

Reinforcement Learning

Lecture 8

Value Function Approximation (v and q)

DeepQ Learning

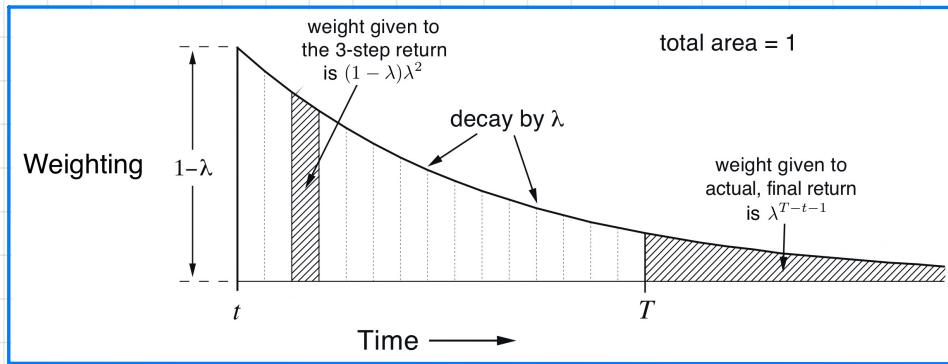
Robert Peharz

Institute of Theoretical Computer Science
Graz University of Technology
Winter Term 2023/24

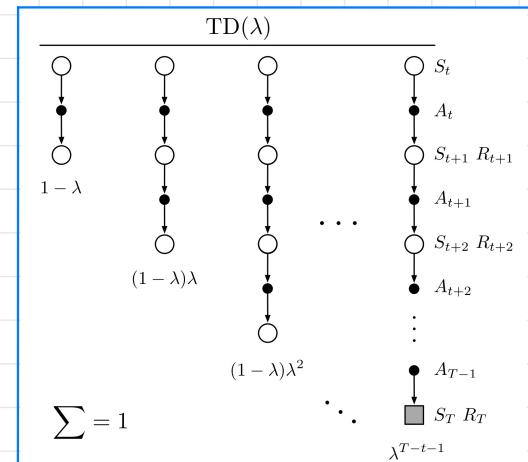
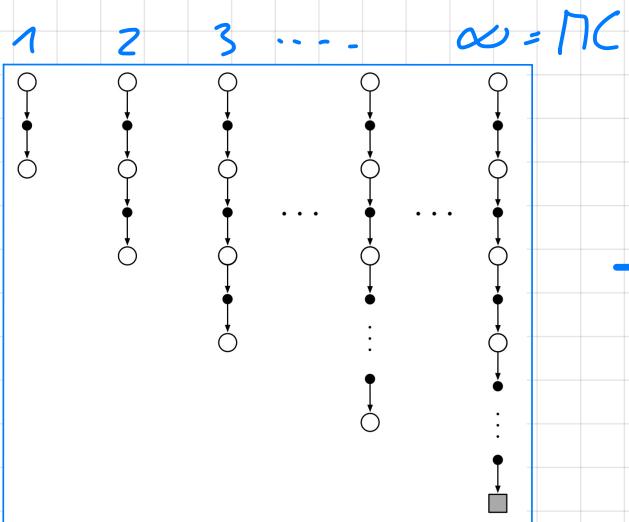
Recap: λ -Return

$$\sum_{n=1}^{\infty} (1-\lambda) \lambda^{n-1} = 1 \quad 0 \leq \lambda \leq 1$$

stick-breaking weights



n-step TD



n-step return

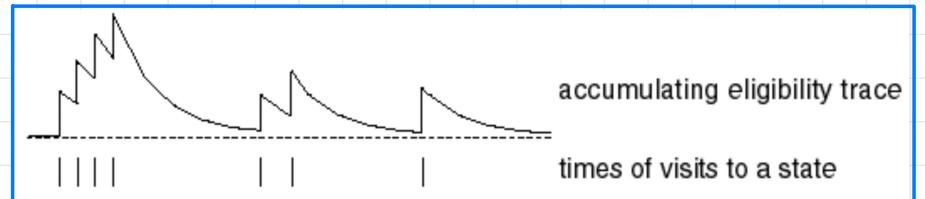
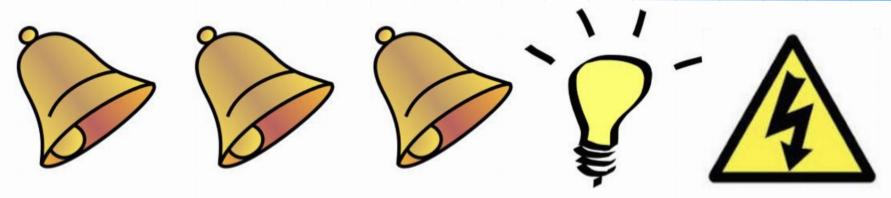
$$\hat{G}_t^n := \sum_{k=0}^{n-1} \gamma^k R_{t+k+1} + \gamma^n V(S_{t+n})$$

using V estimate
= "bootstrapping"

λ -Return

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} \hat{G}_t^n$$

Recap: TD(λ)



$TD(\lambda)$

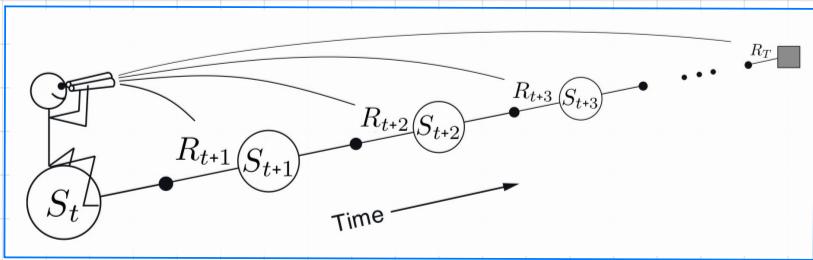
$$e(s) \leftarrow \gamma \rho e(s) + \mathbb{I}[S=s] \quad \nabla s$$

$$\delta \leftarrow R + \gamma v(S') - v(S) \quad // \text{TD error}$$

$$v(s) \leftarrow v(s) + \alpha \delta e(s) \quad \nabla s$$

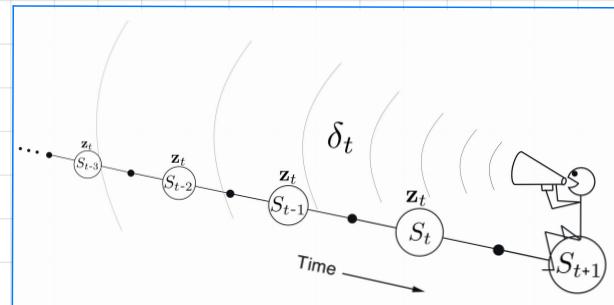
eligibility traces

λ -Return



\approx

$TD(\lambda)$



exactly equivalent, if $TD(\lambda)$ performs accumulated update at end of episode

Value Function Approximation

Why Function Approximation?

- Reinforcement Learning = $\begin{cases} \text{learning } V_{\pi} & (\text{prediction}) \text{ or} \\ \text{learning } q^* & (\text{control}) \end{cases}$
 - So far, we were considering tabular methods, i.e. we were explicitly storing $V(s)$ or $q(s, a)$

$$V = \begin{bmatrix} & \\ & \\ & \\ & \end{bmatrix} \Bigg\} |S^j|$$

$$q = \begin{bmatrix} & & \\ & & \\ & & \\ & & \end{bmatrix} \dots \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} \Bigg\} |S^j|$$

$$\underbrace{\quad}_{|\mathcal{A}|}$$

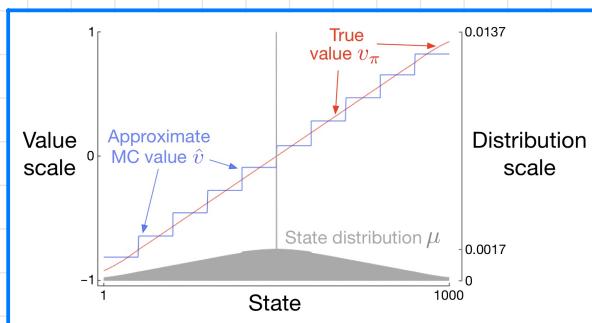
- For "interesting" problems, these are too many values
 - Game of Go: 10^{170} states, continuous state spaces, ...
 - Even when we can store the value functions, updating all states/actions might be too expensive
→ generalize over many states



Function Approximators

Images: Sutton & Barto

- aggregated states



- Linear

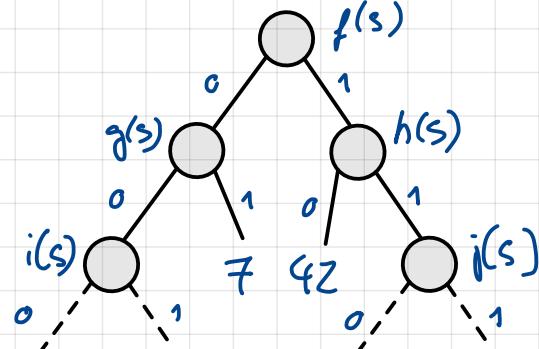
$$v(s) = \sum_{i=1}^K w_i x_i(s)$$

$$q(s,a) = \sum_{i=1}^K w_i x_i(s,a)$$

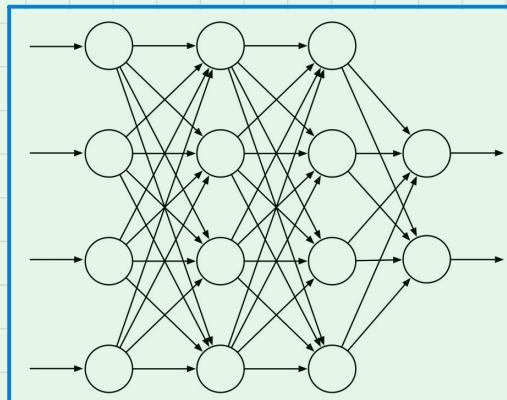
}

x_i ... fixed non-linear features
of s or s,a pair
 w_i ... tunable weights

- Decision Trees



- Neural Networks



Most prominent
nowadays

Differentiable Approximators

- we will consider approximators that are differentiable

$$v(s) \approx \hat{v}(s, \underline{w})$$

$$q(s, a) \approx \hat{q}(s, a, \underline{w})$$

$$\underline{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

- where \underline{w} is a weight vector

- usually $|\underline{w}| = d \ll |S'|$ ($|S'| \times |\mathcal{A}|$)

- we assume we can compute the gradient

$$\nabla_{\underline{w}} \hat{v}(s, \underline{w}) := \begin{pmatrix} \frac{\partial \hat{v}(s, \underline{w})}{\partial w_1} \\ \frac{\partial \hat{v}(s, \underline{w})}{\partial w_2} \\ \vdots \\ \frac{\partial \hat{v}(s, \underline{w})}{\partial w_d} \end{pmatrix}$$

$$\nabla_{\underline{w}} \hat{q}(s, a, \underline{w}) := \begin{pmatrix} \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_1} \\ \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_2} \\ \vdots \\ \frac{\partial \hat{q}(s, a, \underline{w})}{\partial w_d} \end{pmatrix}$$

- in particular, neural networks and linear models
- neural nets allow for automatic differentiation (backprop)

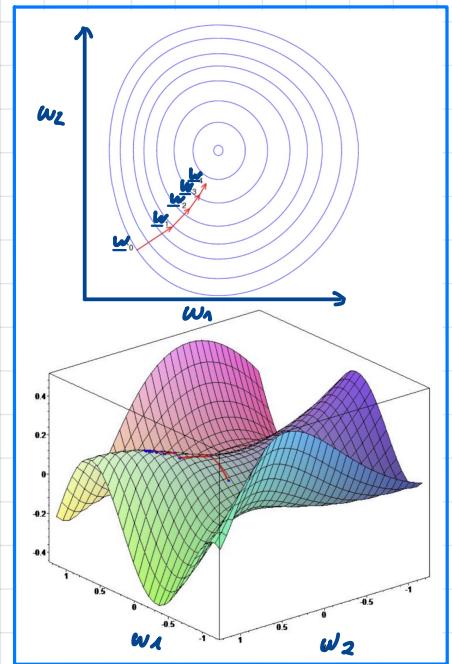
Gradient Descent for Supervised Learning

Image: D. Silver

- gradient is the direction of steepest increase
- to minimize a function $f(\underline{w})$, repeatedly take small step in direction of negative gradient

Gradient Descent:

```
init  $\underline{w}$ 
repeat
     $\underline{w} \leftarrow \underline{w} - \alpha \nabla_{\underline{w}} f(\underline{w})$  stepsize
    until  $\nabla_{\underline{w}} f(\underline{w}) \approx 0$ 
```



in supervised learning, we have a dataset $\mathcal{D} = \{(\underline{x}_1, y_1), (\underline{x}_2, y_2), \dots, (\underline{x}_N, y_N)\}$ where \underline{x}_i is an input vector and y_i is a corresponding target.

y_i might be a discrete value (classification)
or continuous value (regression)

goal: learn to predict y from \underline{x} , i.e. $y \approx f(\underline{x}, \underline{w})$

$$\min_{\underline{w}} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N l(f(\underline{x}_i, \underline{w}), y_i)$$

↑
training loss

where $l(\cdot, \cdot)$ is a loss function comparing $f(\underline{x}, \underline{w})$ and y

Stochastic Gradient Descent (SGD)

- training loss \mathcal{L} and its gradient need full pass over N samples
- if N is large, vanilla gradient descent becomes slow
- SGD: Use a small batch of K randomly selected samples:

$$\mathcal{L} := \frac{1}{N} \sum_{i=1}^N \ell(f(x_i, w), y_i) \approx \frac{1}{K} \sum_{i=1}^K \ell(f(x_i, \underline{w}), y_i) =: \mathcal{L}_K$$

$$\nabla_{\underline{w}} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \nabla_w \ell(\dots) \approx \frac{1}{K} \sum_{i=1}^K \nabla_w \ell(\dots) = \nabla_{\underline{w}} \mathcal{L}_K \quad K \ll N$$

- batch loss \mathcal{L}_K and batch gradient $\nabla_{\underline{w}} \mathcal{L}_K$ are noisy but unbiased

$$\mathbb{E}[\mathcal{L}_K] = \mathcal{L} \quad \mathbb{E}[\nabla_{\underline{w}} \mathcal{L}_K] = \nabla_{\underline{w}} \mathcal{L}$$

let α_i be the SGD stepsize at iteration i

if $\sum_{i=1}^{\infty} \alpha_i = \infty$ and $\sum_{i=1}^{\infty} \alpha_i^2 < \infty$: SGD converges to
a local minimum of \mathcal{L} .
(e.g. $\alpha = 1/i$)

Value Function Approximation via SGD

- treat value function approximation as a supervised problem
- using quadratic loss (least squares)

$$l(\hat{v}(s, \underline{w}), v_{\pi}(s)) = \frac{1}{2} (v_{\pi}(s) - \hat{v}(s, \underline{w}))^2 \quad (\text{P}) \text{rediction}$$

$$\nabla_{\underline{w}} l = -(v_{\pi}(s) - \hat{v}(s, \underline{w})) \nabla_{\underline{w}} \hat{v}(s, \underline{w})$$

$$l(\hat{q}(s, a, \underline{w}), q_{\pi}(s, a, \underline{w})) = \frac{1}{2} (q_{\pi}(s, a, \underline{w}) - \hat{q}(s, a, \underline{w}))^2 \quad (\text{C}) \text{ontrol}$$

$$\nabla_{\underline{w}} l = -(q_{\pi}(s, a) - \hat{q}(s, a, \underline{w})) \nabla_{\underline{w}} \hat{q}(s, a, \underline{w})$$

- follow policy and perform SGD updates: $\underline{w} \leftarrow \underline{w} - \alpha \nabla_{\underline{w}} l$
- underlying objectives: $\mathbb{E}_{\pi} \left[\frac{1}{2} (v_{\pi}(s) - \hat{v}(s, \underline{w}))^2 \right] \quad \mathbb{E}_{\pi} \left[\frac{1}{2} (q_{\pi}(s, a) - \hat{q}(s, a, \underline{w}))^2 \right]$
- what is the problem with this algorithm?

Monte Carlo Function Approximation

Recall and compare

Tabular MC update:

$$v(s) \leftarrow v(s) + \alpha (G_t - v(s))$$

$$\nabla_{\underline{w}} l = -(v_{\pi}(s) - \hat{v}(s, \underline{w})) \nabla_{\underline{w}} \hat{v}(s, \underline{w}) \quad (P)$$

$$\nabla_{\underline{w}} l = -(q_{\pi}(s, a) - \hat{q}(s, a, \underline{w})) \nabla_{\underline{w}} \hat{q}(s, a, \underline{w}) \quad (C)$$

- the target values $v_{\pi}(s)$ ($q_{\pi}(s, a)$) are not known
- if they were, no need to approximate them ...
- idea: use MC to get unbiased estimate of gradients

initialize \underline{w}

repeat // for each episode

generate episode with $\hat{\pi}$ (ϵ -greedy(q))

$$g \leftarrow 0$$

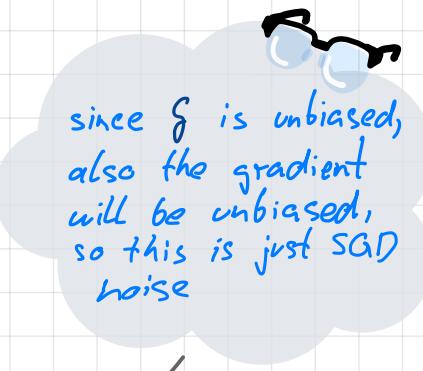
for $t = T-1 \dots 0$:

$$g \leftarrow R_{t+1} + \gamma g$$

$$\underline{w} \leftarrow \underline{w} + \alpha (g - \hat{v}(s_t, \underline{w})) \nabla_{\underline{w}} \hat{v}(s_t, \underline{w})$$

$$(\underline{w} \leftarrow \underline{w} + \alpha (g - \hat{q}(s_t, a_t, \underline{w})) \nabla_{\underline{w}} \hat{q}(s_t, a_t, \underline{w}))$$

since g is unbiased,
also the gradient
will be unbiased,
so this is just SGD
noise



TD(0) and Sarsa Function Approximation

Instead of estimating $V_{\pi}(s)$ ($q_{\pi}(s, a)$) with MC, plug in TD estimates:

$$\nabla_{\underline{w}} \ell = - (V_{\pi}(s) - \hat{V}(s, \underline{w})) \nabla_{\underline{w}} \hat{V}(s, \underline{w}) \quad (P)$$

$$\approx - (R_{t+1} + \gamma \underbrace{\hat{V}(s', \underline{w})}_{\text{"ignored" in the gradient computation}} - \hat{V}(s, \underline{w})) \nabla_{\underline{w}} \hat{V}(s, \underline{w})$$

"ignored" in the gradient computation
→ "semi-gradient" method

$$\nabla_{\underline{w}} \ell = - (q_{\pi}(s, a) - \hat{q}(s, a, \underline{w})) \nabla_{\underline{w}} \hat{q}(s, a, \underline{w}) \quad (C)$$

$$\approx - (R_{t+1} + \gamma \underbrace{\hat{q}(s', a', \underline{w})}_{\text{"semi-gradient" method}} - \hat{q}(s, a, \underline{w})) \nabla_{\underline{w}} \hat{q}(s, a, \underline{w})$$

→ "semi-gradient" method

- bootstrapping \Leftrightarrow "semi-gradient"
- TD does not work reliably with non-linear approximators



TD(0) and Sarsa Function Approximation

initialize \underline{w}

repeat

 initialize S

 select $A \sim \tilde{\pi}(a|S)$ (ϵ -greedy(q))

 repeat

 take action A , observe R, S'

 if S' is terminal:

$$\underline{w} \leftarrow \underline{w} + \alpha [R - \hat{v}(S, \underline{w})] \nabla_{\underline{w}} \hat{v}(S, \underline{w})$$

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

 break

 // start next episode

 select $A' \sim \tilde{\pi}(a'|S)$ (ϵ -greedy(q))

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

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

$S, A \leftarrow S', A'$

(command alternatives in color  are for control, i.e. Sarsa)

Linear Function Approximation

- assume a set of d features x_1, x_2, \dots, x_d where

$$x_i : S \rightarrow \mathbb{R} \text{ (for } \hat{v}) \quad x_i : S \times \mathcal{A} \rightarrow \mathbb{R} \text{ (for } \hat{q})$$

$$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix}$$

- the linear approximator is given as

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

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

- gradients: $\nabla_{\underline{w}} \hat{v} = \underline{x}$

$$\nabla_{\underline{w}} \hat{q} = \underline{x}$$

- objective is expected quadratic loss

$$\mathbb{E}_{\pi} \left[\frac{1}{2} (\hat{v}(s, \underline{w}) - v_{\pi}(s))^2 \right] \quad \mathbb{E}_{\pi} \left[\frac{1}{2} (\hat{q}(s, a, \underline{w}) - q_{\pi}(s, a))^2 \right]$$

- for linear approximations: quadratic in \underline{w} , hence convex

- any local minimum is also a global minimum

Linear Function Approximation cont'd

Linear function approximators subsume tabular methods,
by using indicator features:

$$x_i(s) := \underline{\mathbb{I}}[s = s_i] \quad s_i \in S \quad (\text{for } \hat{v})$$

$$x_i(s, a) := \underline{\mathbb{I}}[(s, a)_i = (s, a)] \quad (s, a)_i \in S \times A \quad (\text{for } \hat{q})$$

Note how the previous two algorithms reduce to tabular methods!

Monte Carlo

tabular: $v(s) \leftarrow v(s) + \alpha (G_t - v(s))$

gradient: $\underline{w} \leftarrow \underline{w} + \alpha (g - \hat{v}(s_t, \underline{w})) \nabla_{\underline{w}} \hat{v}(s_t, \underline{w})$

$$\underline{x} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{pmatrix} s_t$$

TD

tabular: $v(s) \leftarrow v(s) + \alpha (R_{t+1} + \gamma v(s') - v(s))$

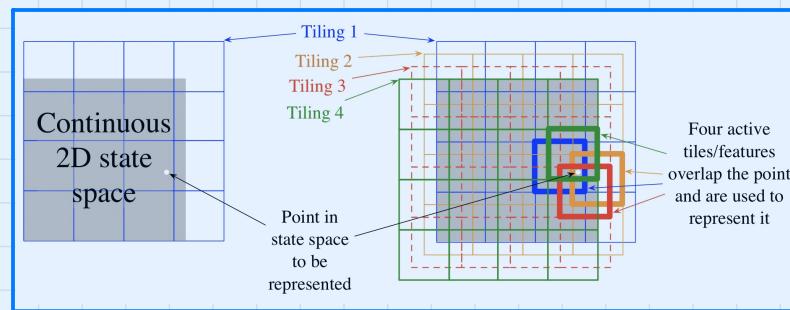
gradient: $\underline{w} \leftarrow \underline{w} + \alpha [R + \gamma \hat{v}(s', \underline{w}) - \hat{v}(s, \underline{w})] \nabla_{\underline{w}} \hat{v}(s, \underline{w})$

Mountain Car with Linear Function

[Sutton & Barto, p. 245]

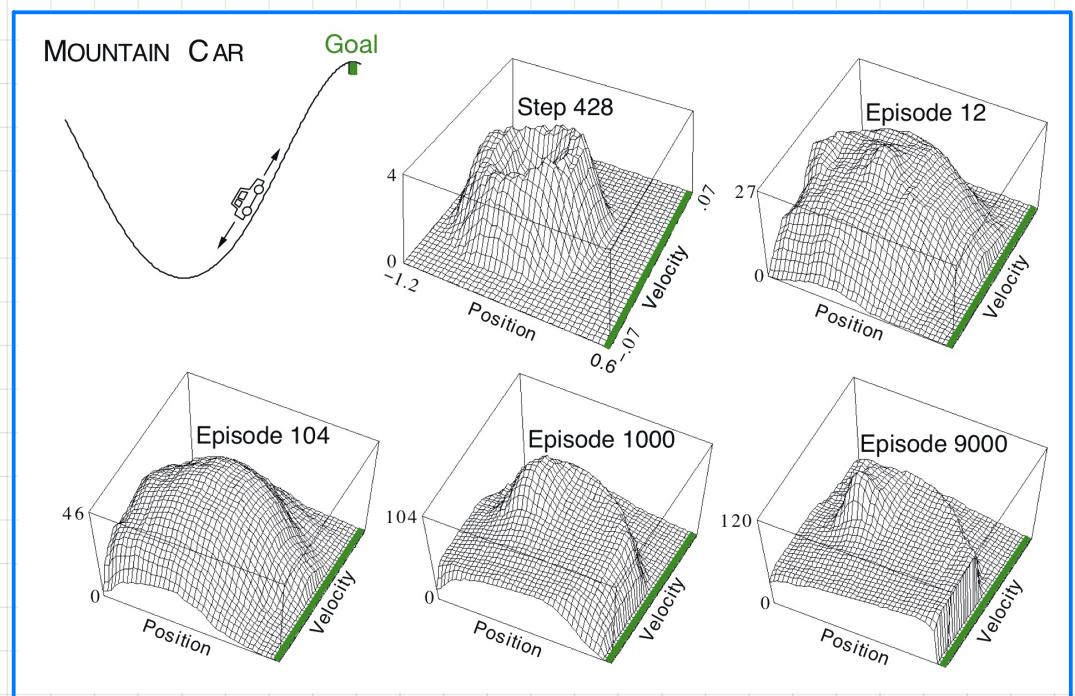
- car with a too weak engine shall go up a hill
→ needs to gather momentum on the other side
- continuous state space: position $x \in [-1.2, 0.6]$, velocity $\dot{x} \in [-0.7, 0.7]$
- actions $A_t = \{-1, 1, 0\}$ (full throttle left/right, no throttle)
- (simplified) physics $\dot{x}_{t+1} = \text{clip}[\dot{x}_t + 0.001 A_t - 0.0025 \cos(3x_t)]$
- features: "tile coding"

- negative value function learned by Sarsa (-1 reward per time step)



Space is discretized into several displaced tilings. For each tiling, we have a vector indicating the containing tile. Overall vector is concatenation of these vectors.

$$16 \times 4 = 64 \text{ binary features}$$



TD(λ) and Sarsa(λ) Function Approximation

initialize \underline{w}

repeat

$$\underline{e} \leftarrow \underline{0}$$

// for each episode

// vector of d zeros

initialize S

$$A \sim \hat{\pi}(a|S) \quad [\epsilon\text{-greedy}(q)]$$

repeat

take action A , observe R, S'

$$A' \sim \hat{\pi}(a|S') \quad [\epsilon\text{-greedy}(q)]$$

$$\underline{e} \leftarrow \gamma \lambda \underline{e} + \nabla_{\underline{w}} \hat{V}(S, \underline{w})$$

$$[\underline{e} \leftarrow \gamma \lambda \underline{e} + \nabla_{\underline{w}} \hat{q}(S, A, \underline{w})]$$

$$\delta \leftarrow R + \gamma \hat{V}(S', \underline{w}) - \hat{V}(S, \underline{w})$$

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

$$\underline{w} \leftarrow \underline{w} + \alpha \delta \underline{e}$$

$$S, A \leftarrow S', A'$$

note that this is
consistent with tabular case

(command alternatives in color  are for control, i.e. Sarsa(λ))

Q-Learning with Function Approximation

initialize w

repeat

initialize S

repeat

select $A \sim \tilde{\pi}(a|s)$ [ϵ -greedy(q)]

take action A , observe R, S'

if s' is terminal:

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

break

// start next episode

$$\underline{w} \leftarrow \underline{w} + \alpha [R + \gamma \max_{a'} \hat{q}(S', a', \underline{w}) - \hat{q}(S, A, \underline{w})] \triangleright_{\underline{w}} \hat{q}(S, A, \underline{w})$$

$$S \leftarrow S'$$

... unstable algorithm



Convergence of Algorithms

Prediction

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

Control

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

(✓) = chatters around near-optimal value function

* not covered here

Playing Atari with Deep Reinforcement Learning

arXiv:1312.5602v1 [cs.LG] 19 Dec 2013

nature

LETTER

[doi:10.1038/nature14236](https://doi.org/10.1038/nature14236)

Human-level control through deep reinforcement learning

Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

Deep Q-Learning

Main Take-aways

- end-to-end Q-learning, involving a ConvNet
- critical component: experience replay on buffered states
- partially-observed MDP: stack last 4 frames
- close-to or super-human play on many Atari games
- however, long-term and symbolic dependencies are still to be learned ("Montezuma's Revenge")