

# Function Approximation

Ayça Özçelikkale

Dept. of Electrical Engineering, Uppsala University

Course Information

Motivation

Linear Value Function Approximation

How to Find the Best Approximation?

Prediction with Approximation

Control with Approximation

Concluding Remarks

## Course Information

Motivation

Linear Value Function Approximation

How to Find the Best Approximation?

Prediction with Approximation

Control with Approximation

Concluding Remarks

# New Teacher

Ayça Özçelikkale ( How to pronounce my name? Read “I-cha” in English).

- ▶ Senior Lecturer at the Division of Signals and Systems, Dept. of Electrical Engineering
- ▶ **Key research themes:**
  - ▶ statistical estimation, online estimation and optimization, machine learning including reinforcement learning
  - ▶ applications of signal processing and machine learning in communications, sensor networks and production logistics
- ▶ **Professional activities:**
  - ▶ Technical committee member of various academic conferences
  - ▶ Past editor of IEEE Wireless Communications Letters

## Course Activities with Ayça

I will hold the lectures for the following subjects:

- ▶ Function Approximation
- ▶ Model-based RL
- ▶ Policy-gradient Approaches

These subjects will be covered in the lectures L6-L9.

The corresponding tinkering notebooks are Notebook 4-6.

Basic Assignment 3, Basic Assignment 4 and a part of the Advanced Assignment will primarily focus on these subjects.

Course Information

**Motivation**

Linear Value Function Approximation

How to Find the Best Approximation?

Prediction with Approximation

Control with Approximation

Concluding Remarks

# Motivation: Why do we study function approximation?

We would like to be able to apply RL to the following problems:

- ▶ Problems with large state space
  - ▶ Games such as Go ( $10^{170}$  states)
- ▶ Problems with continuous state-space
  - ▶ Robotic control tasks
  - ▶ Many real-life decision problems

Why can't we use the methods in the previous lectures?

- ▶ Previously: Every state  $s$  (or state-action) pair has an individual entry  $v(s)$  or  $q(s, a)$
- ▶ It is difficult to learn values of states and/or actions individually
- ▶ No generalization: Learning the values individually ignores the fact that observed states can provide information about unseen states
  - ▶ Consider how the previous methods will learn the value functions for a very large maze!

# Function Approximation

**IDEA:** We will approximate  $v(s)$  or  $q(s, a)$  using **function approximation**. The number of free parameters in the approximation will be typically smaller than the number of states/state-action pairs.

There exist various function approximation methods:

- ▶ Linear Approximation
- ▶ Neural Networks
- ▶ Decision Trees
- ▶ Nearest Neighbor

In this course, we primarily focus on **differentiable** function approximators. In particular, we work on the following:

- ▶ Linear Approximation (Today's class)
- ▶ Neural Networks (Project for 7.5hp course)



Course Information

Motivation

**Linear Value Function Approximation**

How to Find the Best Approximation?

Prediction with Approximation

Control with Approximation

Concluding Remarks

# Feature Vectors

**Idea:** Associate the state  $s$  with a **feature vector**  $\mathbf{x}(s) \in \mathbb{R}^d$ :

$$\mathbf{x}(s) = \begin{bmatrix} x_1(s) \\ x_2(s) \\ \vdots \\ x_d(s) \end{bmatrix}$$

- ▶ Each  $x_i(\cdot)$  is a function from the state space to  $\mathbb{R}$ .
- ▶ The entire function  $x_i(\cdot)$  is called a **feature**.
- ▶ Features are also referred as **basis functions** for linear approximation.

**Example:** Edges in an Image (This is an example of contextually easily interpretable features)

**Example:** Polynomials

**Example:** Fourier Basis

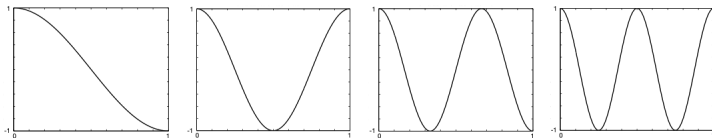
## Example: Fourier Cosine Basis

**Why Fourier Basis?:** Fourier basis is easy to use and has shown to perform well in a range of RL problems.

Let  $s \in [0, 1]$ . We consider one-dimensional order- $n$  Fourier cosine basis with  $n + 1$  features:  $x_i(s) = \cos(i\pi s)$ ,  $s \in [0, 1]$  for  $i = 0, \dots, n$ .

Let  $n = 4$ . Feature vector can be written as follows:

$$\begin{bmatrix} x_0(s) \\ x_1(s) \\ \vdots \\ x_4(s) \end{bmatrix} = \begin{bmatrix} \cos(0\pi s) \\ \cos(1\pi s) \\ \vdots \\ \cos(4\pi s) \end{bmatrix}$$



$x_i$ ,  $i = 1, \dots, 4$  for approximating functions over  $s \in [0, 1]$

**Self-study Exercise:** Plot  $x_0$ .

# Linear Value Function Approximation

We represent the value function by a linear combination of the features

$$\hat{v}(s, \mathbf{w}) = \sum_{i=1}^d w_i x_i(s) = \mathbf{w}^T \mathbf{x}(s) = \mathbf{x}(s)^T \mathbf{w}$$

where

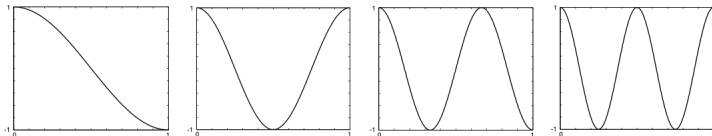
$$\mathbf{x}(s) = \begin{bmatrix} x_1(s) \\ x_2(s) \\ \vdots \\ x_d(s) \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

The approximation is **linear in the weights**.

**Example:** How do we form the approximation with order  $n$  Fourier cosine features?

**Self-study Exercise:** Suppose that you know  $v(s)$  for  $s_1, \dots, s_m$ . Suggest a method to find the weights  $\mathbf{w}$  so that  $v(s) \approx \hat{v}(s, \mathbf{w})$ .

$$\hat{v}(s, \mathbf{w}) = \sum_{i=0}^{n+1} w_i x_i(s) = \mathbf{w}^T \mathbf{x}(s) = \mathbf{x}(s)^T \mathbf{w}$$



$x_i$ ,  $i = 1, \dots, 4$  for approximating functions over  $s \in [0, 1]$

Note: If you think that you need to know that how much time we spend at each  $s_i$ , assume that time spent for all states are equal.

Course Information

Motivation

Linear Value Function Approximation

**How to Find the Best Approximation?**

Preliminaries: Gradient Descent

Value Function Approximation with Gradient Descent

Prediction with Approximation

Control with Approximation

Concluding Remarks

## Course Information

## Motivation

## Linear Value Function Approximation

### How to Find the Best Approximation?

#### Preliminaries: Gradient Descent

Value Function Approximation with Gradient Descent

### Prediction with Approximation

Monte-Carlo for Estimating  $v_{\pi}(s)$

TD(0) for Estimating  $v_{\pi}(s)$

### Control with Approximation

Action-Value Function Approximation

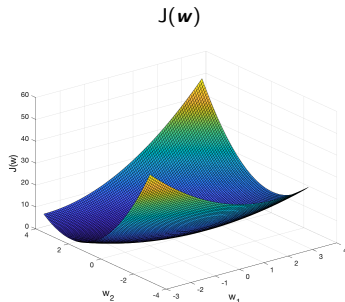
Example: Mountain Car

## Concluding Remarks

# Gradient Descent: A general idea for minimizing functions

Why do we study gradient descent?

- ▶ To motivate **stochastic** gradient descent!
- ▶ Using **stochastic gradient descent (SGD)**, we will motivate the RL methods for approximating  $v$  and  $q$ .



**What is gradient descent?** Gradient descent is an iterative procedure to find minimum of functions.

**Problem statement:** Let  $J(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$  be a scalar valued function with a vector input. We would like to find a  $\mathbf{w}$  that minimizes  $J(\mathbf{w})$  iteratively.



## Notation- Gradient

Let  $J(\mathbf{w}) : \mathbb{R}^d \rightarrow \mathbb{R}$  be a scalar valued function with a vector input. Gradient of  $J(\mathbf{w})$  is given by

$$\nabla J \triangleq \frac{\partial J}{\partial \mathbf{w}} \triangleq \begin{bmatrix} \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} \\ \vdots \\ \frac{\partial J}{\partial w_d} \end{bmatrix} \in \mathbb{R}^{d \times 1}$$

The gradient evaluated at a particular value  $\bar{\mathbf{w}}$  is denoted by the following:

$$\nabla J(\bar{\mathbf{w}}) = \left. \frac{\partial J}{\partial \mathbf{w}} \right|_{\mathbf{w}=\bar{\mathbf{w}}}$$

The gradient  $\nabla J(\bar{\mathbf{w}})$  gives the **direction of maximum ascent** around  $\bar{\mathbf{w}}$ .

**Self-study Exercise:** Let  $J(\mathbf{w}) = w_1^2 + 3w_2$ . Find  $\nabla J(\bar{\mathbf{w}})$  for  $\bar{\mathbf{w}} = \begin{bmatrix} 5 \\ 2 \end{bmatrix}$ .

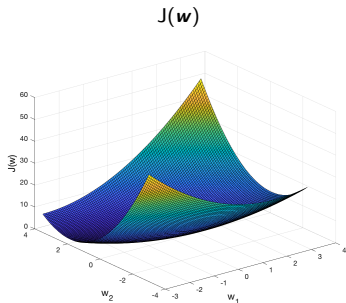
# Gradient Descent Idea

**Gradient Descent Idea:** If you want to minimize a function  $J(\mathbf{w})$  iteratively, start from an initial point  $\mathbf{w}_i$  ( $i$  is the iteration index!), and move in the direction of maximum “descent”:

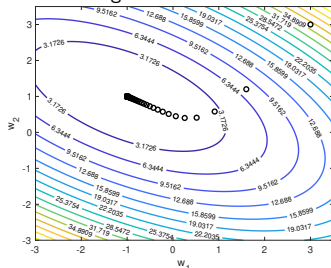
$$\begin{aligned}\mathbf{w}_{i+1} &= \mathbf{w}_i + \Delta \mathbf{w}_i \\ &= \mathbf{w}_i - \frac{1}{2} \alpha \nabla J(\mathbf{w}_i)\end{aligned}$$

where  $\alpha > 0$  is the step size. (Here,  $\frac{1}{2}$  is just to make our notation neater since it will cancel out later.)

**Example:** Find  $\mathbf{w}$  that minimizes  $J(\mathbf{w})$  plotted below.



Contour Plot of  $J(\mathbf{w})$  with example trajectory of gradient descent



Course Information

Motivation

Linear Value Function Approximation

**How to Find the Best Approximation?**

Preliminaries: Gradient Descent

**Value Function Approximation with Gradient Descent**

Prediction with Approximation

Monte-Carlo for Estimating  $v_{\pi}(s)$

TD(0) for Estimating  $v_{\pi}(s)$

Control with Approximation

Action-Value Function Approximation

Example: Mountain Car

Concluding Remarks

# Value Function Approximation with Gradient Descent

- ▶ Let  $J(\mathbf{w})$  denote the mean-square error between the true value function  $v_\pi(s)$  and our approximation  $\hat{v}(s, \mathbf{w})$ :

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

- ▶ We would like to find the optimal  $\mathbf{w}$  that minimizes  $J(\mathbf{w})$  iteratively using gradient descent:

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

- ▶ It is difficult to know  $\mathbb{E}_\pi[v_\pi(S) - \hat{v}(S, \mathbf{w})]$ !
- ▶ **Stochastic Gradient Descent Idea:** Use a sample instead of the expectation

$$\mathbf{w}_{i+1} = \mathbf{w}_i + \alpha(v_\pi(S) - \hat{v}(S, \mathbf{w}_i)) \nabla \hat{v}(S, \mathbf{w}_i)$$

- ▶ Hence, the update  $\Delta \mathbf{w}_i$  is given by:  $\Delta \mathbf{w}_i = \alpha(v_\pi(S) - \hat{v}(S, \mathbf{w}_i)) \nabla \hat{v}(S, \mathbf{w}_i)$

$$\Delta \mathbf{w}_i = \text{step size} \times \text{prediction error} \times \text{gradient of value function approximation}$$

## Linear Value Function Approximation with Gradient Descent

- ▶ Let  $J(\mathbf{w})$  denote the mean-square error between the true value function  $v_\pi(s)$  and our approximation  $\hat{v}(s, \mathbf{w})$ :

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

- ▶ We would like to find the optimal  $\mathbf{w}$  that minimizes  $J(\mathbf{w})$  iteratively using gradient descent:

$$\begin{aligned}\mathbf{w}_{i+1} &= \mathbf{w}_i - \frac{1}{2} \alpha \nabla J(\mathbf{w}_i) \\ &= \mathbf{w}_i + \alpha \mathbb{E}_\pi[v_\pi(s) - \mathbf{x}(S)^T \mathbf{w}_i] \mathbf{x}(S)\end{aligned}$$

- ▶ It is difficult to know  $\mathbb{E}_\pi[v_\pi(S) - \mathbf{x}(S)^T \mathbf{w}_i]!$
- ▶ **Stochastic Gradient Descent Idea:** Use a sample instead of the expectation

$$\mathbf{w}_{i+1} = \mathbf{w}_i + \alpha((v_\pi(S) - \hat{v}(S, \mathbf{w}_i)) \mathbf{x}(S))$$

- ▶ Hence, the update  $\Delta \mathbf{w}_i$  is given by  $\Delta \mathbf{w}_i = \alpha(v_\pi(S) - \hat{v}(S, \mathbf{w}_i)) \mathbf{x}(S)$

$$\Delta \mathbf{w}_i = \text{step size} \times \text{prediction error} \times \text{gradient of value function approximation}$$

Course Information

Motivation

Linear Value Function Approximation

How to Find the Best Approximation?

**Prediction with Approximation**

Monte-Carlo for Estimating  $v_{\pi}(s)$

TD(0) for Estimating  $v_{\pi}(s)$

Control with Approximation

Concluding Remarks

## Linear Value Function Approximation with Gradient Descent Without Knowing $v_\pi(S)$

- ▶ We obtained

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

i.e., we assumed that true value function is known, but we do not know  $v_\pi(s)$ !

- ▶ We substitute a target for  $v_\pi(s)$ :

- ▶ For Monte-Carlo, the target is the return  $G_t$

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

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

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

## Course Information

## Motivation

## Linear Value Function Approximation

### How to Find the Best Approximation?

- Preliminaries: Gradient Descent

- Value Function Approximation with Gradient Descent

### Prediction with Approximation

- Monte-Carlo for Estimating  $v_{\pi}(s)$

- TD(0) for Estimating  $v_{\pi}(s)$

### Control with Approximation

- Action-Value Function Approximation

- Example: Mountain Car

## Concluding Remarks



# Monte-Carlo with Linear Value Function Approximation

- ▶ Return  $G_t$  is an unbiased, noisy estimate of the true value  $v_\pi(s)$
- ▶ Below Monte-Carlo algorithm converges to a global optimum (if  $\alpha$  is reduced over time appropriately).

Input: Inputs: Policy  $\pi$ , features  $\mathbf{x}(\cdot)$ .

Initialization: Set step size  $\alpha > 0$ . Initialize weights  $\mathbf{w} \in \mathbb{R}^d$ .

repeat (for each episode)

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

repeat : (for each step of the episode  $t = 0, 1, \dots, T-1$ )

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha(G_t - \hat{v}(S_t, \mathbf{w}))\mathbf{x}(S_t)$$

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

## Example: Function Approximation for 1000-state random walk task

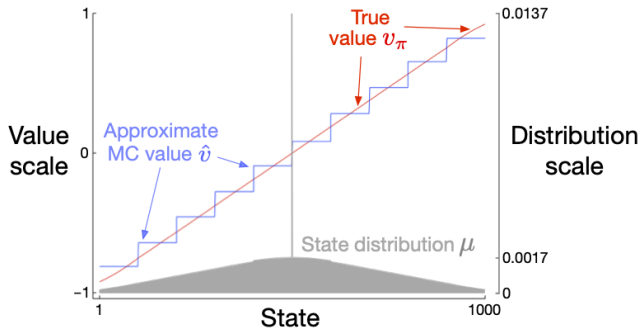
### Task:

- ▶ A corridor of non-terminal 1000-states and one terminal state at each end
- ▶ State transitions are to one of the 100 neighboring states to its left, or to one of the 100 neighboring states to its right, all with equal probability.
  - ▶ Around the edges, all the probability due to fewer neighbours go the terminating on the edge. Example: State 950 has a 0.25 chance of terminating on the right terminating state within one step
- ▶ All episodes start at state 500
- ▶ All transitions have the reward of 0, except to the left terminating state with -1 and to the right terminating state with +1.

**Function Approximation by state-aggregation:** States are grouped together, with one approximated value (one component of the weight vector  $\mathbf{w}$ ) for each group. Here, the 1000 states were partitioned into 10 groups of 100 states each.

**Self-study Exercise:** What are the basis functions for the above state-aggregation setting?

## Example: Function Approximation for 1000-state random walk task



Value approximation and state distribution

State distribution  $\mu(s)$ : fraction of time spent on  $s$  under policy  $\pi$ . By definition,  $\mu(s) \geq 0$  and  $\sum_{s \in \mathcal{S}} \mu(s) = 1$ . We have  $\sum_{s \in \mathcal{S}} \mu(s)(v_\pi(s) - \hat{v}(s, \mathbf{w}))^2 = \mathbb{E}_\pi[(v_\pi(S) - \hat{v}(S, \mathbf{w}))^2]$ .

## Course Information

## Motivation

## Linear Value Function Approximation

### How to Find the Best Approximation?

Preliminaries: Gradient Descent

Value Function Approximation with Gradient Descent

### Prediction with Approximation

Monte-Carlo for Estimating  $v_{\pi}(s)$

**TD(0) for Estimating  $v_{\pi}(s)$**

### Control with Approximation

Action-Value Function Approximation

Example: Mountain Car

## Concluding Remarks

# TD(0) with Linear Value Function Approximation

- ▶ The TD-target  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$  is a biased estimate of the true value  $v_\pi(s)$
- ▶ In contrast to Monte-Carlo methods, we do not need to wait for the end of episode.
- ▶ Below TD(0) algorithm converges to a local optimum.

Input: Inputs: Policy  $\pi$ , features  $\mathbf{x}(\cdot)$ .

Initialization: Set step size  $\alpha > 0$ . Initialize weights  $\mathbf{w} \in \mathbb{R}^d$ .

repeat (for each episode)

    Initialize  $S$ .

repeat : (for each step of the episode)

        Choose  $A \sim \pi(\cdot|S)$

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

$\mathbf{w} \leftarrow \mathbf{w} + \alpha(R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w}))\mathbf{x}(S)$

$S \leftarrow S'$

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

Course Information

Motivation

Linear Value Function Approximation

How to Find the Best Approximation?

Prediction with Approximation

**Control with Approximation**

Action-Value Function Approximation

Example: Mountain Car

Concluding Remarks

## Course Information

## Motivation

## Linear Value Function Approximation

### How to Find the Best Approximation?

- Preliminaries: Gradient Descent

- Value Function Approximation with Gradient Descent

### Prediction with Approximation

- Monte-Carlo for Estimating  $v_{\pi}(s)$

- TD(0) for Estimating  $v_{\pi}(s)$

### Control with Approximation

- Action-Value Function Approximation

- Example: Mountain Car

## Concluding Remarks

# Linear Action-Value Function Approximation

**Feature Vectors:** Associate the state-action pair  $(s, a)$  with a feature vector  $\mathbf{x}(s, a) \in \mathbb{R}^d$ :

$$\mathbf{x}(s, a) = \begin{bmatrix} x_1(s, a) \\ x_2(s, a) \\ \vdots \\ x_d(s, a) \end{bmatrix}$$

We represent the action-value function by a linear combination of the features

$$\hat{q}(s, a, \mathbf{w}) = \sum_{i=1}^d w_i x_i(s, a) = \mathbf{w}^T \mathbf{x}(s, a) = \mathbf{x}(s, a)^T \mathbf{w}$$



# Action-Value Function Approximation

- ▶ **Aim:** Approximate the action value function  $\hat{q}(S, A, \mathbf{w}) \approx q_\pi(S, A)$
- ▶ We're now interested in minimizing  $J(\mathbf{w}) = \mathbb{E}_\pi[(q_\pi(S, A) - \hat{q}(S, A, \mathbf{w}))^2]$
- ▶ The update: **step size**  $\times$  **prediction error**  $\times$  **gradient of action-value function approx.**

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

- ▶ We substitute a target for  $q_\pi(S, A)$ . For Sarsa, we have

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

- ▶ Actions from continuous spaces or from large discrete sets are more difficult to handle than such state spaces.

## Course Information

## Motivation

## Linear Value Function Approximation

### How to Find the Best Approximation?

Preliminaries: Gradient Descent

Value Function Approximation with Gradient Descent

### Prediction with Approximation

Monte-Carlo for Estimating  $v_{\pi}(s)$

TD(0) for Estimating  $v_{\pi}(s)$

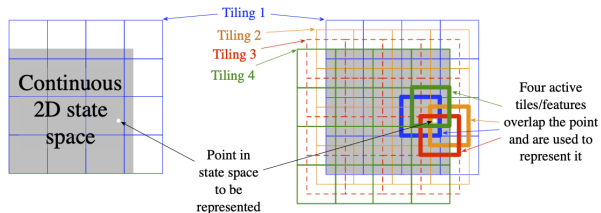
### Control with Approximation

Action-Value Function Approximation

Example: Mountain Car

## Concluding Remarks

# Tile Coding



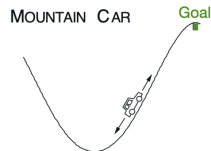
Multiple, over-lapping grid-tilings, each off-set by a fraction of a tile-width on a bounded two-dimensional space.

- ▶ The number of features that are active at one time is the same for any state. Setting the step size  $\alpha$  is easier.
- ▶ Since the features are binary vectors, calculation of approximate value function is easy.

# Mountain Car

## Mountain Car Task:

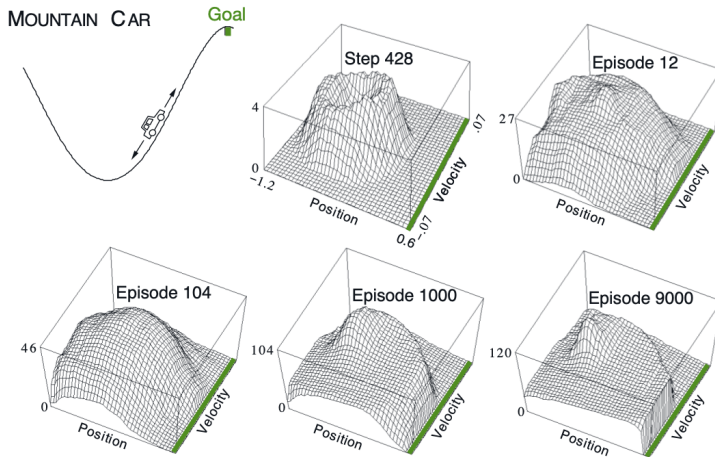
- ▶ States: Position  $x_t \in [-1.2, 0.6]$ , velocity  $\dot{x}_t \in [-0.07, 0.07]$
- ▶ Actions: + full forward (+1), -full reverse (-1), zero throttle (0).
- ▶ Reward: -1 for all transitions.
- ▶ Car starts randomly  $x_t \in [-0.6, -0.4]$ .
- ▶ Episode terminates when the car passes the goal.
- ▶ Caveat: Car cannot climb the hill the goal resides without backing on the hill on the left.



## Approach: Sarsa with linear approximation with tile coding

- ▶ Initialization: Initial state-action values are set to an optimistic 0 (all true values are negative in this task).
- ▶  $\epsilon$ -greedy with  $\epsilon = 0$

# Mountain Car with Sarsa with linear approximation



Cost-to-go function ( $-\max_a \hat{q}(s, a, \mathbf{w})$ ) for Mountain Car task

Course Information

Motivation

Linear Value Function Approximation

How to Find the Best Approximation?

Prediction with Approximation

Control with Approximation

Concluding Remarks

# Concluding Remarks I

**Motivation:** With function approximation, we would like to be able to

- ▶ deal with large/continuous state/action-spaces
- ▶ use less computation/memory/experience
- ▶ obtain better generalization properties

## Linear Approximation

- ▶ Value Function:  $\hat{v}(s, \mathbf{w}) = \sum_{i=1}^d w_i x_i(s)$
- ▶ Action-Value Function:  $\hat{q}(s, a, \mathbf{w}) = \sum_{i=1}^d w_i x_i(s, a)$

## Concluding Remarks II

**How to find the “best” approximation?** We combine two ideas:

- ▶ Stochastic Gradient Descent
  - ▶ Update: step size  $\times$  prediction error  $\times$  gradient of (action-)value function approx.

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

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

- ▶ Replace  $v_\pi(s, \mathbf{w})$  or  $q(s, a, \mathbf{w})$  with an appropriate update target:
  - ▶ MC for predicting  $v_\pi(s, \mathbf{w})$ :  $G_t$
  - ▶ TD(0) for predicting  $v_\pi(s, \mathbf{w})$ :  $R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})$
  - ▶ Sarsa for predicting  $q_\pi(s, a, \mathbf{w})$ :  $R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w})$

**More?**

- ▶ Important Non-linear Approximation Approach: Artificial Neural Networks.



# References

- ▶ Sutton, Barto, Reinforcement Learning: An Introduction, 2nd edition
- ▶ Bertsekas, Reinforcement Learning and Optimal Control, 1st edition
- ▶ Tsitsiklis and Van Roy. An Analysis of Temporal-Difference Learning with Function Approximation. 1997