

Carnegie Mellon

School of Computer Science

Deep Reinforcement Learning and Control

Function Approximation for (on policy) Prediction and Control

Lecture 8, CMU 10-403

Katerina Fragkiadaki



Used Materials

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

Large-Scale Reinforcement Learning

- ▶ Reinforcement learning has been used to solve large problems, e.g.
 - Backgammon: 10^{20} states
 - Computer Go: 10^{170} states
 - Helicopter: continuous state space

- ▶ Tabular methods clearly do not work

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**

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

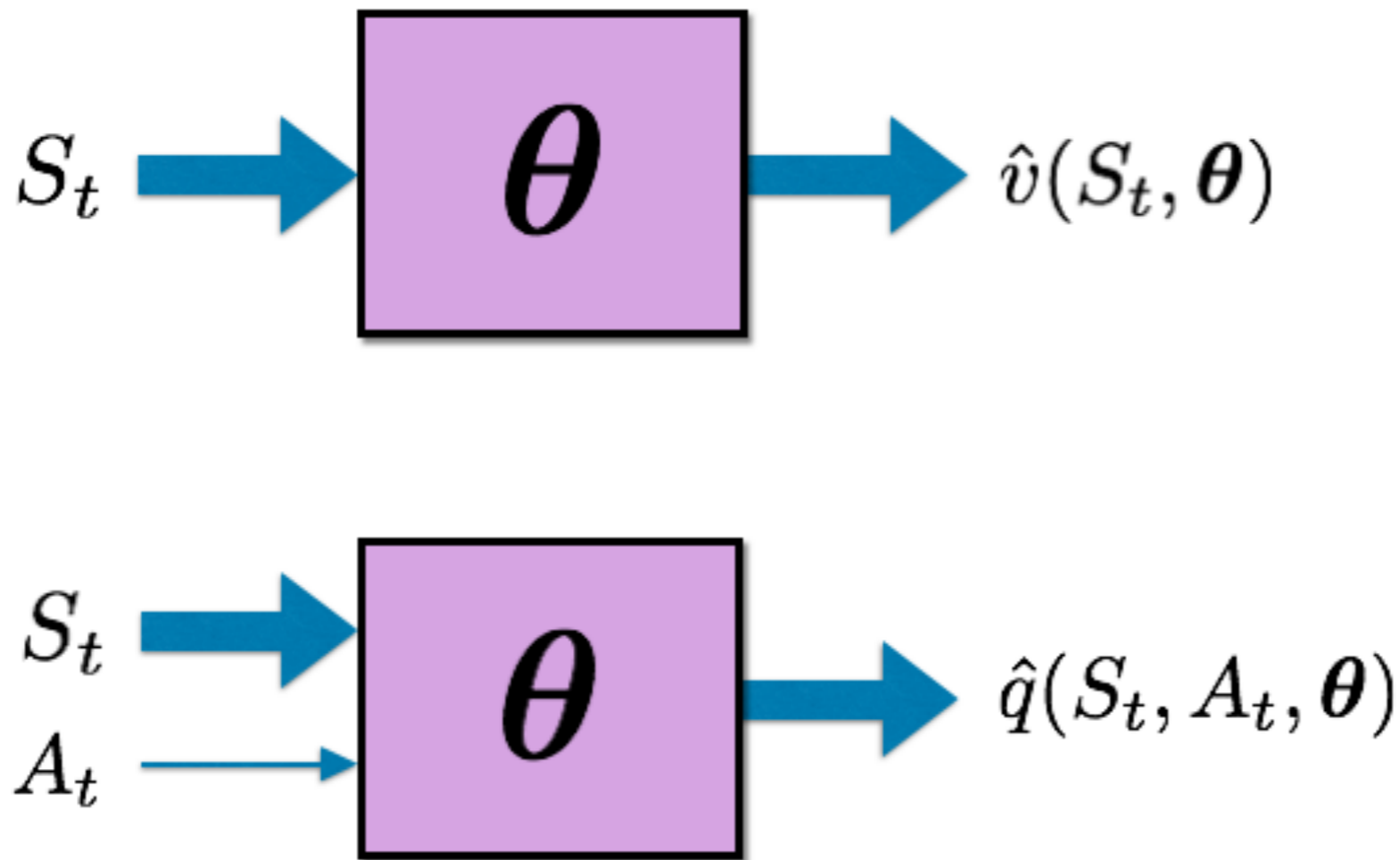
or

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

- Generalize from seen states to unseen states

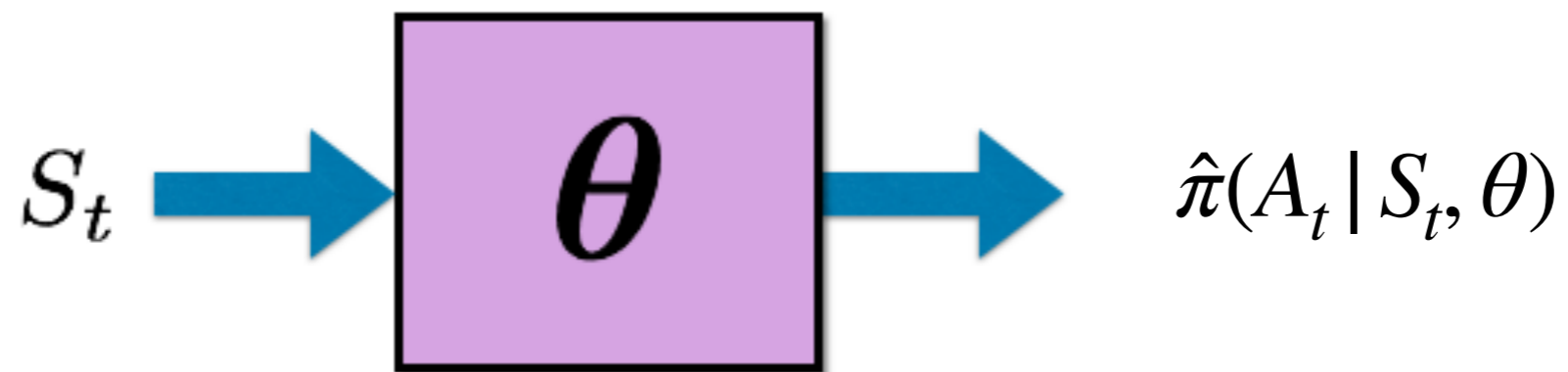
Value Function Approximation (VFA)

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



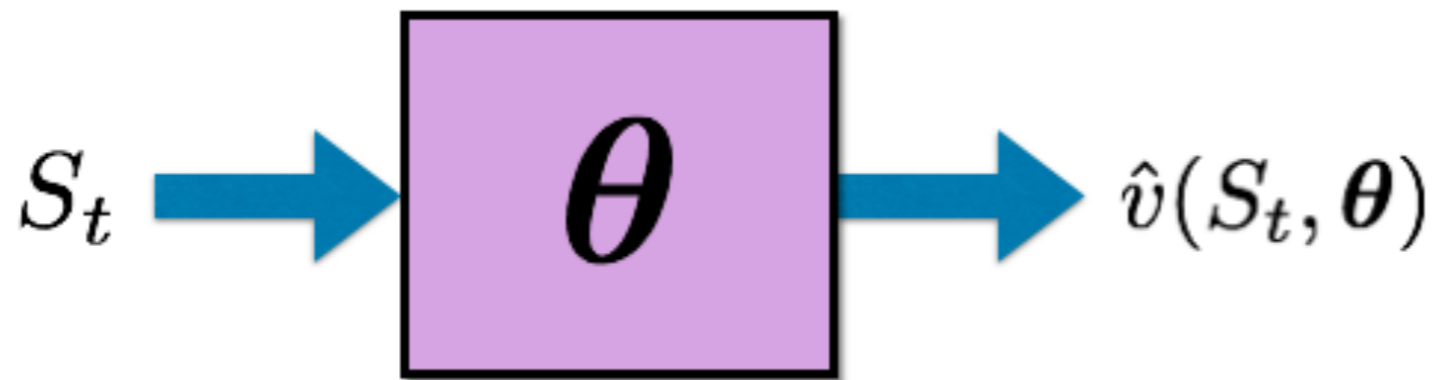
Value Function Approximation (VFA)

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

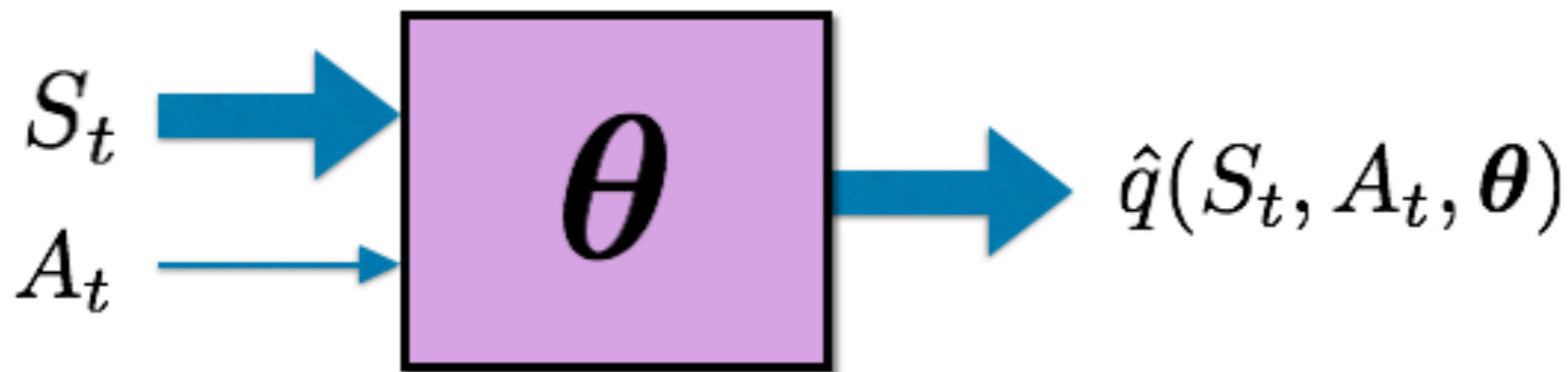


Value Function Approximation (VFA)

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



$$|\theta| \ll \ll |\mathcal{S}|$$



When we update the parameters θ , the values of many states change simultaneously!

Which Function Approximation?

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

Which Function Approximation?

▶ There are many **function approximators**, e.g.

- **Linear combinations of features**

- **Neural networks**

- Decision tree

- Nearest neighbour

- Fourier / wavelet bases

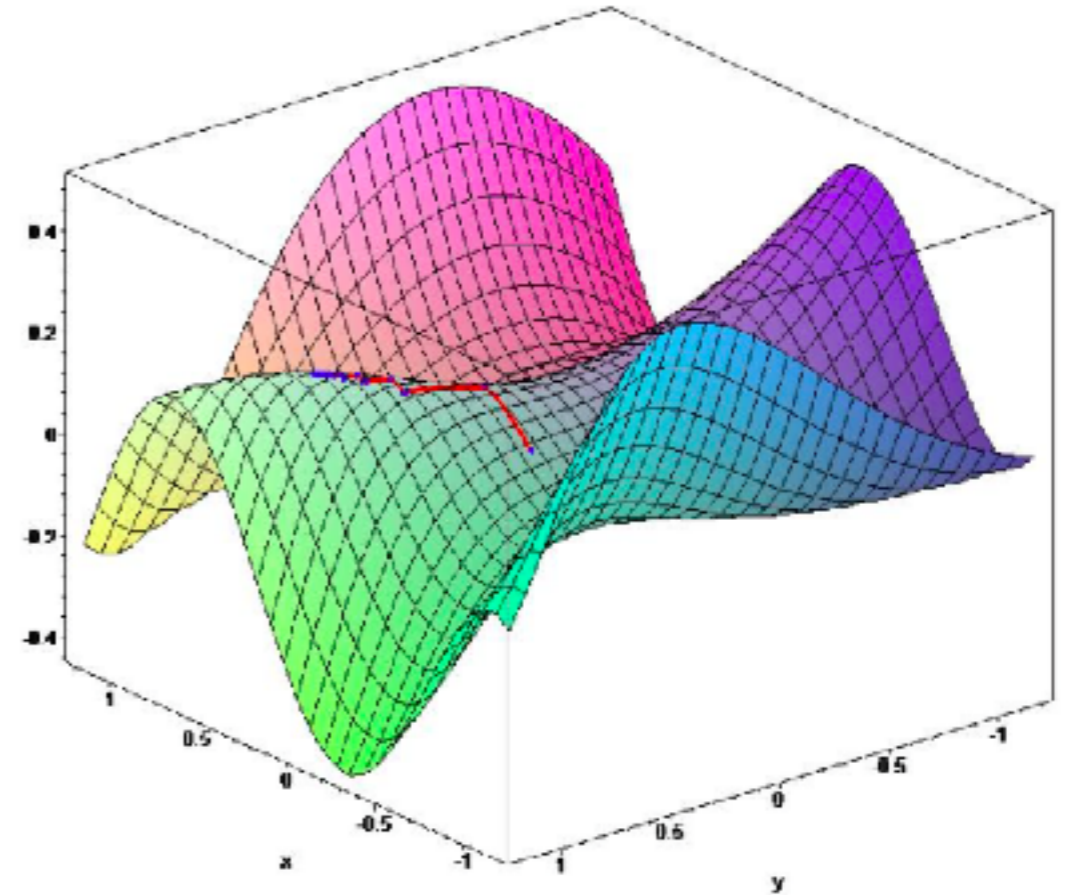
- ...

▶ **differentiable function approximators**

Gradient Descent

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

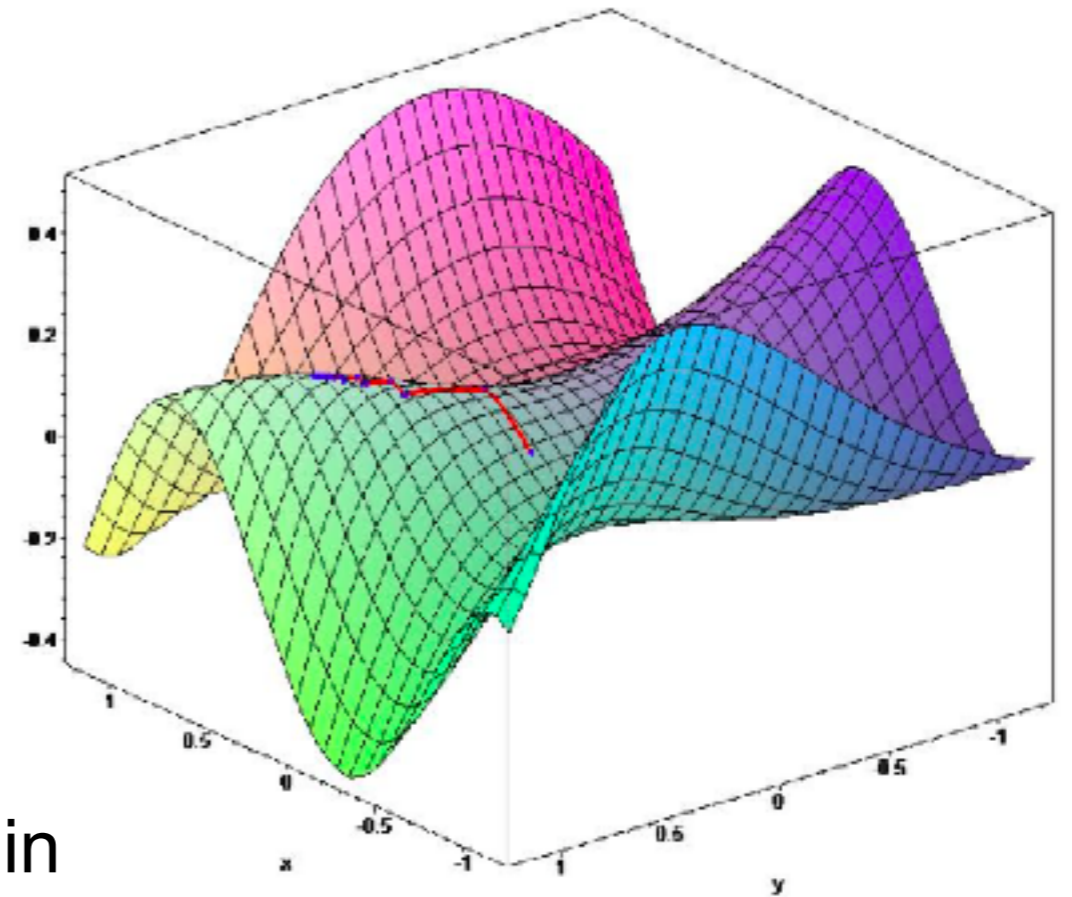
$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_n} \end{pmatrix}$$



Gradient Descent

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

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_n} \end{pmatrix}$$



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

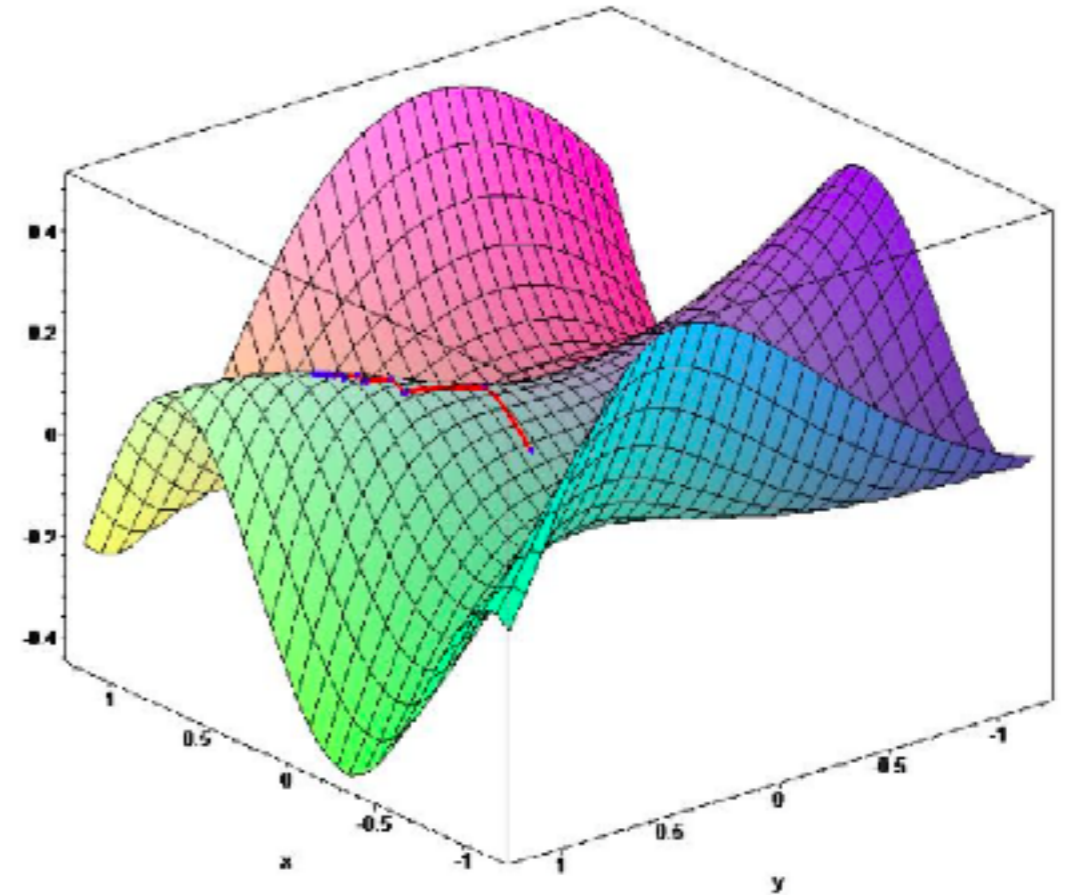
$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

Step-size

Gradient Descent

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

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_n} \end{pmatrix}$$



- ▶ Starting from a guess \mathbf{w}_0
- ▶ We consider the sequence $\mathbf{w}_0, \mathbf{w}_1, \mathbf{w}_2, \dots$ s.t. :

$$\mathbf{w}_{n+1} = \mathbf{w}_n - \frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w}_n)$$

- ▶ We then have $J(\mathbf{w}_0) \geq J(\mathbf{w}_1) \geq J(\mathbf{w}_2) \geq \dots$

Our objective

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

Our objective

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

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

Our objective

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

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

Let $\mu(S)$ denote how much time we spend in each state s under policy π , then:

$$J(w) = \sum_{n=1}^{|\mathcal{S}|} \mu(S) [v_\pi(S) - \hat{v}(S, \mathbf{w})]^2 \quad \sum_{s \in \mathcal{S}} \mu(S) = 1$$

Very important choice: it is OK if we cannot learn the value of states we visit very few times, there are too many states, I should focus on the ones that matter: the RL way of approximating the Bellman equations!

Our objective

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

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

Let $\mu(S)$ denote how much time we spend in each state s under policy π , then:

$$J(w) = \sum_{n=1}^{|\mathcal{S}|} \mu(S) [v_\pi(S) - \hat{v}(S, w)]^2 \quad \sum_{s \in \mathcal{S}} \mu(S) = 1$$

In contrast to:

$$J_2(w) = \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} [v_\pi(S) - \hat{v}(S, w)]^2$$

On-policy state distribution

Let $h(s)$ be the initial state distribution, i.e, the probability that an episode starts at state s , then:

$$\eta(s) = h(s) + \sum_{\bar{s}} \eta(\bar{s}) \sum_a \pi(a | \bar{s}) p(s | \bar{s}, a), \quad \forall s \in \mathcal{S}$$

$$\mu(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}, \quad \forall s \in \mathcal{S}$$

Gradient Descent

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

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

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

Gradient Descent

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

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

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

- ▶ Starting from a guess \mathbf{w}_0

Gradient Descent

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

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

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

- ▶ Starting from a guess w_0
- ▶ We consider the sequence w_0, w_1, w_2, \dots s.t. :

$$w_{n+1} = w_n - \frac{1}{2}\alpha \nabla_w J(w_n)$$

- ▶ We then have $J(w_0) \geq J(w_1) \geq J(w_2) \geq \dots$

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

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

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

- ▶ Gradient descent finds a local minimum:

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

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

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

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 w that give the best fitting value function $v(s,w)$?

▶ Least squares algorithms find parameter vector w minimizing sum-squared error between $v(S_t, w)$ and target values v_t^{π} :

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

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

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

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

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

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

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**

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

- ▶ 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 π 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 π

 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 $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$ 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):

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

We ignore the dependence of the target on w !

We call it semi-gradient methods

TD Learning with VFA

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

Input: the policy π 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 ϵ -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

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

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

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

Can we guess the deep Q learning update rule?

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

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., ε -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

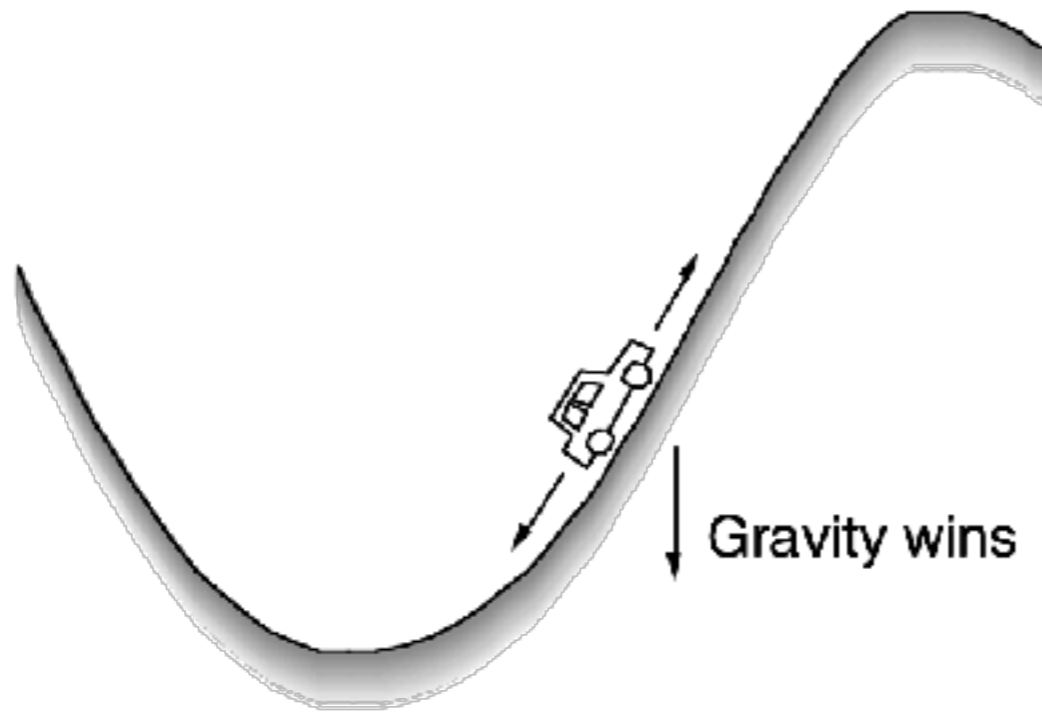
Choose A' as a function of $\hat{q}(S', \cdot, \boldsymbol{\theta})$ (e.g., ε -greedy)

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

$S \leftarrow S'$

$A \leftarrow A'$

Example: The Mountain-Car problem



Minimum-Time-to-Goal Problem

SITUATIONS:

car's position and velocity

ACTIONS:

three thrusts: forward, reverse, none

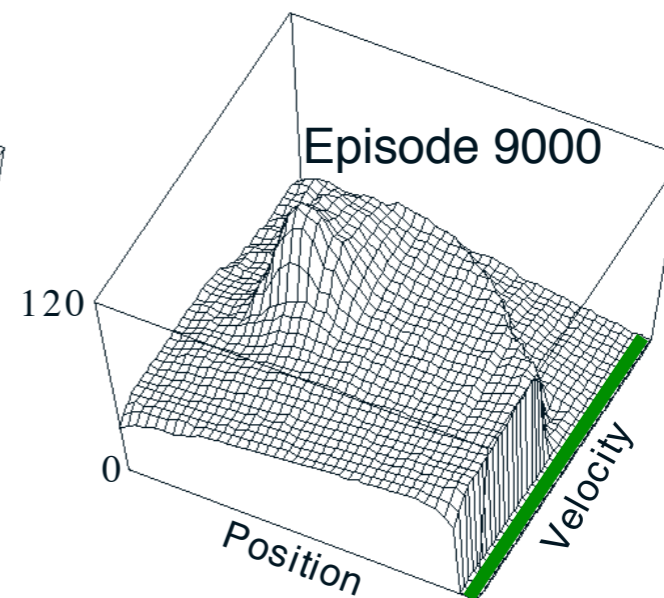
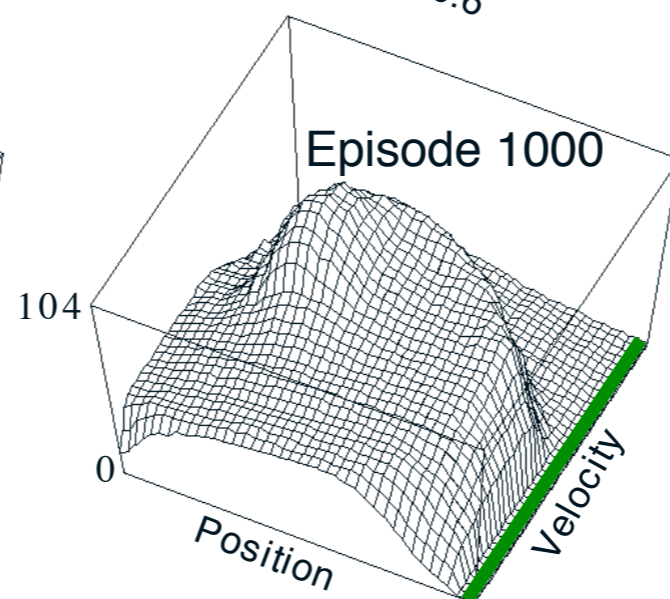
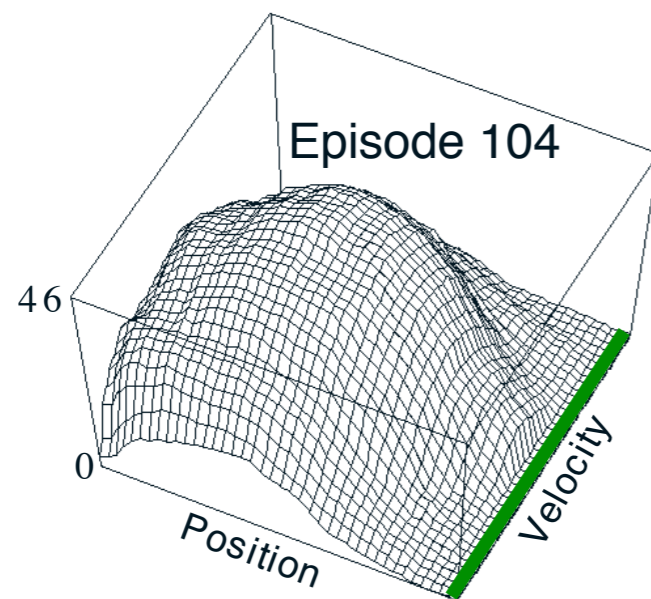
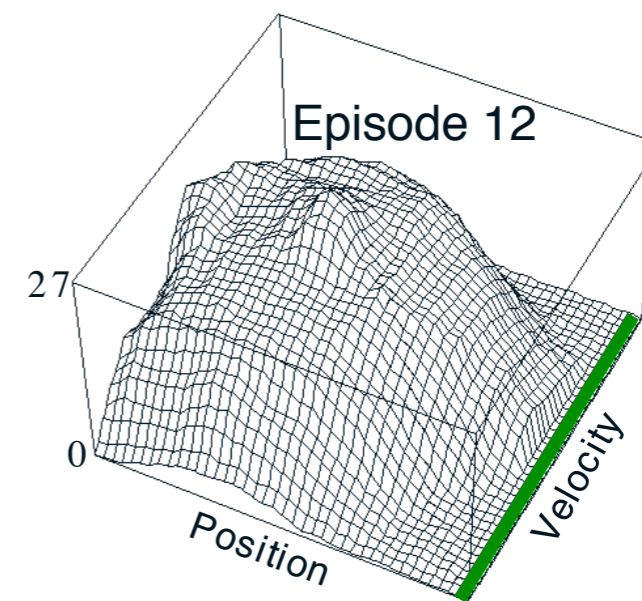
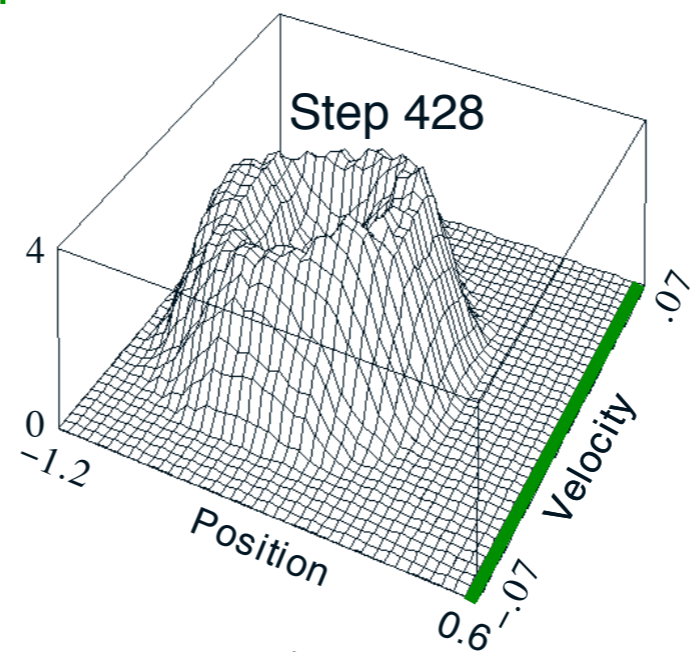
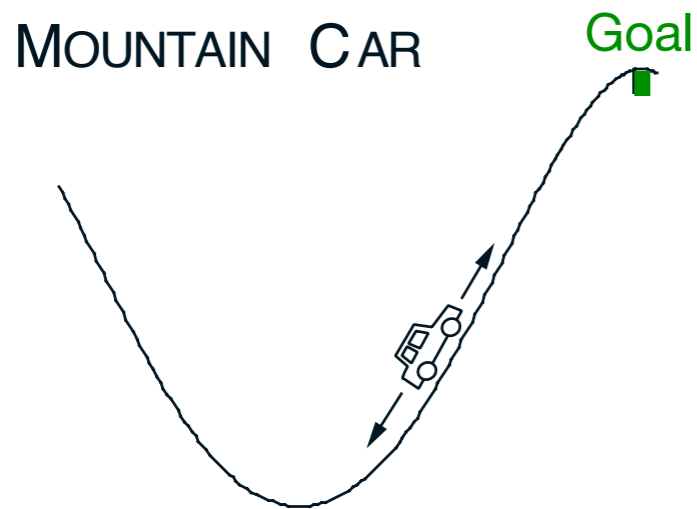
REWARDS:

always -1 until car reaches the goal

Episodic, No Discounting, $\gamma=1$

Example: The Mountain-Car problem

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



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”)

Which Function Approximation?

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

Nearest neighbors

- ▶ Save training examples in memory as they arrive $(s, v(s))$. (state, value)
- ▶ Then, given a new state s' , retrieve closest state examples from the memory and average their values based on similarity:

$$v(s') = \sum_{i=1}^K k(h_{s'}, h_{s_i}) v(s_i)$$

- ▶ Accuracy improves as more data accumulates.
- ▶ Agent's experience has an **immediate affect** on value estimates in the neighborhood of its environment's current state.
- ▶ Parametric methods need to incrementally adjust parameters of a global approximation.