

Reinforcement Learning IV

Model-based Planning and Value Function Approximation

Antoine SYLVAIN

EPITA

2021

Contents

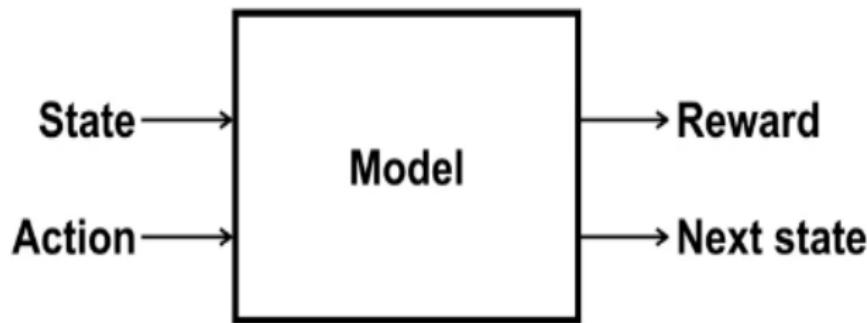
1 Model-based planning

2 Value Function Approximation

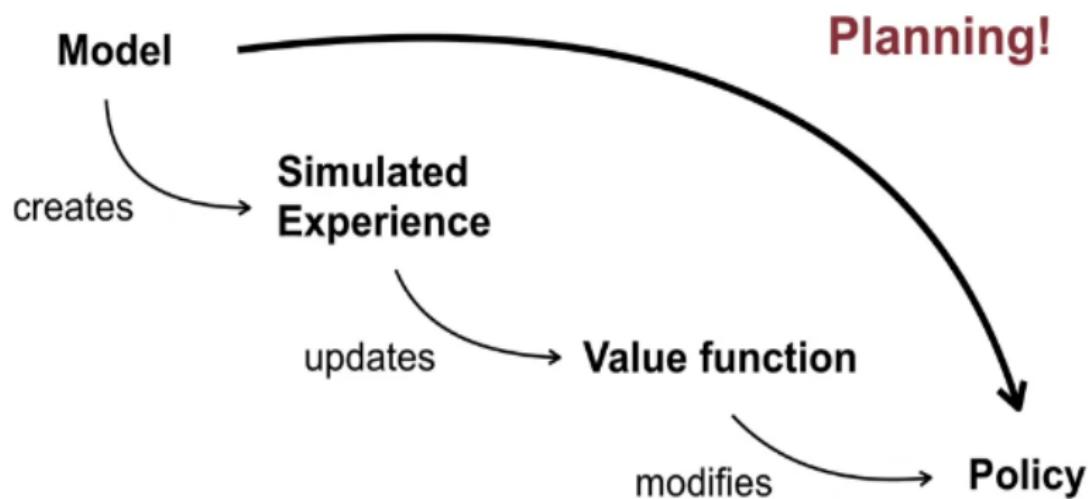
3 Gradient methods

4 Neural Networks

What is a model ?



What is a model ?



Types of models

- We distinguish two types of models:

Types of models

- We distinguish two types of models:
 - **Sample models:** produces an actual outcome drawn from some underlying probabilities
e.g.: generating a sequence of heads and tails by flipping a coin

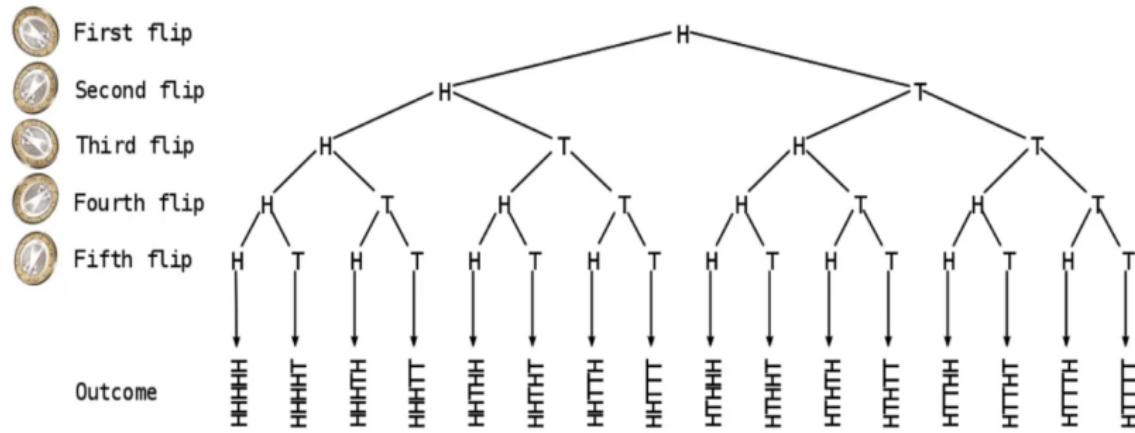
Types of models

- We distinguish two types of models:
 - **Sample models:** produces an actual outcome drawn from some underlying probabilities
e.g.: generating a sequence of heads and tails by flipping a coin
 - **Distribution models:** specifies completely the likelihood or probability of every outcome
e.g.: generating the tree of possible outcomes of a sequence of coins flipping

Sample model

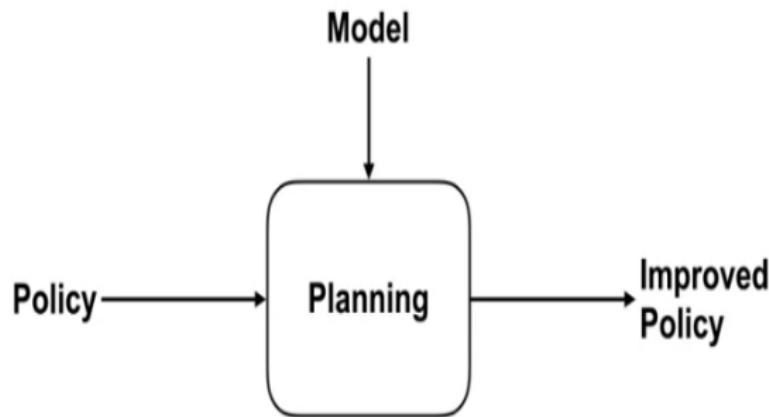
Head, Tail, Tail, Head, Tail

Distribution model



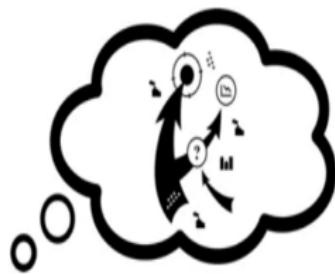
$$p = \frac{1}{32}$$

Planning



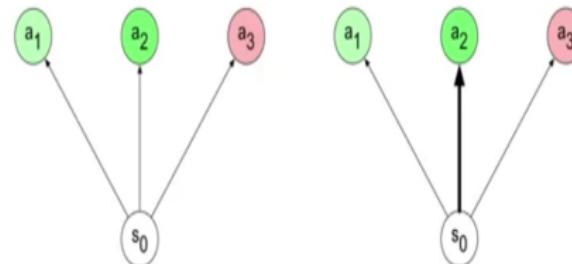
Planning

Planning



$$Q_{\pi}(s, a)$$

$$\pi(s_0) = a_2$$



AB Example

Two states A, B ; no discounting; 8 episodes of experience

A, 0, B, 0

B, 1

B, 1

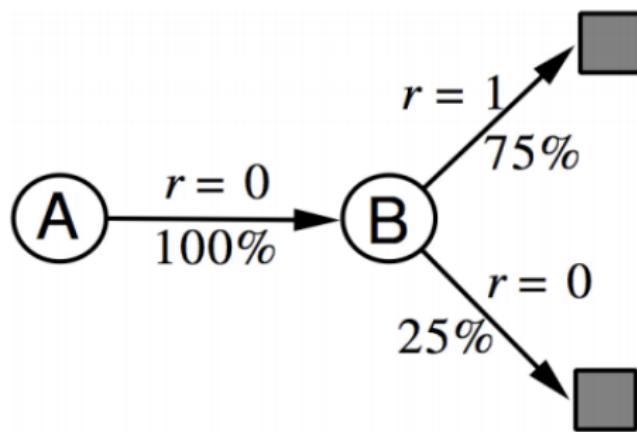
B, 1

B, 1

B, 1

B, 1

B, 0



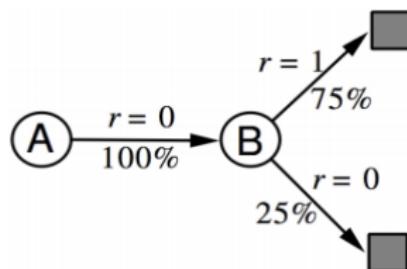
We have constructed a **table lookup model** from the experience

AB Example with a model

- Construct a table-lookup model from real experience
- Apply model-free RL to sampled experience

Real experience

A, 0, B, 0
 B, 1
 B, 0



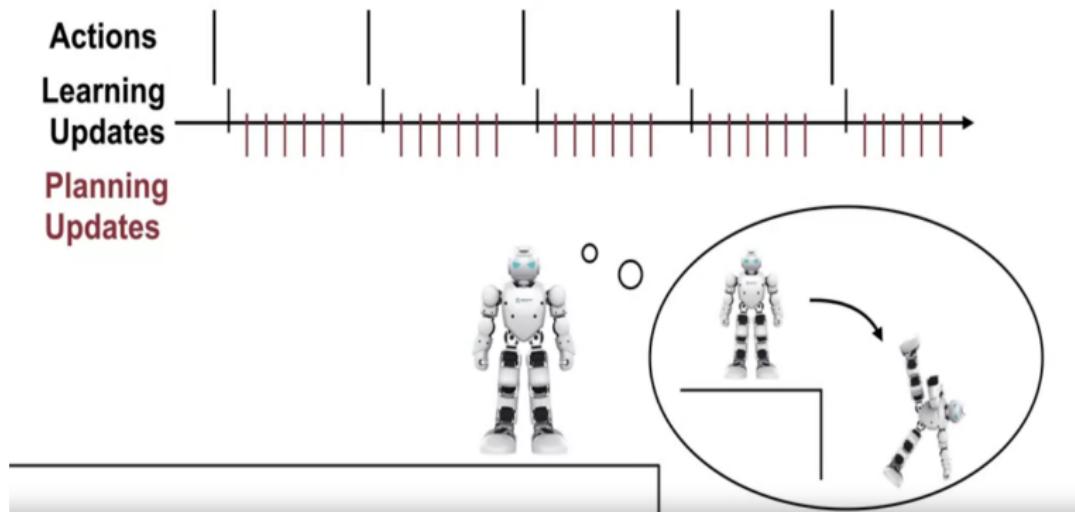
Sampled experience

B, 1
 B, 0
 B, 1
 A, 0, B, 1
 B, 1
 A, 0, B, 1
 B, 1
 B, 0

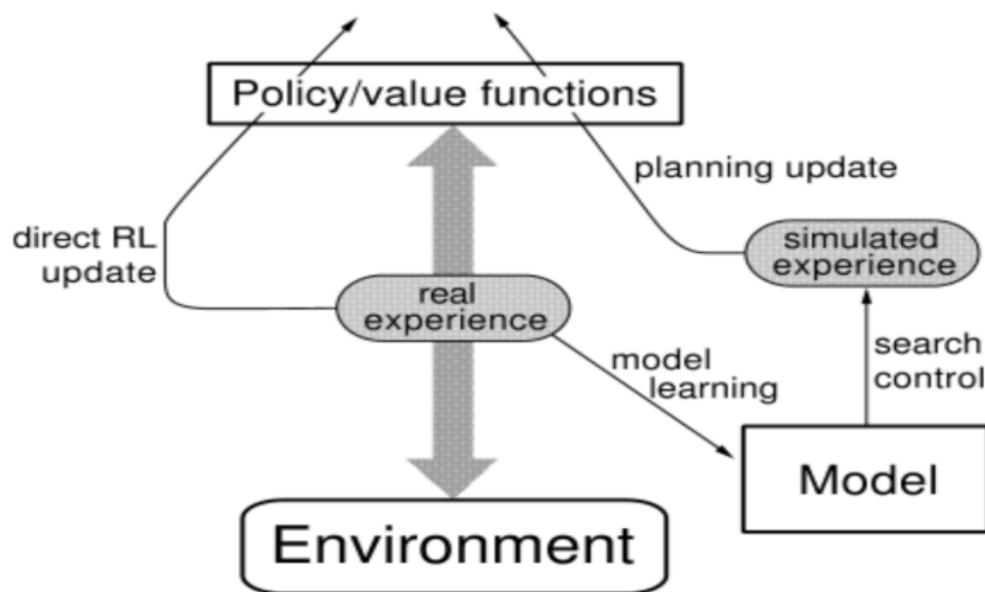
e.g. Monte-Carlo learning: $V(A) = 1, V(B) = 0.75$

Planning

Advantages of Planning



Dyna Architecture



Dyna-Q Algorithm



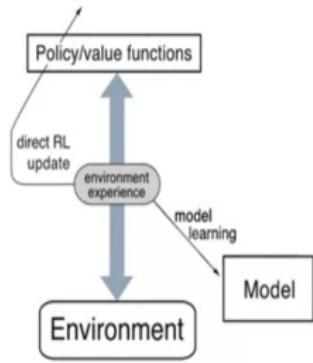
Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- $S \leftarrow$ current (nonterminal) state
- $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- Take action A ; observe resultant reward, R , and state, S'
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

Dyna-Q Algorithm



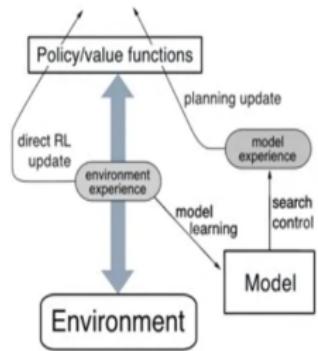
Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- (c) Take action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)

Dyna-Q Algorithm



Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Loop forever:

- $S \leftarrow$ current (nonterminal) state
- $A \leftarrow \varepsilon\text{-greedy}(S, Q)$
- Take action A ; observe resultant reward, R , and state, S'
- $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
- $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- Loop repeat n times:

$S \leftarrow$ random previously observed state

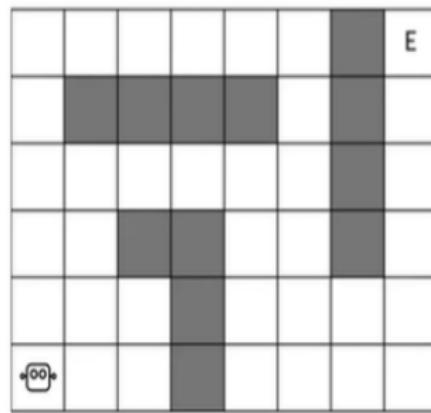
$A \leftarrow$ random action previously taken in S

$R, S' \leftarrow Model(S, A)$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

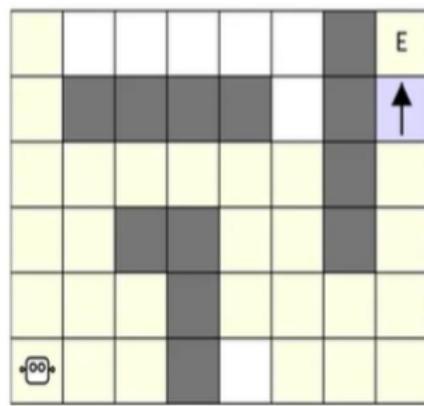
Dyna-Q Maze Example

Number of actions taken: 0



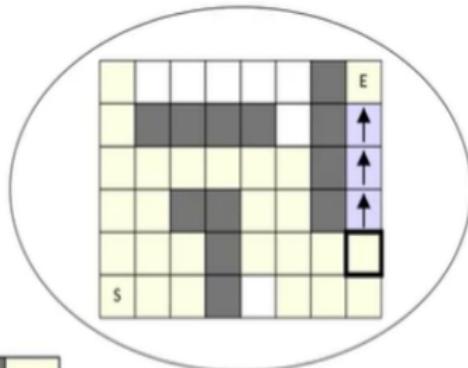
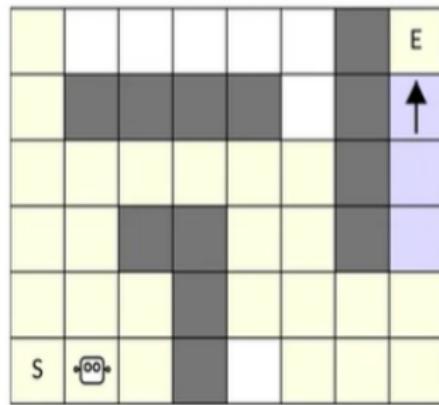
Dyna-Q Maze Example

Number of actions taken: 184



Dyna-Q Maze Example

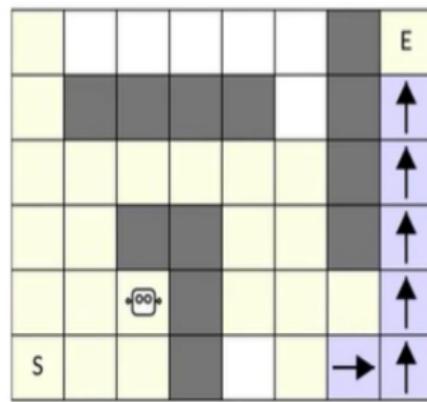
Number of steps planned: 100
Number of actions taken: 185



Dyna-Q Maze Example

Number of steps planned: 400

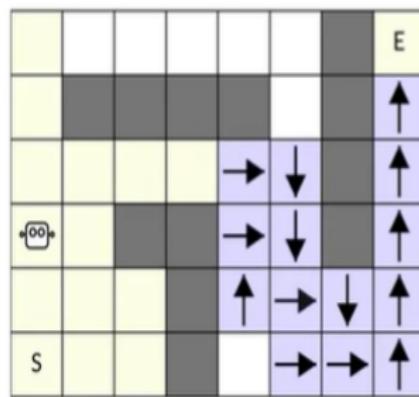
Number of actions taken: 188



Dyna-Q Maze Example

Number of steps planned: 1200

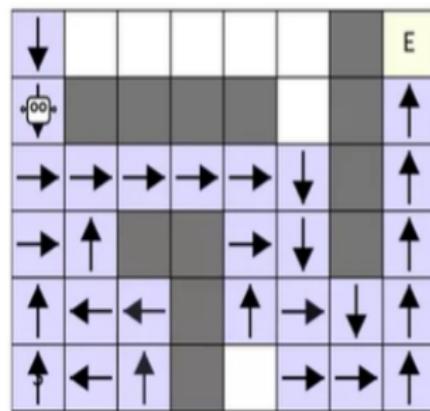
Number of actions taken: 196



Dyna-Q Maze Example

Number of steps planned: 1800

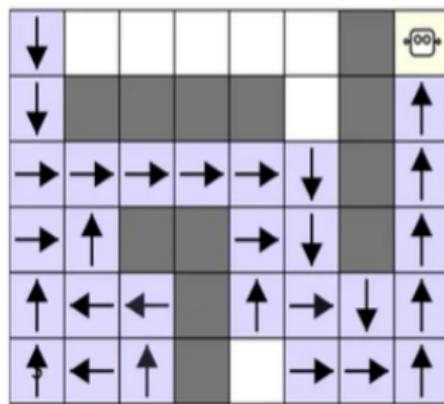
Number of actions taken: 202



Dyna-Q Maze Example

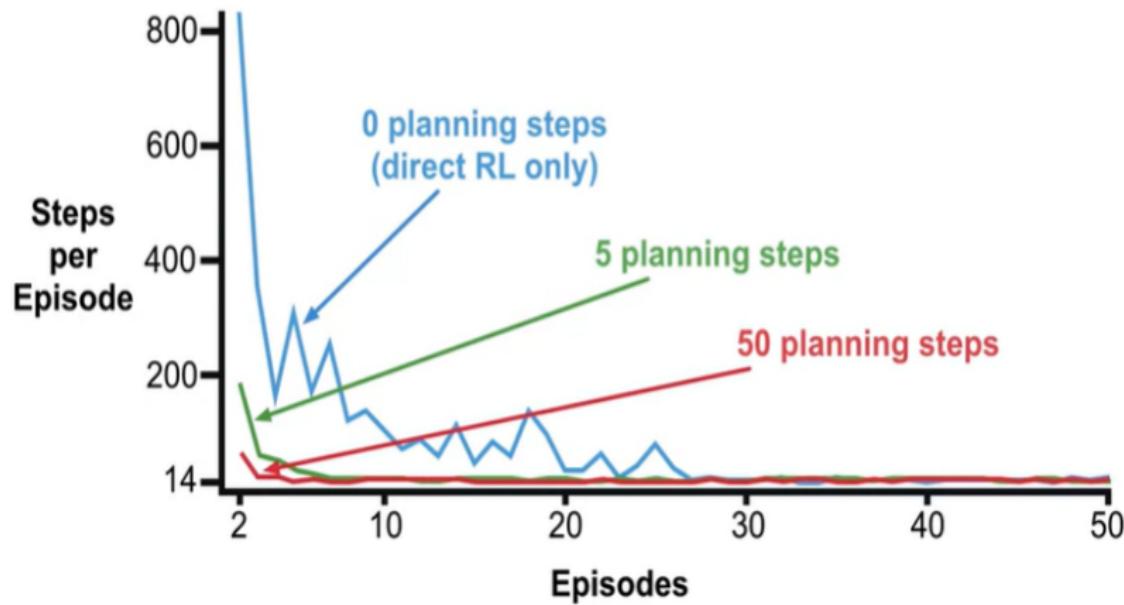
Number of steps planned: 3800

Number of actions taken: 223



Dyna-Q Planning Efficiency

More planning → faster learning



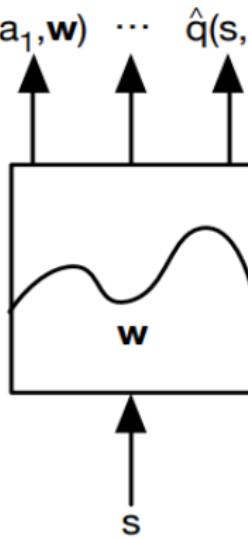
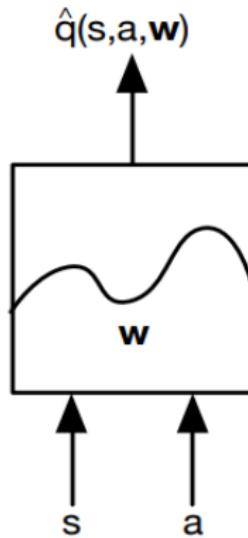
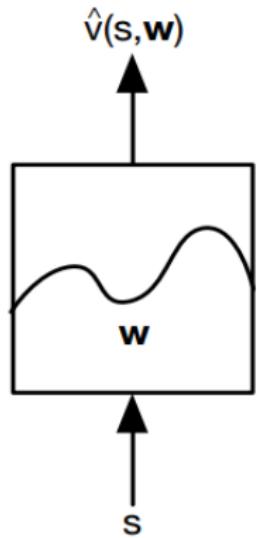
Contents

- 1 Model-based planning
- 2 Value Function Approximation
- 3 Gradient methods
- 4 Neural Networks

Value Function Approximation

- Estimates value function with function approximation:
 $\hat{v}(s, w) \approx v_\pi(s)$
 $\hat{q}(s, a, w) \approx q_\pi(s, a)$
- Generalise from seen states to unseen states
- Update parameter w (e.g. with Monte-Carlo or Temporal-Difference Learning)

Types of Approximators



Function Approximators

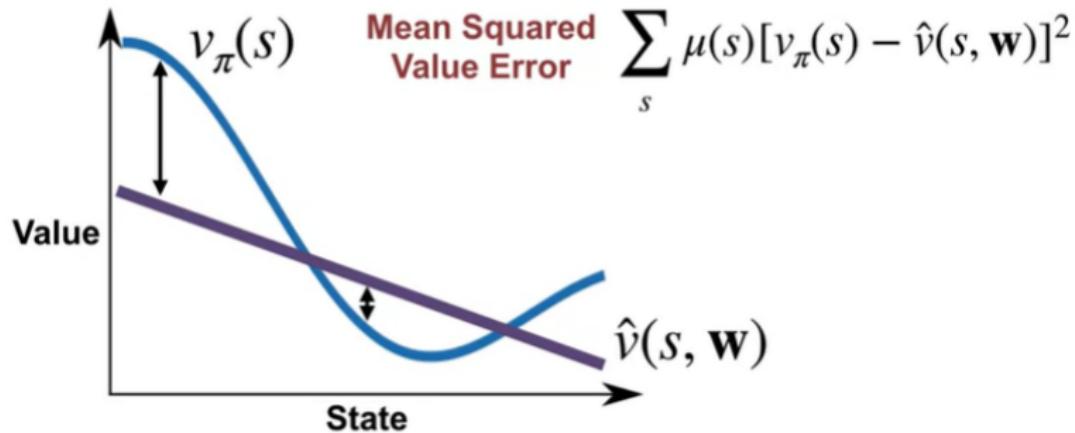
- Linear combination of features
- Neural Networks
- Decision Trees
- Nearest neighbour
- Fourier / wavelets base
- ...

Function Approximators

- Linear combination of features
- Neural Networks
- Decision Trees
- Nearest neighbour
- Fourier / wavelets base
- ...
- Then, a training method

Value Error Objective

The Mean Squared Value Error Objective

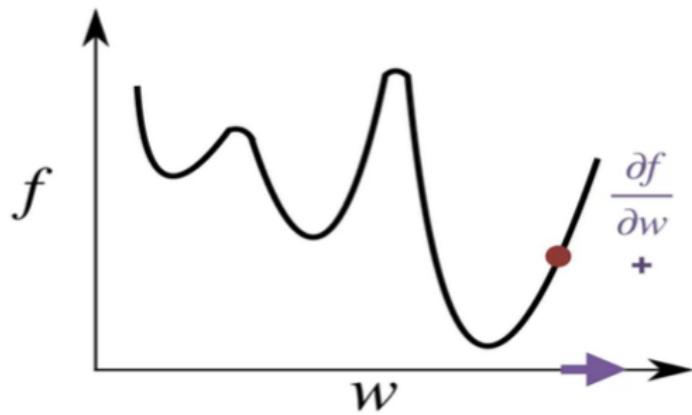


Contents

- 1 Model-based planning
- 2 Value Function Approximation
- 3 Gradient methods
- 4 Neural Networks

Gradient Derivatives

Understanding Derivatives



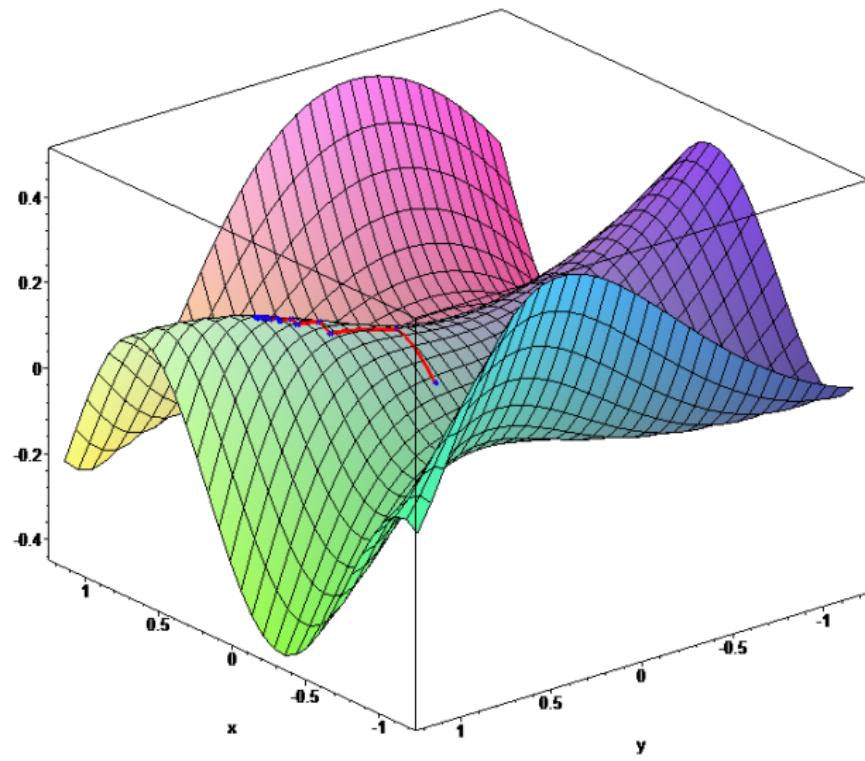
Gradient: Derivatives in multiple dimensions

$$\mathbf{w} \doteq \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_d \end{bmatrix} \quad \nabla f \doteq \begin{bmatrix} \frac{\partial f}{\partial w_1} \\ \frac{\partial f}{\partial w_2} \\ \dots \\ \frac{\partial f}{\partial w_d} \end{bmatrix}$$

The gradient gives the direction of steepest ascent

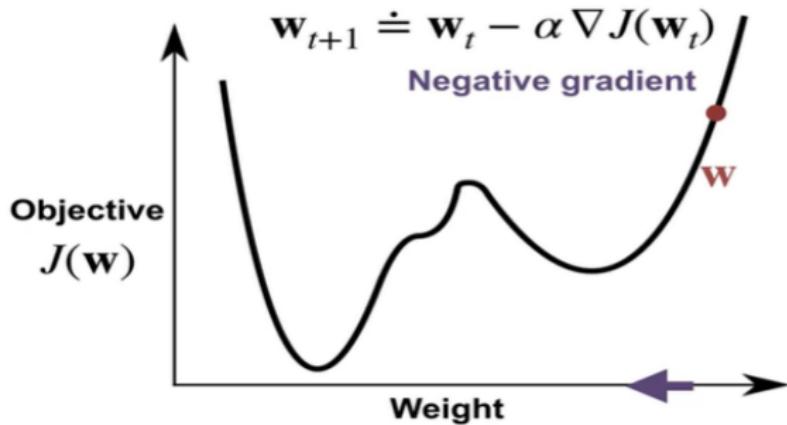
The direction to change w_2 in order to increase f
How quickly f changes

Gradient: Derivatives in multiple dimensions



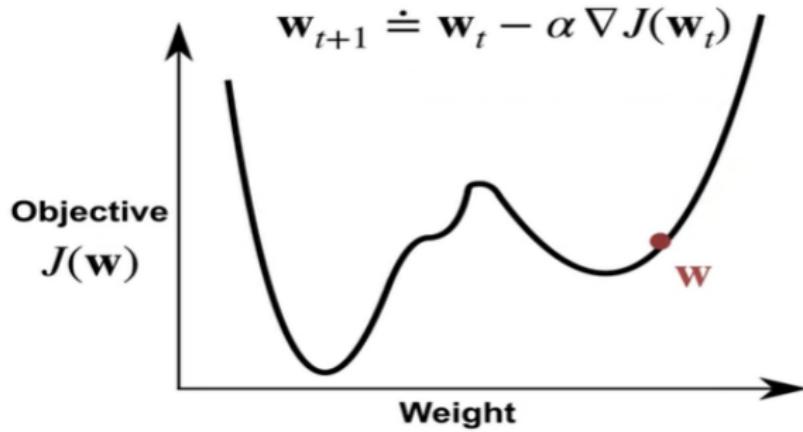
Gradient Descent

Gradient Descent



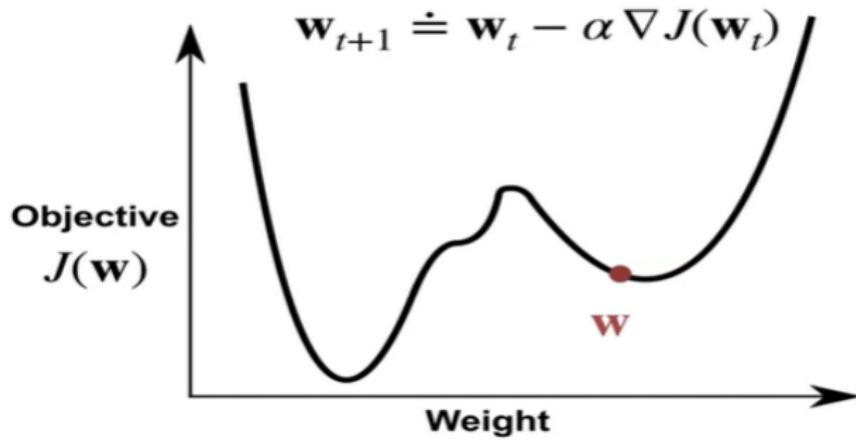
Gradient Descent

Gradient Descent



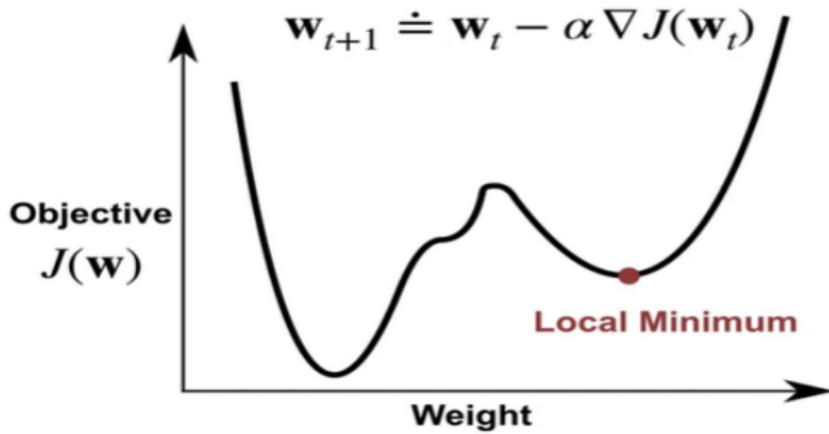
Gradient Descent

Gradient Descent



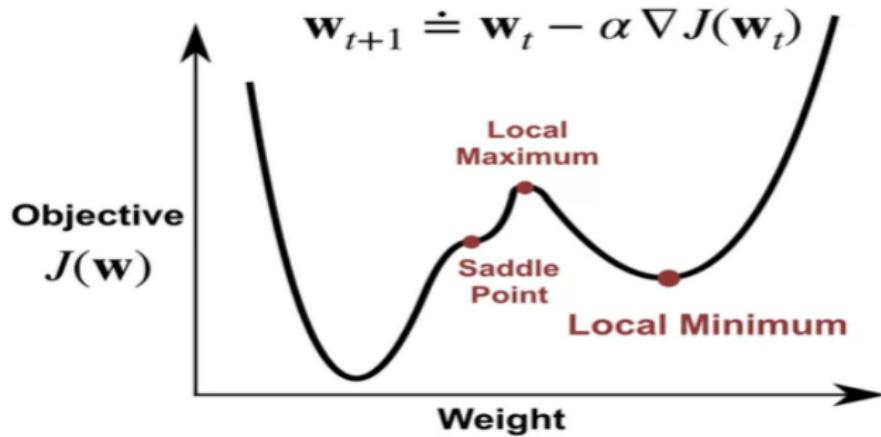
Gradient Descent

Gradient Descent



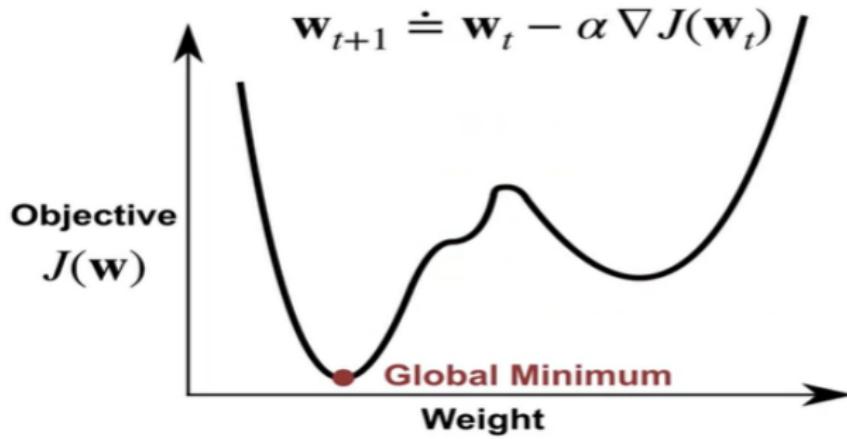
Gradient Descent

Gradient Descent



Gradient Descent

Gradient Descent



Gradient Monte-Carlo

Gradient Monte Carlo

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [G_t - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

Recall that

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi} [G_t \mid S_t = s]$$

$$\begin{aligned} & \mathbb{E}_{\pi} [2[v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})] \\ &= \mathbb{E}_{\pi} [2[G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})] \end{aligned}$$

Gradient MC Algorithm

Input: the policy π to be evaluated

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

Algorithmic parameter: step size $\alpha > 0$

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

Loop (for each episode):

Generate an episode $S_0, A_0, R_1, S_1, A_1, \dots, R_\tau, S_\tau$ using π

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

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

State Aggregation with Monte-Carlo

Random Walk Example



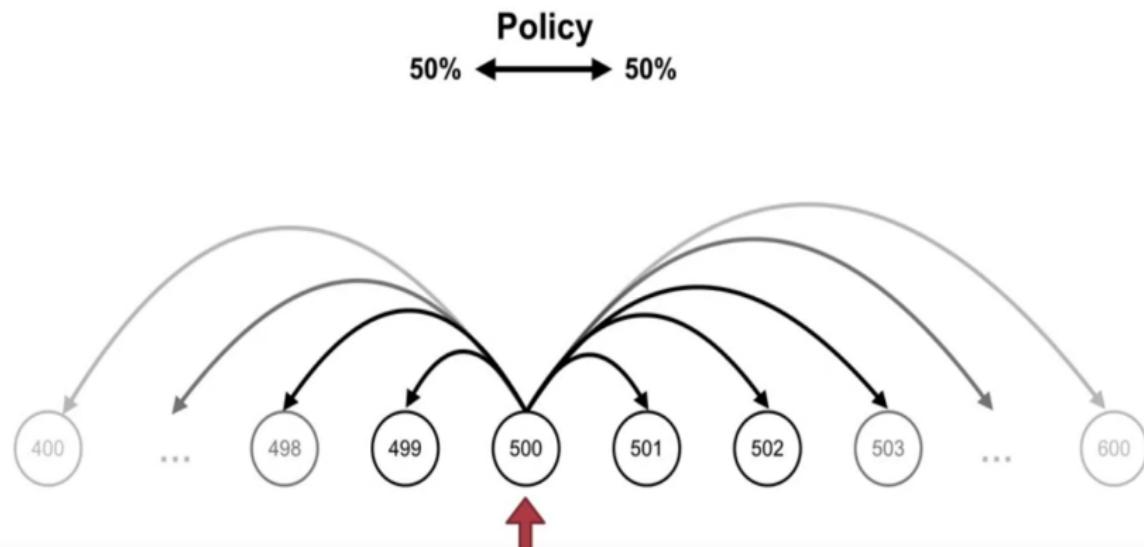
State Aggregation with Monte-Carlo

Random Walk Example



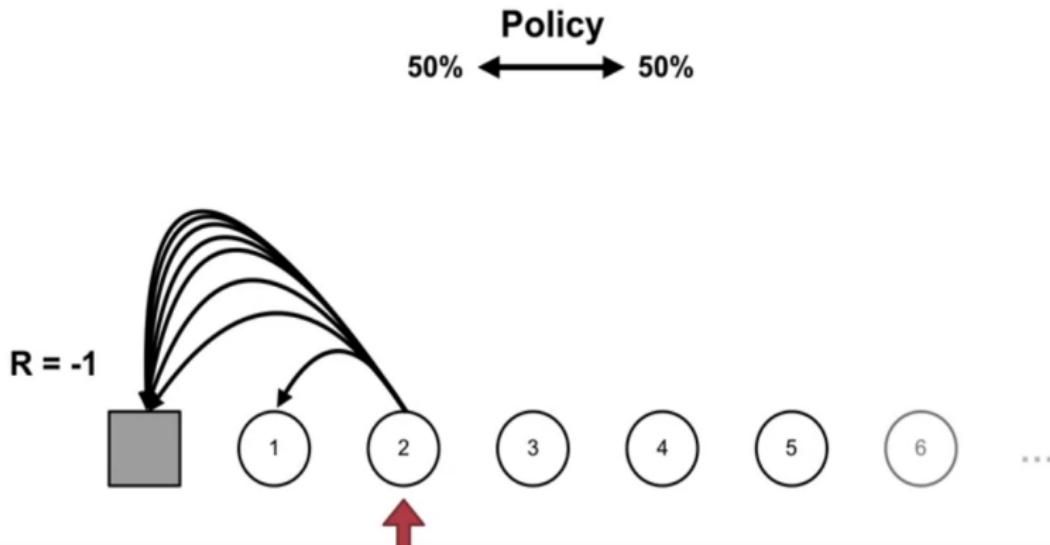
State Aggregation with Monte-Carlo

Random Walk Example



State Aggregation with Monte-Carlo

Random Walk Example



State Aggregation with Monte-Carlo

Random Walk Example

Policy
50% \longleftrightarrow 50%

$$\gamma = 1$$

$R = +1$



State Aggregation with Monte-Carlo

State Aggregation

| State | Value |
|-------|-------|
| s_1 | 3 |
| s_2 | 3 |
| s_3 | 3 |
| s_4 | 3 |
| s_5 | 0 |
| s_6 | 0 |
| s_7 | 0 |
| s_8 | 0 |

$$\left. \begin{array}{l} \left. \begin{array}{ll} \text{State} & \text{Value} \\ \hline s_1 & 3 \\ s_2 & 3 \\ s_3 & 3 \\ s_4 & 3 \\ \hline s_5 & 0 \\ s_6 & 0 \\ s_7 & 0 \\ s_8 & 0 \end{array} \right\} \mathbf{x}(s) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} & \hat{v}(s, \mathbf{w}) = w_1 \\ \left. \begin{array}{ll} s_5 & 0 \\ s_6 & 0 \\ s_7 & 0 \\ s_8 & 0 \end{array} \right\} \mathbf{x}(s) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} & \hat{v}(s, \mathbf{w}) = w_2 \end{array} \right.$$

State Aggregation with Monte-Carlo

How to Compute the Gradient for Monte Carlo with State Aggregation

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

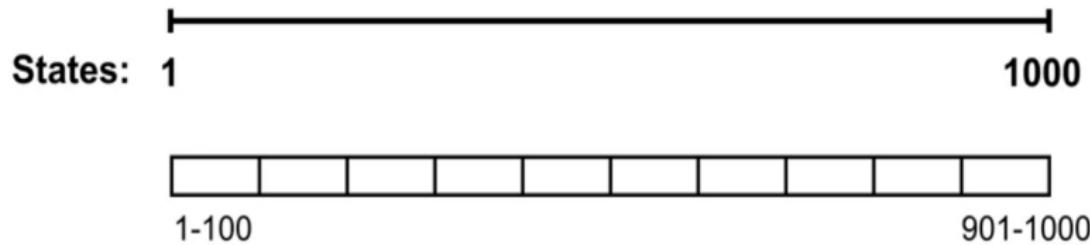
$$w_i \leftarrow w_i + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] 0$$

$$w_j \leftarrow w_j + \alpha [G_t - \hat{v}(S_t, \mathbf{w})] 1$$

$$\nabla \hat{v}(S_t, \mathbf{w}) = \mathbf{x}(S_t) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

State Aggregation with Monte-Carlo

Constructing a State Aggregation for the Random Walk



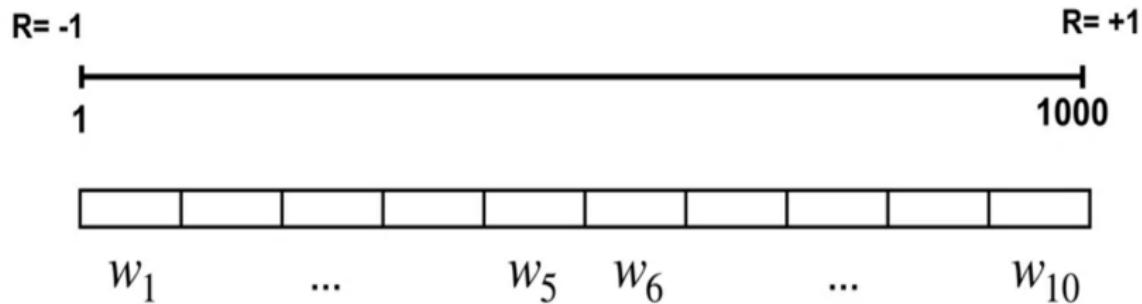
State Aggregation with Monte-Carlo

Monte Carlo Updates for a Single Episode

Return: 1 , 1 , 1 ,..., 1

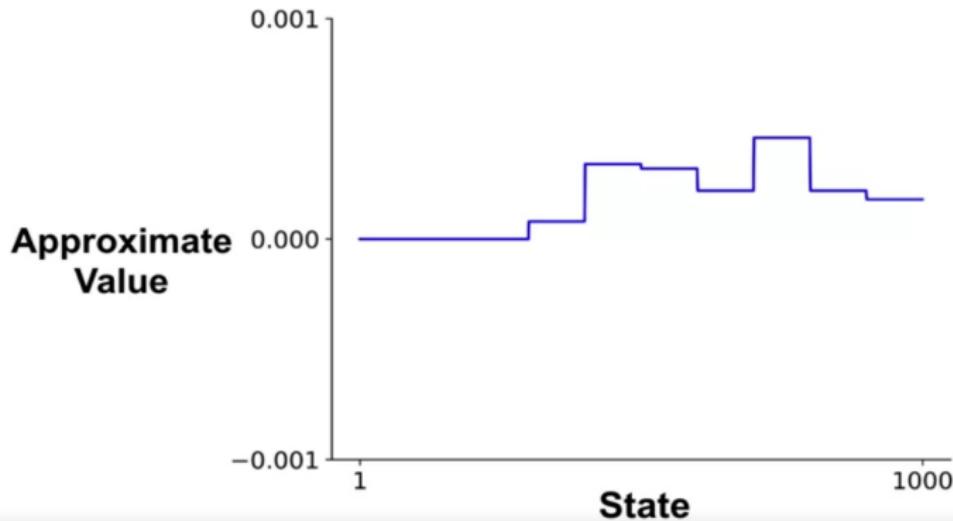
$$\alpha = 2 * 10^{-5}$$

Visited states: 500, 423, 482, ..., 936



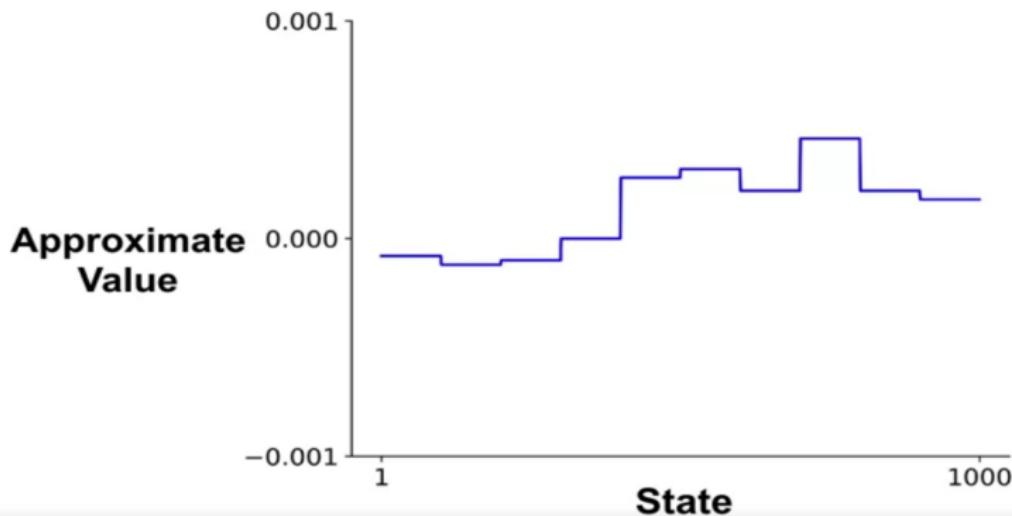
State Aggregation with Monte-Carlo

Value Estimates After One Episode



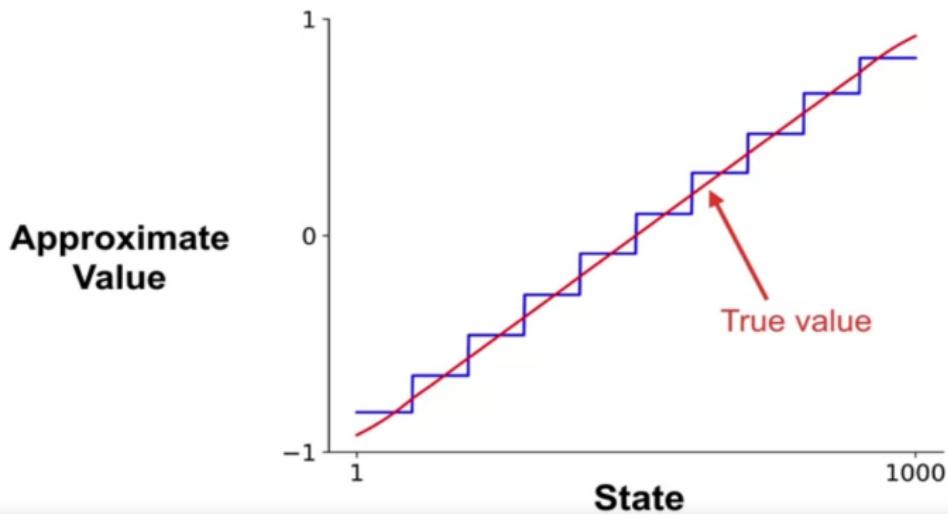
State Aggregation with Monte-Carlo

Value Estimates After Two Episodes



State Aggregation with Monte-Carlo

Final Value Estimates



Semi-Gradient TD

$$w \leftarrow w + \alpha [U_t - \hat{v}(S, w)] \triangledown \hat{v}(S, w)$$

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

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

Input: the policy π to be evaluated

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

Algorithmic parameter: step size $\alpha > 0$

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

Loop for each episode:

 Initialize S

 Loop for each episode step:

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

 Take action A , observe R, S'

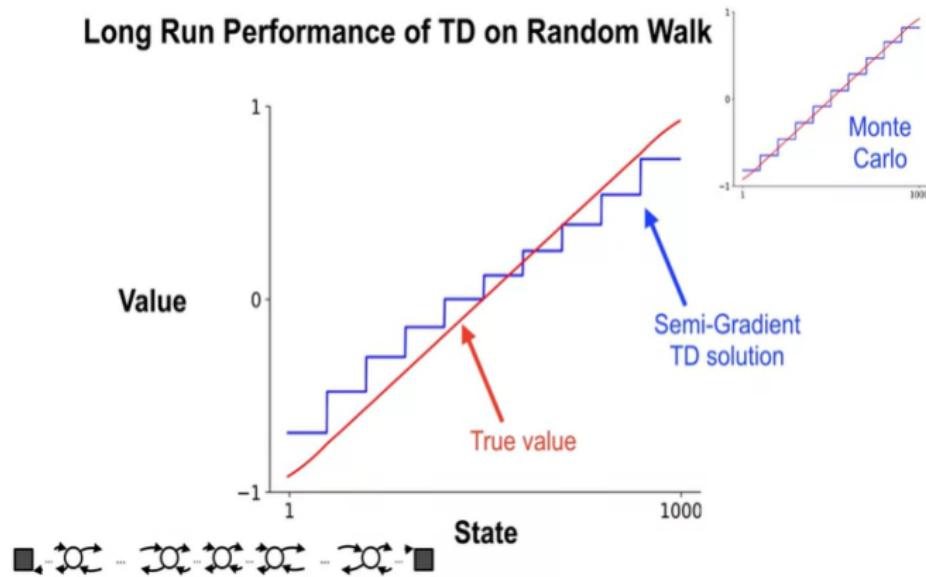
$w \leftarrow w + \alpha[R + \gamma \hat{v}(S', w) - \hat{v}(S, w)] \triangledown \hat{v}(S, w)$

$S \leftarrow S'$

 Until S is terminal

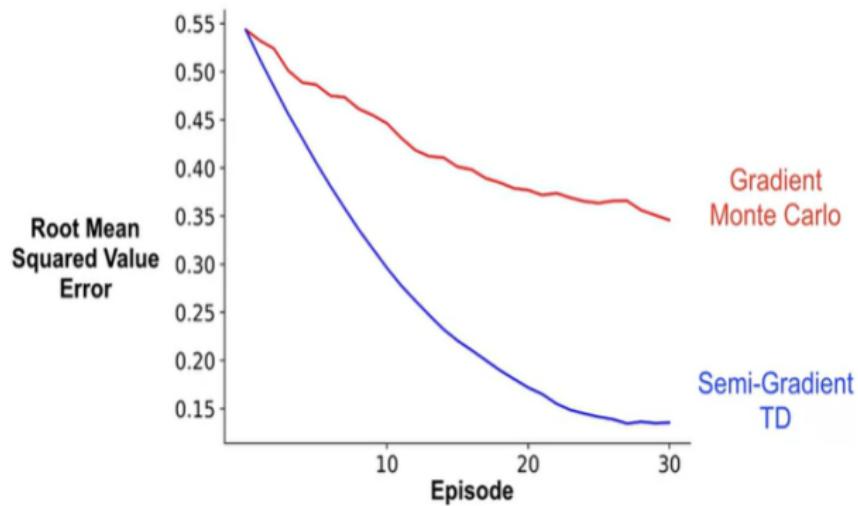
Long-run TD vs MC

Long Run Performance of TD on Random Walk



TD vs MC learning speed

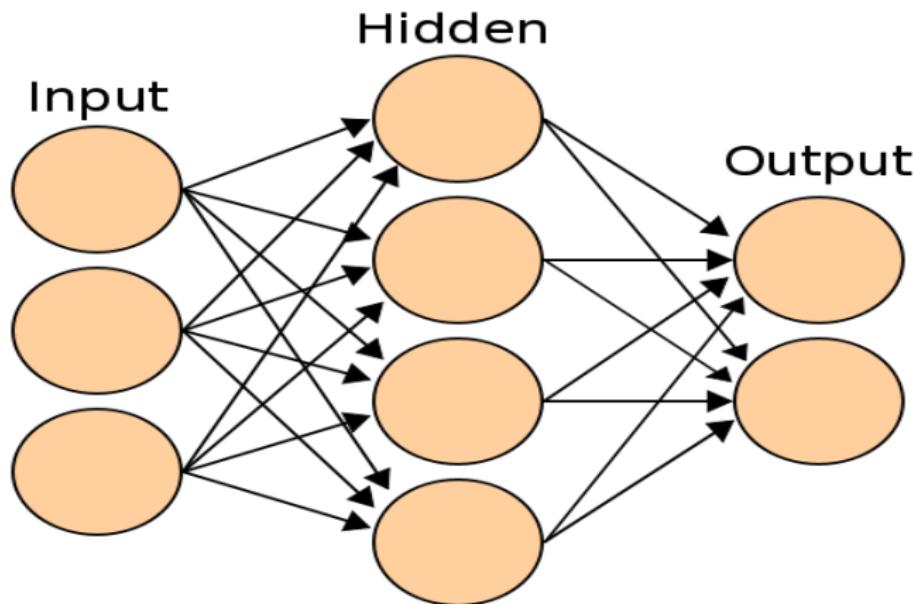
TD converges faster



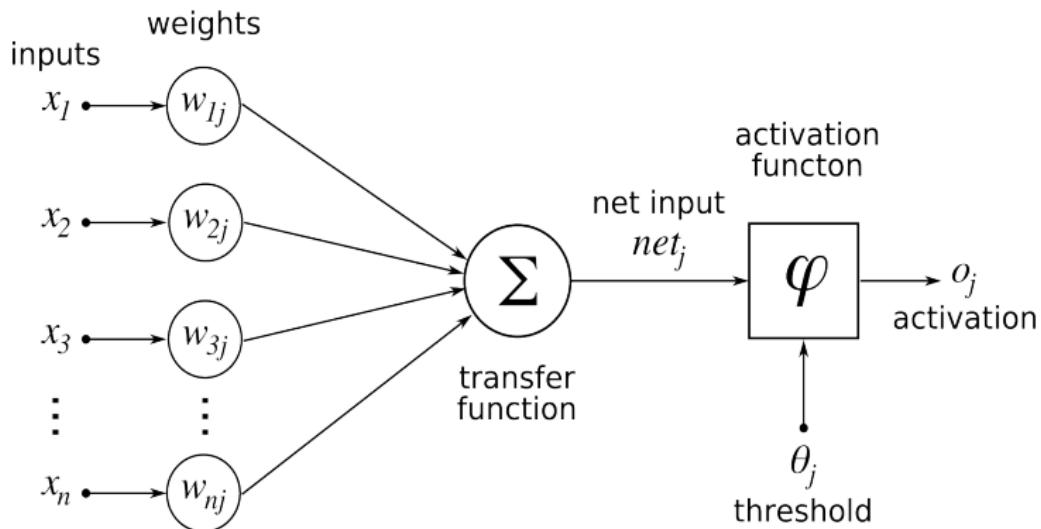
Contents

- 1 Model-based planning
- 2 Value Function Approximation
- 3 Gradient methods
- 4 Neural Networks

Neural Networks



Neural Networks



Training algorithm

- **Backpropagation**
- Evolutionary algorithms
- Pruning
- ...

Deep Q-Learning

- In Q-Learning, we stored every possible state-action value
- If the environment is too big, it becomes impossible in practice

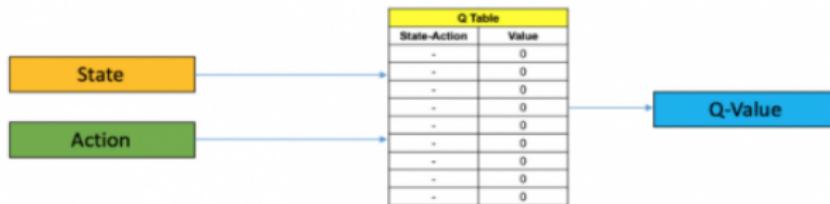
Deep Q-Learning

- In Q-Learning, we stored every possible state-action value
- If the environment is too big, it becomes impossible in practice
- In Deep Q-Learning, we want to use a neural network to approximate the Q-value of all possible action

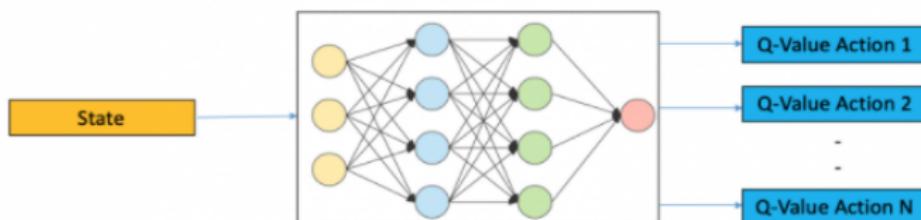
Deep Q-Learning

- In Q-Learning, we stored every possible state-action value
- If the environment is too big, it becomes impossible in practice
- In Deep Q-Learning, we want to use a neural network to approximate the Q-value of all possible action
- In a Deep Q-Network, state will be the input
- Action will be the output

Deep Q-Learning



Q Learning



Deep Q Learning