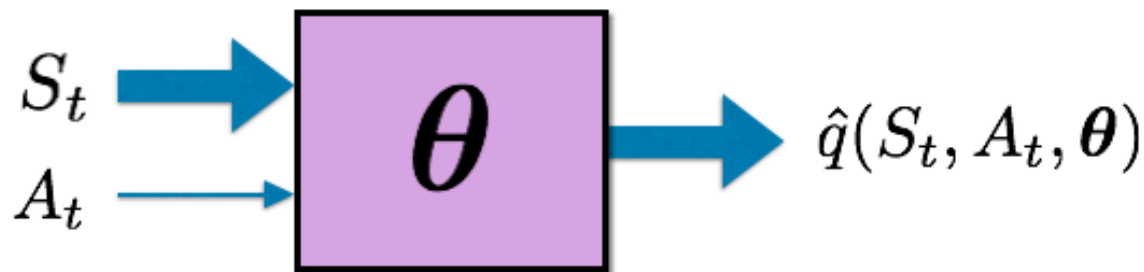
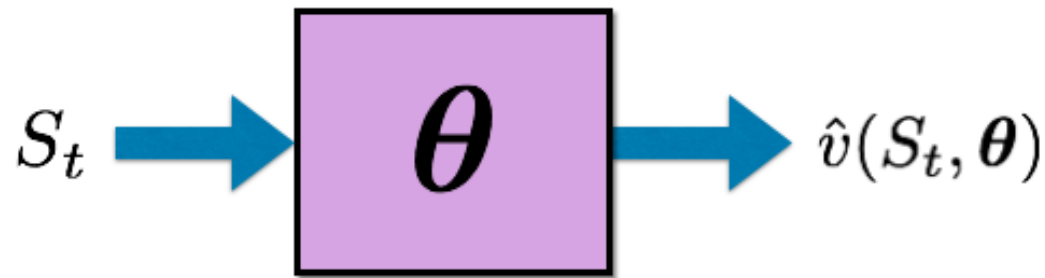
$$\hat{v}(s, \mathbf{w}) \approx v_{\pi}(s)$$

$$\text{or } \hat{q}(s, a, \mathbf{w}) \approx q_{\pi}(s, a)$$

- Generalize from seen states to unseen states

# Value Function Approximation (VFA)

- ▶ Value function approximation (VFA) replaces the table with a general parameterized form:



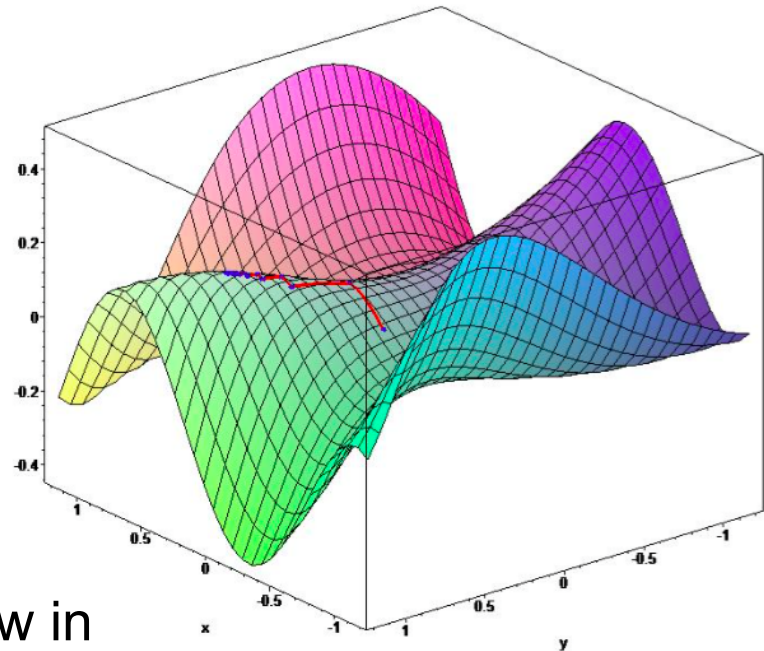
# Which Function Approximation?

- ▶ There are many **function approximators**, e.g.
  - Linear combinations of features
  - Neural networks
  - Decision tree
  - Nearest neighbour
  - Fourier / wavelet bases
  - ...
- ▶ We consider **differentiable function approximators**, e.g.
  - Linear combinations of features
  - Neural networks

# Gradient Descent

- ▶ Let  $J(\mathbf{w})$  be a **differentiable function** of parameter vector  $\mathbf{w}$
- ▶ Define the gradient of  $J(\mathbf{w})$  to be:

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_n} \end{pmatrix}$$



- ▶ To find a local minimum of  $J(\mathbf{w})$ , adjust  $\mathbf{w}$  in direction of the **negative gradient**:

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

Step-size

# Stochastic Gradient Descent

- ▶ **Goal:** find parameter vector  $\mathbf{w}$  minimizing mean-squared error between the **true value function**  $v_\pi(S)$  and its **approximation**  $\hat{v}(S, \mathbf{w})$ :

$$J(\mathbf{w}) = \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, \mathbf{w}))^2]$$

- ▶ Gradient descent finds a local minimum:

$$\begin{aligned}\Delta \mathbf{w} &= -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w}) \\ &= \alpha \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})]\end{aligned}$$

- ▶ **Stochastic gradient descent (SGD)** samples the gradient:

$$\Delta \mathbf{w} = \alpha (v_\pi(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$$

- ▶ Expected update is equal to full gradient update



# Feature Vectors

- ▶ Represent state by a **feature vector**

$$\mathbf{x}(S) = \begin{pmatrix} \mathbf{x}_1(S) \\ \vdots \\ \mathbf{x}_n(S) \end{pmatrix}$$

- ▶ For example
  - Distance of robot from landmarks
  - Trends in the stock market
  - Piece and pawn configurations in chess

# Linear Value Function Approximation (VFA)

- ▶ Represent **value function** by a linear combination of features

$$\hat{v}(S, \mathbf{w}) = \mathbf{x}(S)^\top \mathbf{w} = \sum_{j=1}^n \mathbf{x}_j(S) \mathbf{w}_j$$

- ▶ Objective function is **quadratic in parameters**  $\mathbf{w}$

$$J(\mathbf{w}) = \mathbb{E}_\pi \left[ (v_\pi(S) - \mathbf{x}(S)^\top \mathbf{w})^2 \right]$$

- ▶ Update rule is particularly simple

$$\nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}) = \mathbf{x}(S)$$

$$\Delta \mathbf{w} = \alpha (v_\pi(S) - \hat{v}(S, \mathbf{w})) \mathbf{x}(S)$$

- ▶ **Update** = step-size  $\times$  prediction error  $\times$  feature value
- ▶ Later, we will look at the neural networks as function approximators.

# Incremental Prediction Algorithms

- ▶ We have assumed the **true value function**  $v_{\pi}(s)$  is given by a supervisor
- ▶ But in RL there is no supervisor, only rewards
- ▶ In practice, we substitute a target for  $v_{\pi}(s)$
- ▶ For MC, the target is the **return**  $G_t$

$$\Delta \mathbf{w} = \alpha(\mathbf{G}_t - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$$

- ▶ For TD(0), the target is the **TD target**:  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$

$$\Delta \mathbf{w} = \alpha(\mathbf{R}_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w})$$

Remember  $\Delta \mathbf{w} = \alpha(v_{\pi}(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$

# Monte Carlo with VFA

- ▶ Return  $G_t$  is an **unbiased**, noisy sample of true value  $v_{\pi}(S_t)$
- ▶ Can therefore apply supervised learning to “**training data**”:

$$\langle S_1, G_1 \rangle, \langle S_2, G_2 \rangle, \dots, \langle S_T, G_T \rangle$$

- ▶ For example, using **linear Monte-Carlo policy evaluation**

$$\begin{aligned}\Delta \mathbf{w} &= \alpha(\mathbf{G}_t - \hat{v}(S_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S_t, \mathbf{w}) \\ &= \alpha(G_t - \hat{v}(S_t, \mathbf{w})) \mathbf{x}(S_t)\end{aligned}$$

- ▶ Monte-Carlo evaluation converges to a local optimum

# Monte Carlo with VFA

## Gradient Monte Carlo Algorithm for Approximating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

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

Initialize value-function weights  $\boldsymbol{\theta}$  as appropriate (e.g.,  $\boldsymbol{\theta} = \mathbf{0}$ )

Repeat forever:

    Generate an episode  $S_0, A_0, R_1, S_1, A_1, \dots, R_T, S_T$  using  $\pi$

    For  $t = 0, 1, \dots, T - 1$ :

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha [G_t - \hat{v}(S_t, \boldsymbol{\theta})] \nabla \hat{v}(S_t, \boldsymbol{\theta})$$

# TD Learning with VFA

- ▶ The TD-target  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$  is a **biased sample** of true value  $v_{\pi}(S_t)$

- ▶ Can still apply supervised learning to “**training data**”:

$$\langle S_1, R_2 + \gamma \hat{v}(S_2, \mathbf{w}) \rangle, \langle S_2, R_3 + \gamma \hat{v}(S_3, \mathbf{w}) \rangle, \dots, \langle S_{T-1}, R_T \rangle$$

- ▶ For example, using **linear TD(0)**:

$$\begin{aligned} \Delta \mathbf{w} &= \alpha (\mathbf{R} + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w}) \\ &= \alpha \delta \mathbf{x}(S) \end{aligned}$$

# TD Learning with VFA

Semi-gradient TD(0) for estimating  $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^n \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$

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

Repeat (for each episode):

    Initialize  $S$

    Repeat (for each step of episode):

        Choose  $A \sim \pi(\cdot | S)$

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

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

$S \leftarrow S'$

    until  $S'$  is terminal

# Control with VFA

- ▶ Policy evaluation **Approximate** policy evaluation:  $\hat{q}(\cdot, \cdot, \mathbf{w}) \approx q_\pi$
- ▶ Policy improvement  $\epsilon$ -greedy policy improvement



# Action-Value Function Approximation

- ▶ Approximate the **action-value function**

$$\hat{q}(S, A, \mathbf{w}) \approx q_{\pi}(S, A)$$

- ▶ Minimize **mean-squared error** between the true action-value function  $q_{\pi}(S, A)$  and the approximate action-value function:

$$J(\mathbf{w}) = \mathbb{E}_{\pi} [(q_{\pi}(S, A) - \hat{q}(S, A, \mathbf{w}))^2]$$

- ▶ Use **stochastic gradient descent** to find a local minimum

$$-\frac{1}{2} \nabla_{\mathbf{w}} J(\mathbf{w}) = (q_{\pi}(S, A) - \hat{q}(S, A, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w})$$

$$\Delta \mathbf{w} = \alpha (q_{\pi}(S, A) - \hat{q}(S, A, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w})$$

# Linear Action-Value Function Approximation

- ▶ Represent state and action by a **feature vector**

$$\mathbf{x}(S, A) = \begin{pmatrix} \mathbf{x}_1(S, A) \\ \vdots \\ \mathbf{x}_n(S, A) \end{pmatrix}$$

- ▶ Represent action-value function by **linear combination of features**

$$\hat{q}(S, A, \mathbf{w}) = \mathbf{x}(S, A)^\top \mathbf{w} = \sum_{j=1}^n \mathbf{x}_j(S, A) \mathbf{w}_j$$

- ▶ **Stochastic gradient descent** update

$$\nabla_{\mathbf{w}} \hat{q}(S, A, \mathbf{w}) = \mathbf{x}(S, A)$$

$$\Delta \mathbf{w} = \alpha (q_\pi(S, A) - \hat{q}(S, A, \mathbf{w})) \mathbf{x}(S, A)$$

# Incremental Control Algorithms

- ▶ Like prediction, we must substitute a target for  $q_{\pi}(S,A)$
- ▶ For MC, the target is the return  $G_t$

$$\Delta \mathbf{w} = \alpha(\mathbf{G}_t - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

- ▶ For TD(0), the target is the TD target:  $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$

$$\Delta \mathbf{w} = \alpha(\mathbf{R}_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(S_t, A_t, \mathbf{w})$$

# Incremental Control Algorithms

## 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

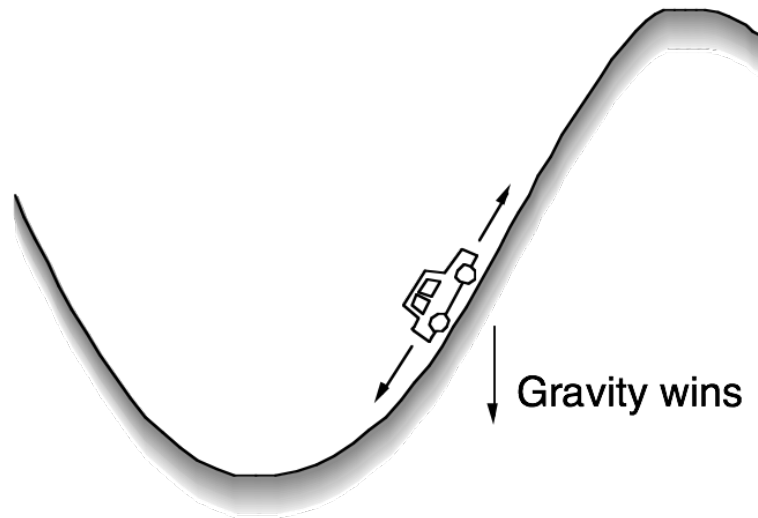
    Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \boldsymbol{\theta})$  (e.g.,  $\varepsilon$ -greedy)

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

$S \leftarrow S'$

$A \leftarrow A'$

# Example: The Mountain-Car problem



Minimum-Time-to-Goal Problem

## SITUATIONS:

car's position and velocity

## ACTIONS:

three thrusts: forward, reverse, none

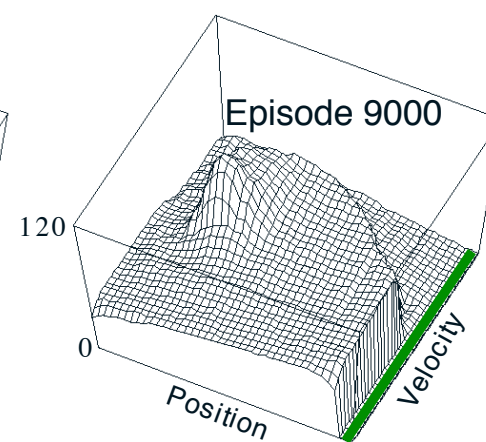
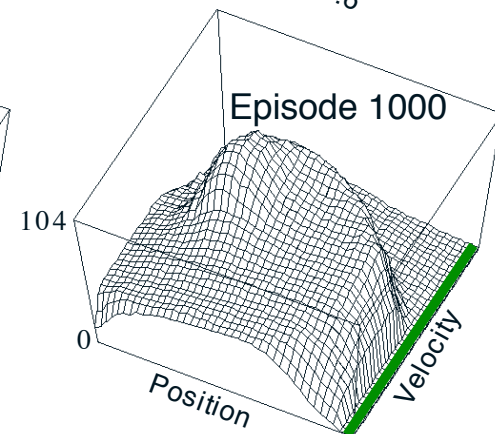
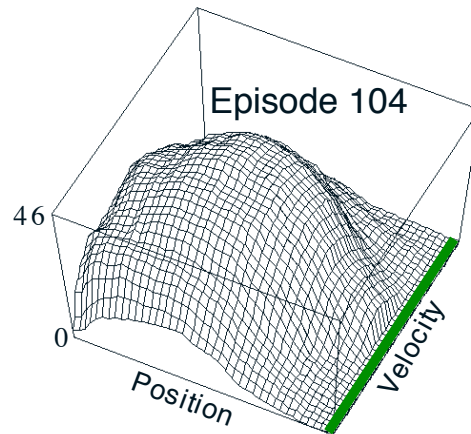
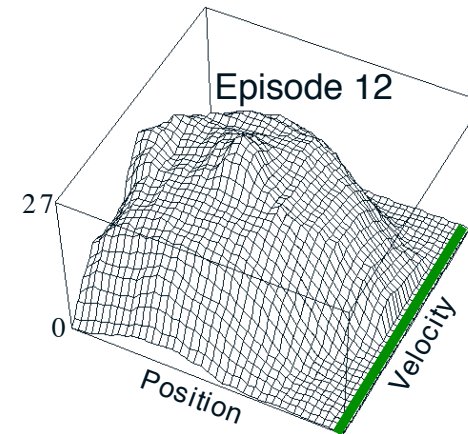
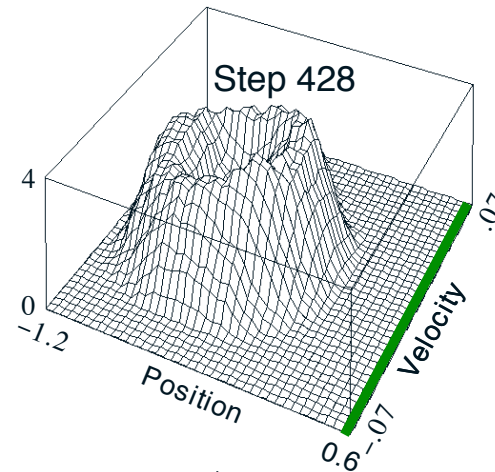
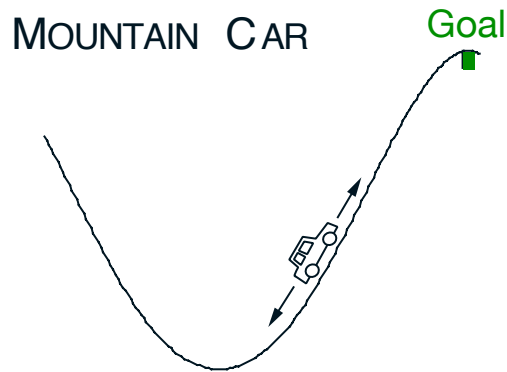
## REWARDS:

always  $-1$  until car reaches the goal

Episodic, No Discounting,  $\gamma=1$

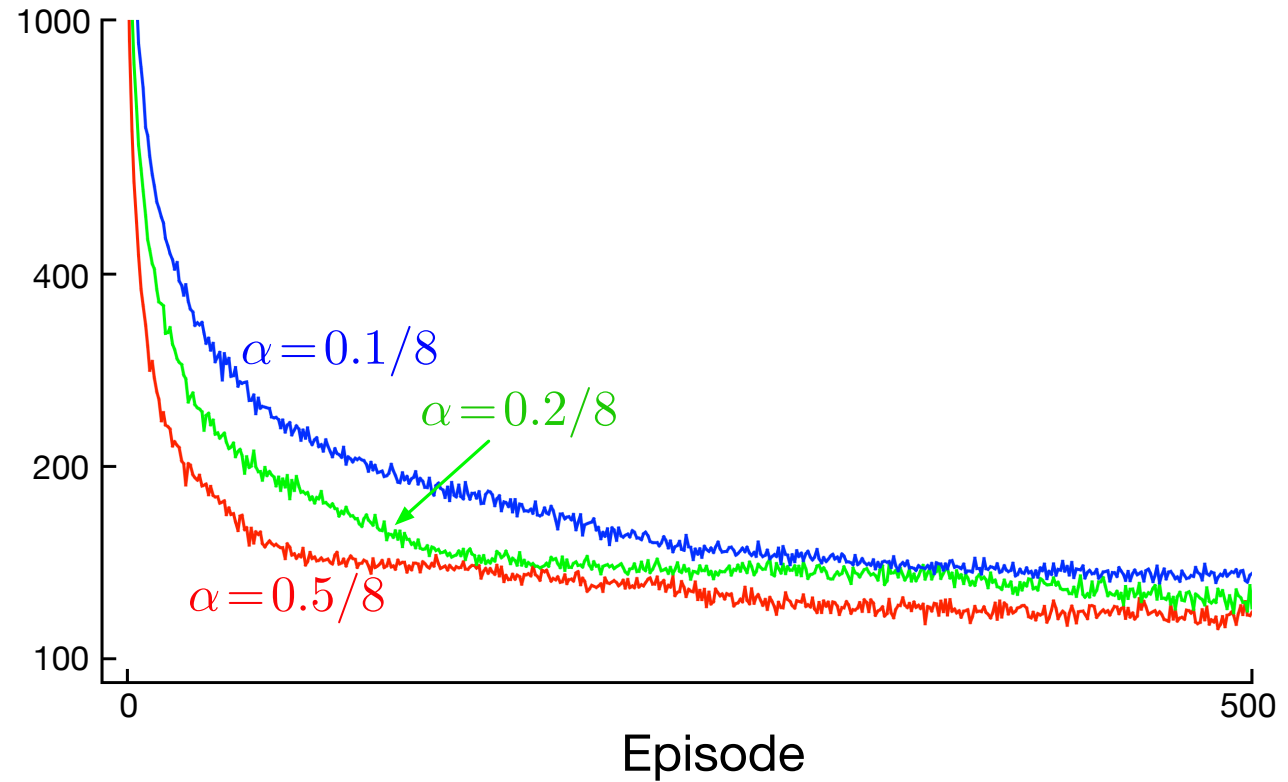
# Example: The Mountain-Car problem

$$- \max_a \hat{q}(s, a, \theta)$$



# Linear Sarsa: Mountain Car

Mountain Car  
Steps per episode  
log scale  
averaged over 100 runs



# Batch Reinforcement Learning

- ▶ Gradient descent is simple and appealing
- ▶ But it is not **sample efficient**
- ▶ Batch methods seek to find the best fitting value function
- ▶ Given the agent's **experience** (“training data”)



# Least Squares Prediction

- ▶ Given **value function approximation**:  $\hat{v}(s, \mathbf{w}) \approx v_{\pi}(s)$
- ▶ And **experience**  $\mathcal{D}$  consisting of  $\langle \text{state}, \text{value} \rangle$  pairs

$$\mathcal{D} = \{ \langle s_1, v_1^{\pi} \rangle, \langle s_2, v_2^{\pi} \rangle, \dots, \langle s_T, v_T^{\pi} \rangle \}$$

- ▶ Find parameters  $\mathbf{w}$  that give the best fitting value function  $v(s, \mathbf{w})$ ?
- ▶ Least squares **algorithms** find parameter vector  $\mathbf{w}$  minimizing sum-squared error between  $v(s_t, \mathbf{w})$  and target values  $v_t^{\pi}$ :

$$\begin{aligned} LS(\mathbf{w}) &= \sum_{t=1}^T (v_t^{\pi} - \hat{v}(s_t, \mathbf{w}))^2 \\ &= \mathbb{E}_{\mathcal{D}} [(v^{\pi} - \hat{v}(s, \mathbf{w}))^2] \end{aligned}$$

# SGD with Experience Replay

- ▶ Given **experience** consisting of  $\langle \text{state}, \text{value} \rangle$  pairs

$$\mathcal{D} = \{ \langle s_1, v_1^\pi \rangle, \langle s_2, v_2^\pi \rangle, \dots, \langle s_T, v_T^\pi \rangle \}$$

- ▶ Repeat
  - Sample state, value from experience

$$\langle s, v^\pi \rangle \sim \mathcal{D}$$

- Apply stochastic gradient descent update

$$\Delta \mathbf{w} = \alpha (v^\pi - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$$

- ▶ Converges to least squares solution
- ▶ We will look at Deep Q-networks later.