

Value Function Approximation

岩延 yany@ucas.ac.cn

Value Function Approximation

Introduction

▶ Large-Scale Reinforcement Learning

- Reinforcement learning can be used to solve large problems, e.g.
 - Computer Go: 10^{170} states
 - Helicopter: continuous state space

▶ 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

▶ How can we scale up the RL methods for prediction and control?

Value Function Approximation

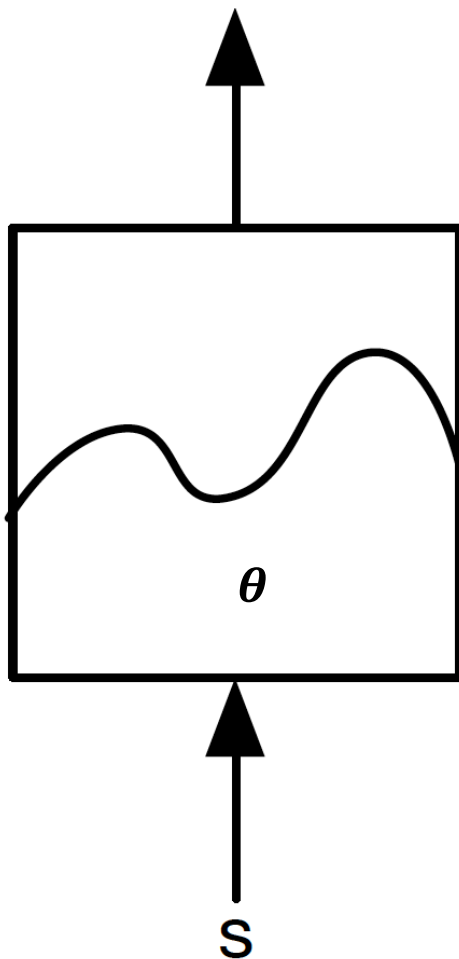
Value Function Approximation

- ▶ **Policy Evaluation (the prediction problem)**
 - for a given policy π , compute the state-value function V_π
- ▶ **So far value functions were stored in lookup tables**
 - Every state s has an entry $V(s)$
 - Or every state-action pair s, a has an entry $Q(s, a)$
- ▶ **Value Function Approximation**
 - Solution for large MDPs
 - the value function estimate at time t , V_t , depends on a parameter vector θ_t
 - Estimate value function with function approximation
 - $V(s, \theta) \approx V_\pi(s)$
 - $Q(s, a, \theta) \approx Q_\pi(s, a)$
 - Generalize from seen states to unseen states
 - Only the parameter vector is updated

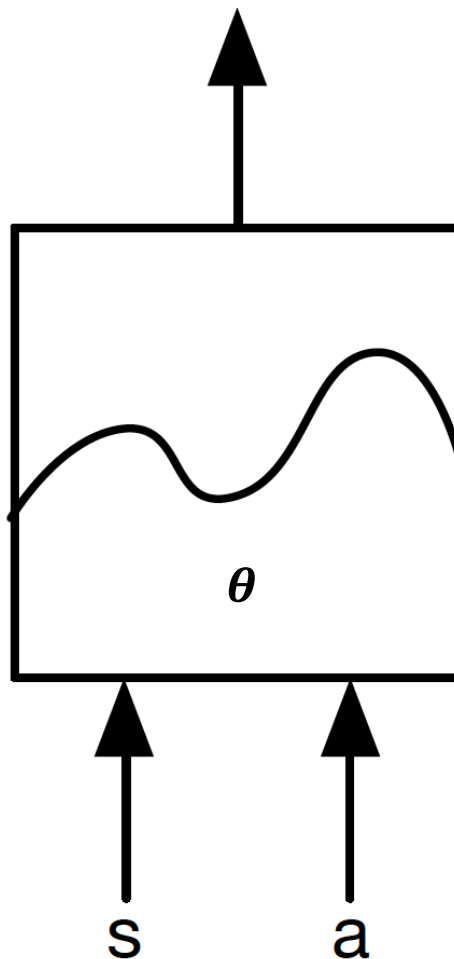
Value Function Approximation

Types of Value Function Approximation

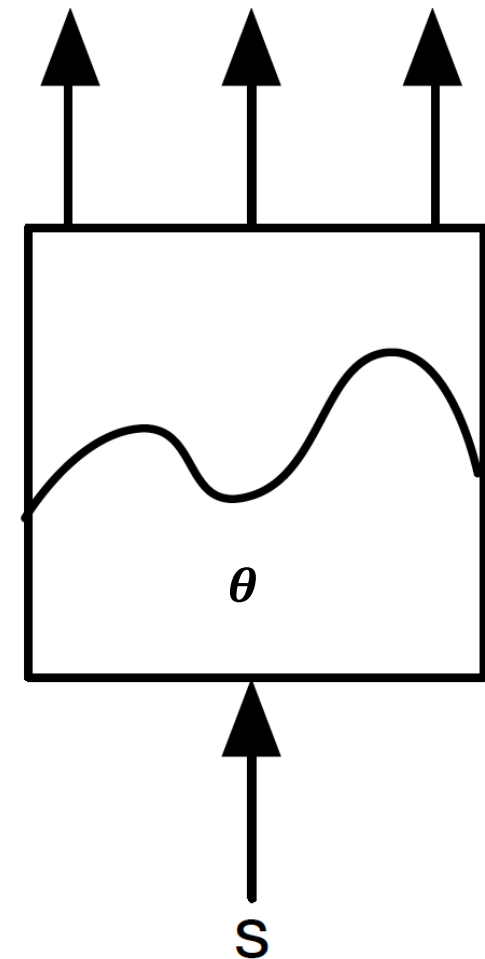
► $V(s, \theta)$



$Q(s, a, \theta)$



$Q(s, a_1, \theta) Q(s, a_2, \theta)$



Value Function Approximation

Backups as Training Examples

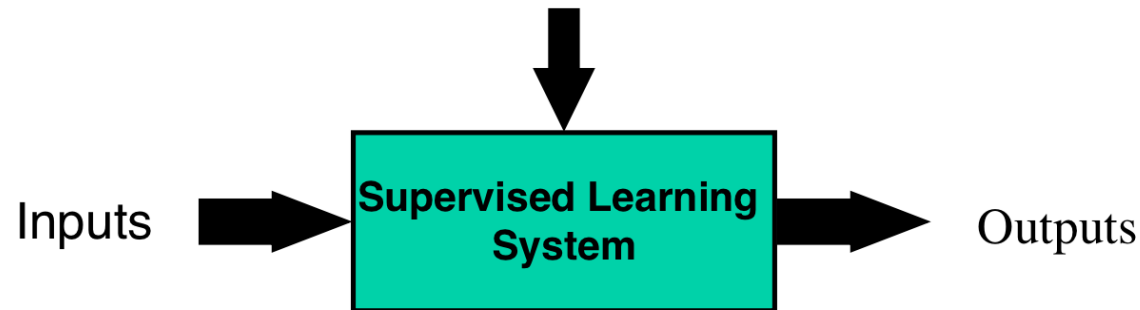
► Backup $s \rightarrow u$

- estimated value function update its value at particular states toward a backed-up value
- TD(0):
 - $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
 - $S_t \rightarrow R_{t+1} + \gamma V(S_{t+1})$
- MC: $S_t \rightarrow G_t$
- n-step TD: $S_t \rightarrow G_{t:t+n}$
- DP: $S_t \rightarrow E[R_{t+1} + \gamma V(S_{t+1}) | S_t = s]$

Value Function Approximation

Adapt Supervised Learning Algorithms

- ▶ Training Info = desired (target) outputs



- ▶ Training example = {input, target output}
- ▶ Error = (target output – actual output)
- ▶ e.g., TD(0) backup:
 - $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$
 - $S_t \rightarrow R_{t+1} + \gamma V(S_{t+1})$
- ▶ As a training example
 - $\{S_t, R_{t+1} + \gamma V(S_{t+1})\}$

Value Function Approximation

Which Function Approximator?

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

- Linear combinations of features
- Neural network
- Decision tree
- Nearest neighbor
- Fourier / wavelet bases
- ...

- ▶ **But RL has some special requirements:**

- usually want to learn while interacting
- ability to handle nonstationarity
- other?

Value Function Approximation

Which Function Approximator?

- ▶ We consider **differentiable** function approximators, e.g.
 - Linear combinations of features
 - Neural network
 - Decision tree
 - Nearest neighbor
 - Fourier / wavelet bases
 - ...
- ▶ Furthermore, we require a training method that is suitable for non-stationary, non-iid data

Value Function Approximation

Performance Measures

► The Prediction Objective

- a common and simple one is the mean-squared error (MSE) over a state distribution μ under policy π
 - $\text{MSE}(\boldsymbol{\theta}) = \sum_{s \in \mathcal{S}} \mu(s) [V_\pi(s) - V(s, \boldsymbol{\theta})]^2$
 - $= \mathbb{E}_\pi [[V_\pi(s) - V(s, \boldsymbol{\theta})]^2]$

► The on-policy distribution

- the distribution created while following the policy being evaluated. Stronger results are available for this distribution.

► Global optimum and local optimum

- $\text{MSE}(\boldsymbol{\theta}^*) \leq \text{MSE}(\boldsymbol{\theta})$

Value Function Approximation

Gradient Descent

- ▶ Let $J(\theta)$ be a differentiable function of parameter vector θ

- $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$

- ▶ Define the gradient of $J(\theta)$

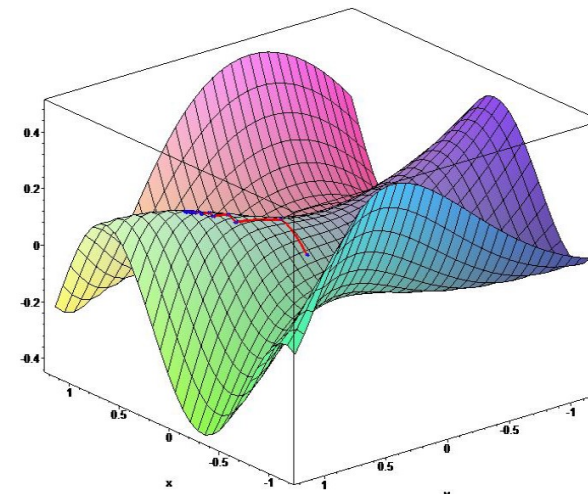
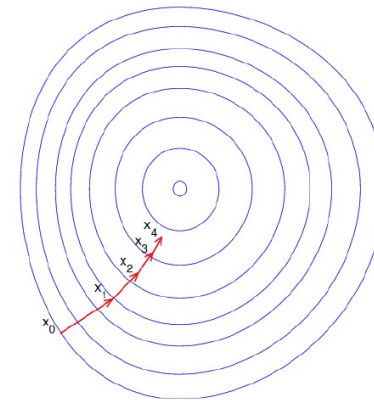
- $\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$

- ▶ To find a local minimum of $J(\theta)$

- ▶ Adjust θ in direction of gradient

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

- $\theta_{t+1} = \theta_t - \frac{1}{2}\alpha\nabla_{\theta}J(\theta)$



Value Function Approximation

Gradient Descent

- ▶ Let $J(\boldsymbol{\theta})$ be a differentiable function of parameter vector $\boldsymbol{\theta}$
 - $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_n)^T$
- ▶ Assume V_t is a (sufficiently smooth) differentiable function
- ▶ $S_t \rightarrow V_\pi(S_t)$
 - Training examples form
 - $\{S_t, V_\pi(S_t)\}$
- ▶ For the MSE object
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{1}{2}\alpha \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{1}{2}\alpha \nabla_{\boldsymbol{\theta}} \text{MSE}(\boldsymbol{\theta}_t)$
 - $= \boldsymbol{\theta}_t - \frac{1}{2}\alpha \nabla_{\boldsymbol{\theta}} \mathbb{E}_\pi [[V_\pi(s) - V(s, \boldsymbol{\theta}_t)]^2]$
 - $= \boldsymbol{\theta}_t + \alpha \mathbb{E}_\pi [V_\pi(s) - V(s, \boldsymbol{\theta}_t)] \nabla_{\boldsymbol{\theta}} V(s, \boldsymbol{\theta}_t)$
 - $= \boldsymbol{\theta}_t + \alpha \sum_{s \in S} \mu(s) [V_\pi(s) - V(s, \boldsymbol{\theta}_t)] \nabla_{\boldsymbol{\theta}} V(s, \boldsymbol{\theta}_t)$
- ▶ Gradient descent finds a local minimum

Value Function Approximation

Stochastic Gradient Descent

- ▶ **Gradient Descent**

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

- ▶ **Stochastic gradient descent samples the gradient**

- $\Delta\theta = \alpha(V_{\pi}(S_t) - V(S_t, \theta))\nabla_{\theta}V(S_t, \theta)$

- ▶ **Expected update is equal to full gradient update**

- ▶ **Each sample gradient is an unbiased estimate of the true gradient**

- $E_{\pi}[V_{\pi}(s) - V(s, \theta)] \nabla_{\theta}V(s, \theta) = \sum_{s \in S} \mu(s)[V_{\pi}(s) - V(s, \theta)] \nabla_{\theta}V(s, \theta)$

- SGD converges to a local minimum of the MSE if α decreases appropriately with t .

Value Function Approximation

General SGD method

- ▶ **Use target $s \rightarrow U_t$ instead of $V_\pi(s)$**
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(U_t - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
- ▶ **If each U_t is an unbiased estimate of $V_\pi(S_t)$**
 - $E[U_t] = V_\pi(S_t)$
 - Gradient descent converges to a local minimum (provided α decreases appropriately)
- ▶ **Example**
 - The MC target $S_t \rightarrow G_t$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - The TD(0) target $S_t \rightarrow R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - The TD(λ) target $S_t \rightarrow G_t^\lambda$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t^\lambda - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$

Value Function Approximation

Stochastic gradient descent

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

Input: the policy π to be evaluated

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

Algorithm parameter: step size $\alpha > 0$

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

Loop forever (for each episode):

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

 Loop for each step of episode, $t = 0, 1, \dots, T - 1$:

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

Value Function Approximation

Semi-gradient Methods

- ▶ **No guarantees if bootstrapping is used as the target U_t**
- ▶ **Bootstrapping U_t is not unbiased**
 - all depend on the current value of the weight vector θ
 - $\theta_{t+1} = \theta_t - \frac{1}{2}\alpha \nabla_{\theta} [U_t - V(S_t, \theta_t)]^2$
 - $= \theta_t + \alpha(U_t - V(S_t, \theta_t)) \nabla_{\theta} V(S_t, \theta_t)$
 - will not produce a true gradient-descent method
- ▶ **Include only a part of the gradient**
 - take into account the effect of changing the weight vector θ on the estimate, but ignore its effect on the target.

Value Function Approximation

Semi-gradient Methods

► semi-gradient TD(0)

$$- S_t \rightarrow R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t)$$

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

Input: the policy π to be evaluated

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

Algorithm parameter: step size $\alpha > 0$

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

Loop for each episode:

 Initialize S

 Loop for each step of episode:

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

 Take action A , observe R, S'

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

$S \leftarrow S'$

 until S is terminal

Value Function Approximation

Linear Function Approximation

► Feature Vectors

- Represent state by a feature vector

- $\mathbf{x}(s) = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$

► Linear methods

- Represent value function by a linear combination of features

- $v(s, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}(s) = \sum_{i=1}^n \theta_i x_i(s)$

Value Function Approximation

Linear Function Approximation

- ▶ **Objective function is quadratic in parameters θ**
 - $\text{MSE}(\theta) = \sum_{s \in \mathcal{S}} \mu(s) [V_\pi(s) - V(s, \theta)]^2$
 - $= \mathbb{E}_\pi [[V_\pi(s) - V(s, \theta)]^2]$
 - $= \mathbb{E}_\pi [[V_\pi(s) - \theta^T \mathbf{x}]^2]$
- ▶ **Stochastic gradient descent converges on global optimum**
- ▶ **Update rule**
 - $\nabla_\theta V(s, \theta) = \mathbf{x}(s)$
 - $\Delta \theta = \alpha (V_\pi(s) - V(s, \theta)) \nabla_\theta V(s, \theta)$
 - $= \alpha (V_\pi(s) - V(s, \theta)) \mathbf{x}(s)$
 - Update = step-size \times prediction error \times feature value

Value Function Approximation

Linear Function Approximation

► semi-gradient TD(0) also converges

- Need different theorem
- Not global optimum, near the local optimum
- The TD(0) target $S_t \rightarrow R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t)) \nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma \boldsymbol{\theta}_t^T \mathbf{x}_{t+1} - \boldsymbol{\theta}_t^T \mathbf{x}_t) \mathbf{x}_t$
 - $\quad = \boldsymbol{\theta}_t + \alpha(R_{t+1} \mathbf{x}_t - \mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^T \boldsymbol{\theta}_t)$
 - For steady state
 - $E(\boldsymbol{\theta}_{t+1} | \boldsymbol{\theta}_t) = \boldsymbol{\theta}_t + \alpha(\mathbf{b} - \mathbf{A} \boldsymbol{\theta}_t)$
 - Where
 - $\mathbf{b} = E[R_{t+1} \mathbf{x}_t]$ and $\mathbf{A} = E[\mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^T]$
- LSTD: Converge to the TD fixed point $\boldsymbol{\theta}_{TD}$
 - $\mathbf{b} - \mathbf{A} \boldsymbol{\theta}_{TD} = \mathbf{0}$
 - $\mathbf{b} = \mathbf{A} \boldsymbol{\theta}_{TD}$
 - $\boldsymbol{\theta}_{TD} = \mathbf{A}^{-1} \mathbf{b}$

Value Function Approximation

Prediction Algorithms

- ▶ Have assumed true value function $V_\pi(s)$ given by supervisor
- ▶ In RL there is no supervisor, only rewards
- ▶ Use target $s \rightarrow U_t$ instead of $V_\pi(s)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(U_t - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}}V(S_t, \boldsymbol{\theta}_t)$
 - The MC target $S_t \rightarrow G_t$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}}V(S_t, \boldsymbol{\theta}_t)$
 - The TD(0) target $S_t \rightarrow R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}}V(S_t, \boldsymbol{\theta}_t)$
 - The TD(λ) target $S_t \rightarrow G_t^\lambda$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t^\lambda - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}}V(S_t, \boldsymbol{\theta}_t)$

Value Function Approximation

Prediction Algorithms

▶ Monte-Carlo with Value Function Approximation

- Return G_t is an unbiased, noisy sample of true value $V_\pi(s)$
- Training data:
 - $\langle S_1, G_1 \rangle, \langle S_2, G_2 \rangle, \dots, \langle S_T, G_T \rangle$
- e.g., linear MC policy evaluation
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - $\quad\quad = \boldsymbol{\theta}_t + \alpha(G_t - V(S_t, \boldsymbol{\theta}_t))\mathbf{x}_t$

▶ MC evaluation converges to a local optimum

▶ Even when using non-linear value function approximation

Value Function Approximation

Prediction Algorithms

► TD Learning with Value Function Approximation

- The TD-target $R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t)$ is a biased sample of true value $V_\pi(s)$
- Can still apply supervised learning to training data:
 - $\langle S_1, R_2 + \gamma V(S_2, \boldsymbol{\theta}_1) \rangle, \langle S_2, R_3 + \gamma V(S_3, \boldsymbol{\theta}_2) \rangle, \dots, \langle S_{T-1}, R_T \rangle$
- e.g., linear TD(0)
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t)) \nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - $\quad = \boldsymbol{\theta}_t + \alpha(R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t)) \mathbf{x}_t$
 - $\quad = \boldsymbol{\theta}_t + \alpha \delta_t \mathbf{x}_t$

► Linear TD(0) converges (close) to global optimum

Value Function Approximation

Prediction Algorithms

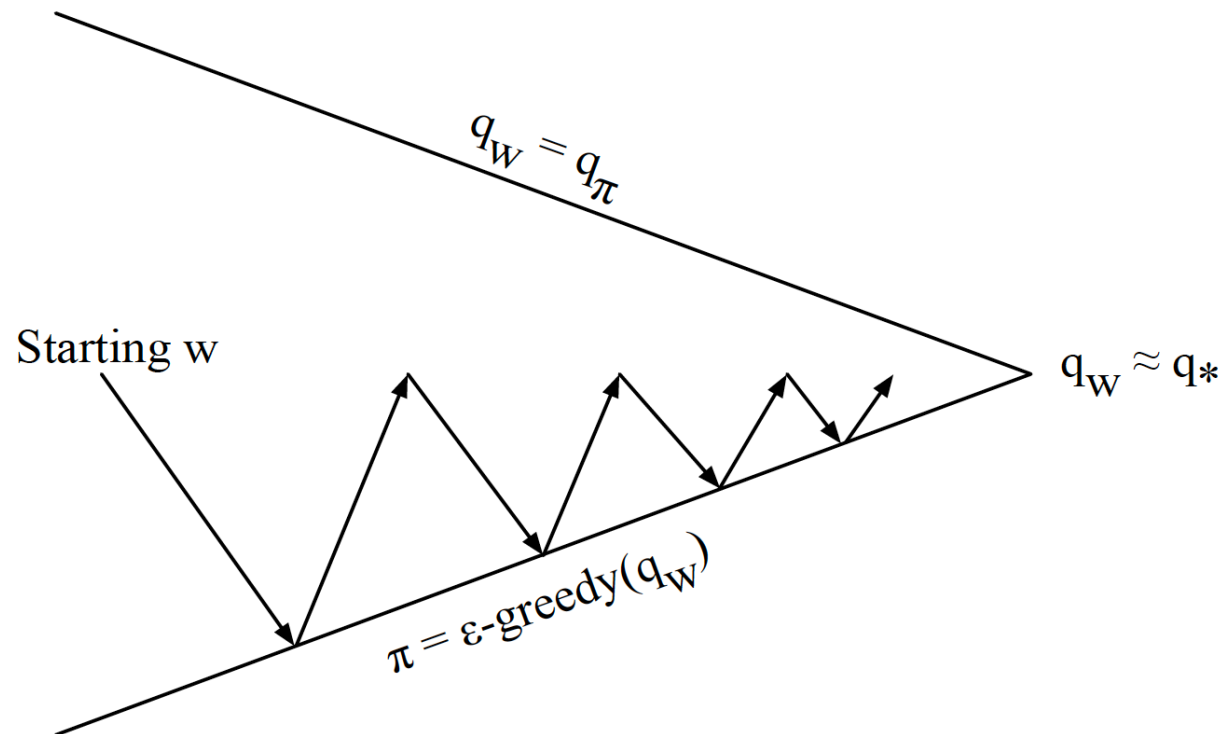
► TD(λ) Learning with Value Function Approximation

- The λ -return G_t^λ is a biased sample of true value $V_\pi(s)$
- Can still apply supervised learning to training data:
 - $\langle S_1, G_1^\lambda \rangle, \langle S_2, G_2^\lambda \rangle, \dots, \langle S_{T-1}, G_{T-1}^\lambda \rangle$
- Forward view linear TD(λ)
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t^\lambda - V(S_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - $\quad = \boldsymbol{\theta}_t + \alpha(G_t^\lambda - V(S_t, \boldsymbol{\theta}_t))\mathbf{x}_t$
- Backward view linear TD(λ)
 - In proportion to TD-error δ_t and eligibility trace $E_t(s)$
 - $\delta_t = R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t)$
 - $E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{x}_t$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha\delta_t \mathbf{E}_t$

Value Function Approximation

Control Algorithms

- ▶ Policy evaluation: Approximate policy evaluation $Q(\cdot, \cdot, \theta) \approx Q_\pi$
- ▶ Policy improvement: ϵ -greedy policy improvement



Value Function Approximation

Control Algorithms

► Action-Value Function Approximation

- Approximate the action-value function
 - $Q(s, a, \boldsymbol{\theta}) \approx Q_{\pi}(s, a)$
- The MSE Objective
 - $\text{MSE}(\boldsymbol{\theta}) = \sum_{s \in S, a \in A} \mu(s) [Q_{\pi}(s, a) - Q(s, a, \boldsymbol{\theta})]^2$
 - $= \mathbb{E}_{\pi} [[Q_{\pi}(s, a) - Q(s, a, \boldsymbol{\theta})]^2]$
- Use SGD to find a local minimum
 - $\Delta \boldsymbol{\theta} = \alpha (Q_{\pi}(s, a) - Q(s, a, \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}} Q(s, a, \boldsymbol{\theta})$

Value Function Approximation

Linear Action-Value Function Approximation

► Feature Vectors

- Represent state by a feature vector

- $\mathbf{x}(s, a) = \begin{pmatrix} x_1(s, a) \\ \vdots \\ x_n(s, a) \end{pmatrix}$

► Linear methods

- Represent value function by a linear combination of features

- $Q(s, a, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}(s, a) = \sum_{i=1}^n \theta_i x_i(s, a)$

► SGD update

- $\nabla_{\boldsymbol{\theta}} Q(s, a, \boldsymbol{\theta}) = \mathbf{x}(s, a)$
- $\Delta \boldsymbol{\theta} = \alpha (Q_{\pi}(s, a) - Q(s, a, \boldsymbol{\theta})) \mathbf{x}$

Value Function Approximation

Control Algorithms

- ▶ Like prediction, use target $s \rightarrow U_t$ instead of $Q_\pi(s, a)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(U_t - Q(S_t, A_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - The MC target $S_t \rightarrow G_t$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(G_t - Q(S_t, A_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - The Sarsa(0) target $S_t \rightarrow R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(R_{t+1} + Q(S_{t+1}, A_{t+1}, \boldsymbol{\theta}_t) - Q(S_t, A_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - The Sarsa(λ) target $S_t \rightarrow Q_t^\lambda$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha(Q_t^\lambda - Q(S_t, A_t, \boldsymbol{\theta}_t))\nabla_{\boldsymbol{\theta}} Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - The backview Sarsa(λ)
 - $\delta_t = R_{t+1} + Q(S_{t+1}, A_{t+1}, \boldsymbol{\theta}_t) - Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - $\mathbf{E}_t(s) = \gamma\lambda\mathbf{E}_{t-1}(s) + \nabla_{\boldsymbol{\theta}} Q(S_t, A_t, \boldsymbol{\theta}_t)$
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha\delta_t\mathbf{E}_t$

Value Function Approximation

Control Algorithms

► Example: Episodic 1-step semi-gradient Sarsa

- Select action and improve policy using an ε -greedy action
- Converges the same ways as TD(0) with same error bound
 - $\text{MSE}(\boldsymbol{\theta}_{TD}) \leq \frac{1}{1-\gamma} \min_{\boldsymbol{\theta}} \text{MSE}(\boldsymbol{\theta})$

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

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$

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

Loop for each episode:

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

 Loop for each step of episode:

 Take action A , observe R, S'

 If S' is terminal:

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

 Go to next episode

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

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

$S \leftarrow S'$

$A \leftarrow A'$

Value Function Approximation

Off-policy Methods with Approximation

- ▶ **Off-policy learning**

- learn a value function for a target policy π
- given data due to a different behavior policy μ

- ▶ **In the prediction case, both policies are static and given**

- ▶ **In the control case, both policies typically change during learning**

- π is the greedy policy
- μ is ε -greedy policy

Value Function Approximation

Off-policy Methods with Approximation

- ▶ **the first challenge of off-policy learning with function approximation**
 - changing the update targets
- ▶ **Semi-gradient Methods**
 - Use per-step importance sampling ratio
 - $\rho_t = \rho_{t:t} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)}$
- ▶ **semi-gradient off-policy TD(0)**
 - $\theta_{t+1} = \theta_t + \alpha \rho_t (R_{t+1} + \gamma V(S_{t+1}, \theta_t) - V(S_t, \theta_t)) \nabla_{\theta} V(S_t, \theta_t)$
 - $= \theta_t + \alpha \rho_t \delta_t \nabla_{\theta} V(S_t, \theta_t)$

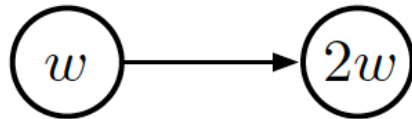
Value Function Approximation

Off-policy Divergence

► the Second challenge of off-policy learning with function approximation

- changing the update distribution: The distribution of updates does not match the on-policy distribution

► Example



- TD error
 - $\delta_t = R_{t+1} + \gamma V(S_{t+1}, \boldsymbol{\theta}_t) - V(S_t, \boldsymbol{\theta}_t)$
 - $= 0 + \gamma 2\boldsymbol{\theta}_t - \boldsymbol{\theta}_t$
 - $= (2\gamma - 1)\boldsymbol{\theta}_t$
- semi-gradient off-policy TD(0)
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \rho_t \delta_t \nabla_{\boldsymbol{\theta}} V(S_t, \boldsymbol{\theta}_t)$
 - $= \boldsymbol{\theta}_t + \alpha 1(2\gamma - 1)\boldsymbol{\theta}_t 1$
 - $= (1 + \alpha(2\gamma - 1))\boldsymbol{\theta}_t$

Value Function Approximation

Instability and Divergence

- ▶ **Danger of instability and divergence arises whenever we combine all of the following three elements**
 - Function approximation
 - Bootstrapping
 - Off-policy training

Value Function Approximation

Convergence

► Convergence of Prediction Algorithms

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	TD(0)	✓	✓	✗
	TD(λ)	✓	✓	✗
Off-Policy	MC	✓	✓	✓
	TD(0)	✓	✗	✗
	TD(λ)	✓	✗	✗

Value Function Approximation

Convergence

► Gradient TD Learning

- TD does not follow the gradient of any objective function
- This is why TD can diverge when off-policy or using non-linear function approximation
- Gradient TD follows true gradient of projected Bellman error

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	TD	✓	✓	✗
	Gradient TD	✓	✓	✓
Off-Policy	MC	✓	✓	✓
	TD	✓	✗	✗
	Gradient TD	✓	✓	✓

Value Function Approximation

Convergence

► Convergence of Control Algorithms

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
Gradient Q-learning	✓	✓	✗

Value Function Approximation

Linear Function Approximation

► Nice properties of linear method

- The gradient is very simple:
 - $\nabla_{\theta} V(s, \theta) = \mathbf{x}(s)$
- For MSE, the error surface is simple
 - quadratic surface with a single minimum
- Linear gradient descent TD(λ) converges:
 - Step size decreases appropriately
 - On-line sampling (states sampled from the on-policy distribution)
 - Converges to parameter vector with property:
 - $\text{MSE}(\theta_{TD}) \leq \frac{1}{1-\gamma} \min_{\theta} \text{MSE}(\theta)$
- Linear gradient descent TD(0) converges:
 - $\text{MSE}(\theta_{TD}) \leq \frac{1-\gamma\lambda}{1-\gamma} \min_{\theta} \text{MSE}(\theta)$

Value Function Approximation

Feature Construction for Linear Methods

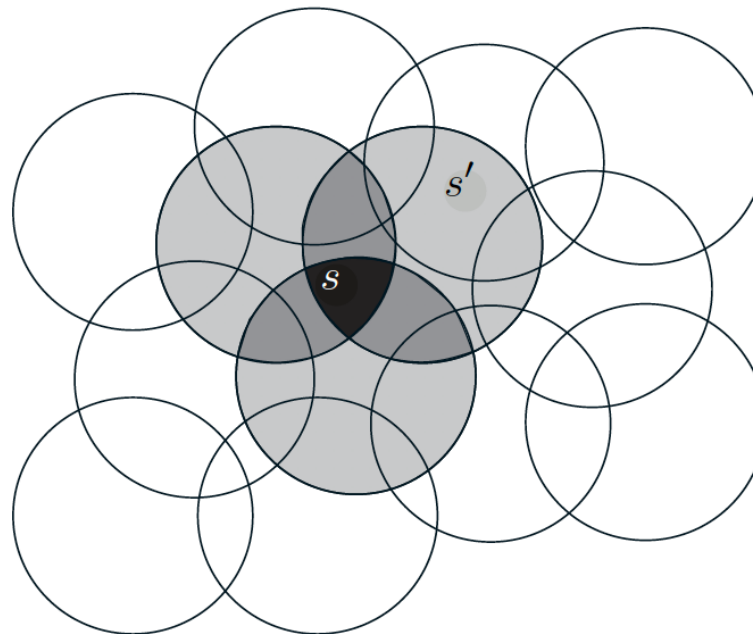
- ▶ **Polynomials**

- e.g., $\mathbf{x}(s) = (1, s_1, s_2, s_1^2, s_2^2, s_1 s_2)^T$

- ▶ **Fourier Basis**

- e.g., $x_i(s) = \cos(i\pi s)$

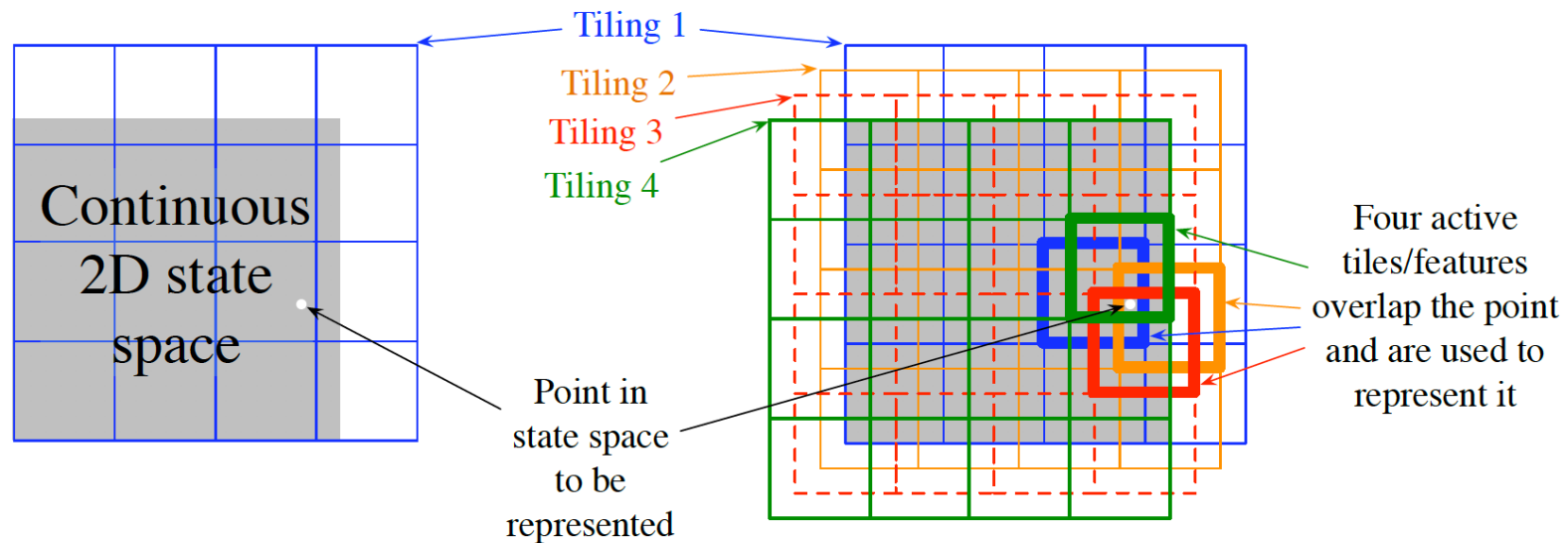
- ▶ **Coarse Coding**



Value Function Approximation

Feature Construction for Linear Methods

► Tile Coding



► Radial Basis Functions

$$- x_i(s) = \exp\left(-\frac{\|s-c_i\|^2}{2\sigma_i^2}\right)$$

Value Function Approximation

Selecting Step-Size Parameters

▶ The classical choice

- $\alpha_t = \frac{1}{t}$
- not appropriate for TD methods or function approximation

▶ Tabular case

- $\alpha_t = \frac{1}{\tau}$
- tabular estimate for a state will approach the mean of its targets after τ experiences with the state

▶ general function approximation

- $\alpha = (\tau E[\mathbf{x}^T \mathbf{x}])^{-1}$

Value Function Approximation

Nonlinear Function Approximation

► Artificial Neural Networks

- widely used for nonlinear function approximation

