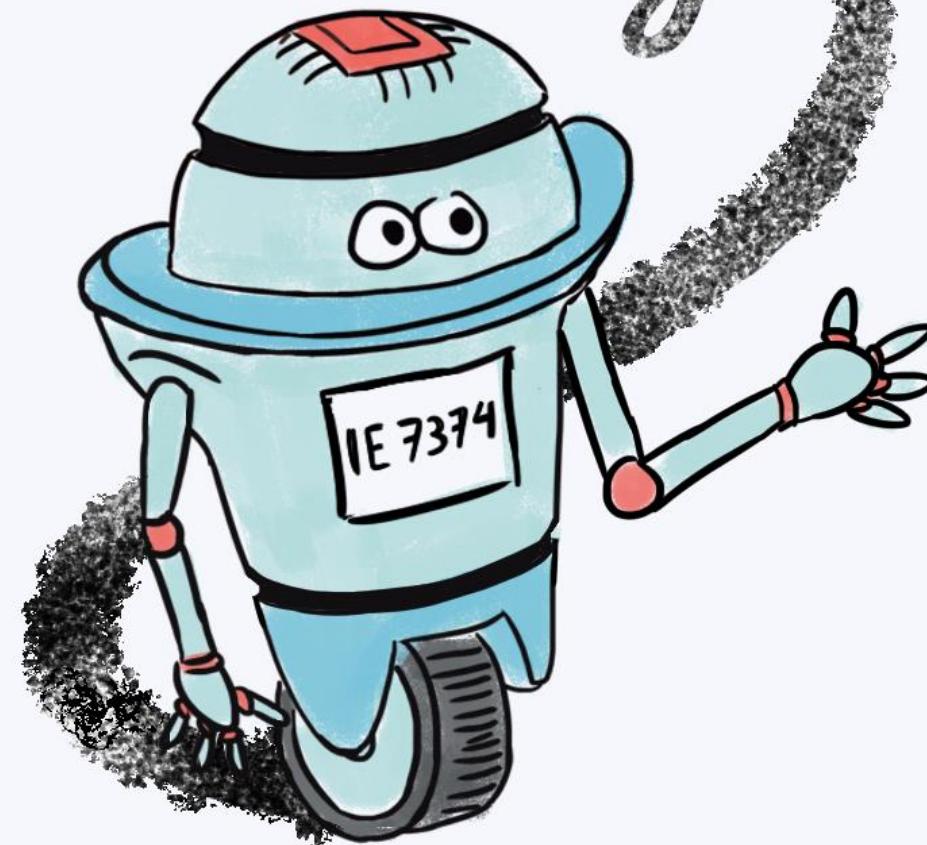
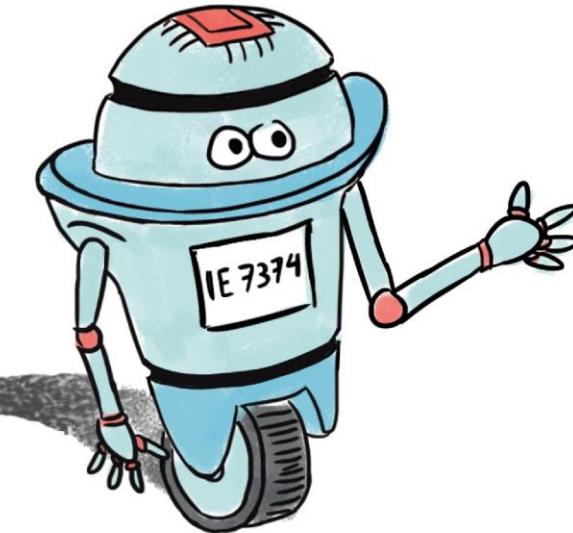


Reinforcement learning



Reinforcement learning 2



On-policy Prediction with Approximation (Part a)

Mohammad Dehghani

Mechanical and Industrial Engineering Department

Northeastern University

Apr-20



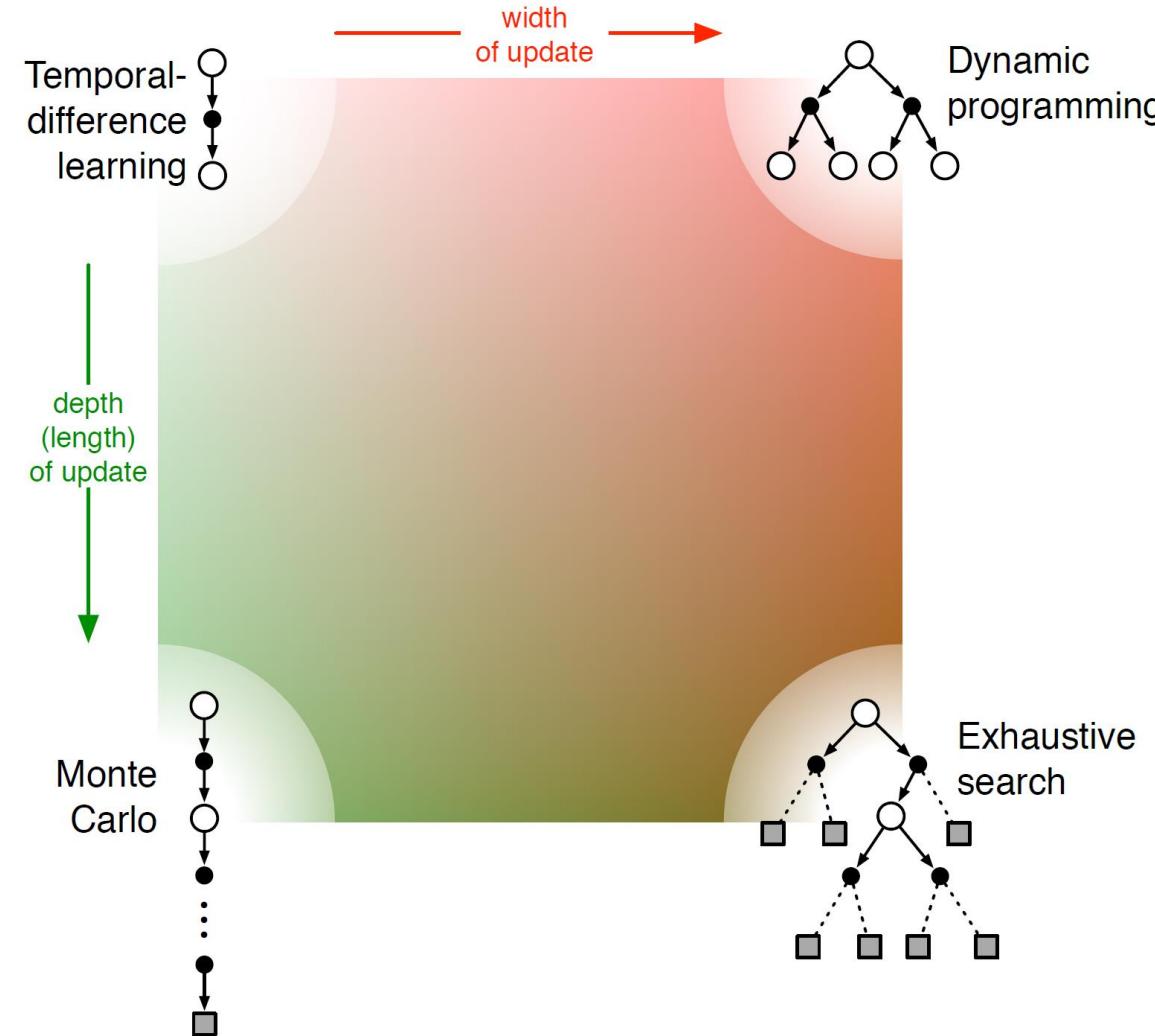
Introduction

Function approximation journey



Introduction

Overview of RL Methods



Introduction

Function approximation in reinforcement learning

- So far we have represented value function by a lookup table
 - Tabular representation
 - » Every state s has an entry $V(s)$, or
 - » Every state-action pair $s; a$ has an entry $Q(s, a)$
- In many of real world problems, the state space is combinatorial and enormous;
 - is much larger than the number of atoms in the universe.
 - » Backgammon: 10^{20} states
 - » Computer Go: 10^{170} states
 - » Helicopter: continuous state space



Introduction

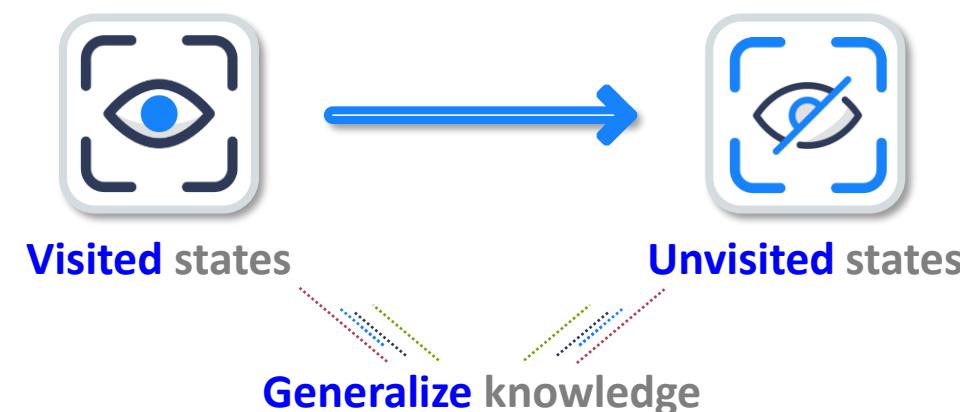
Generalization

- **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.
 - » Tabular methods become inefficient to store the values for all states in a table
 - » In fact, we can't even guarantee we will see the same state more than once.

- **Solution for large MDPs:**

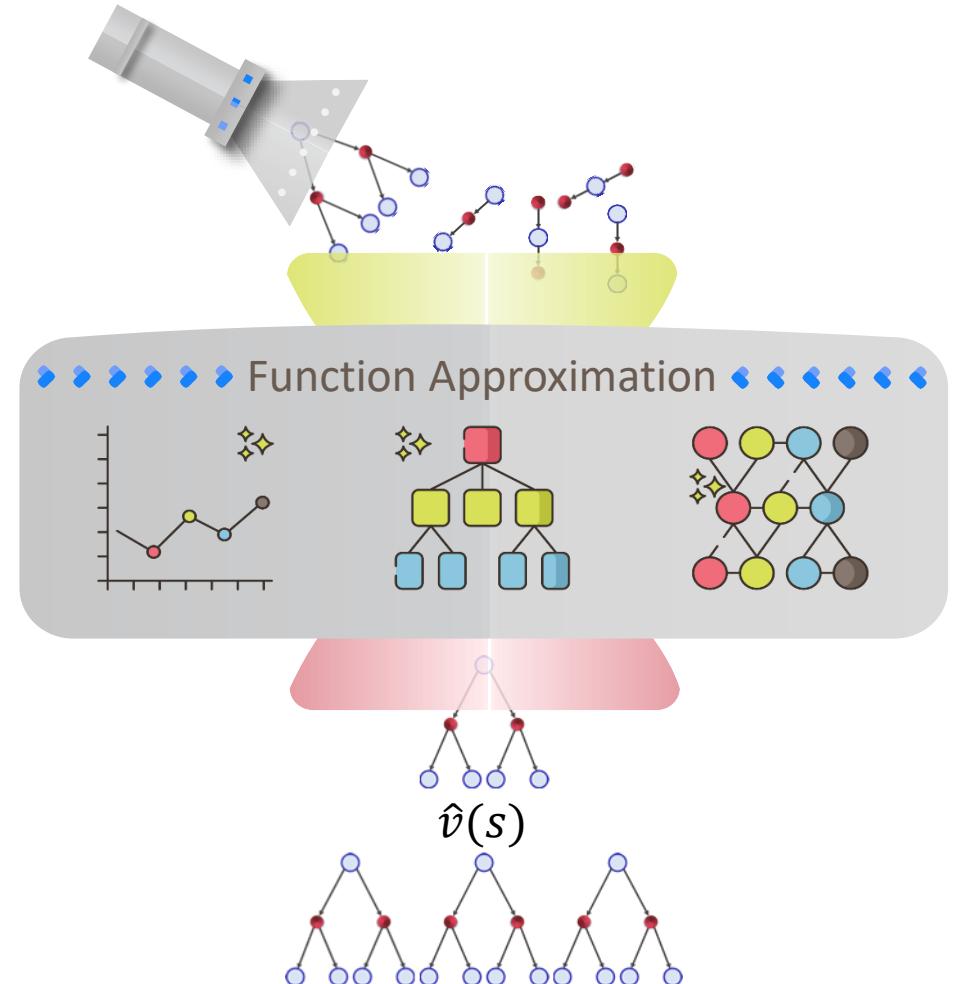
- Generalize from seen states to unseen states



Introduction

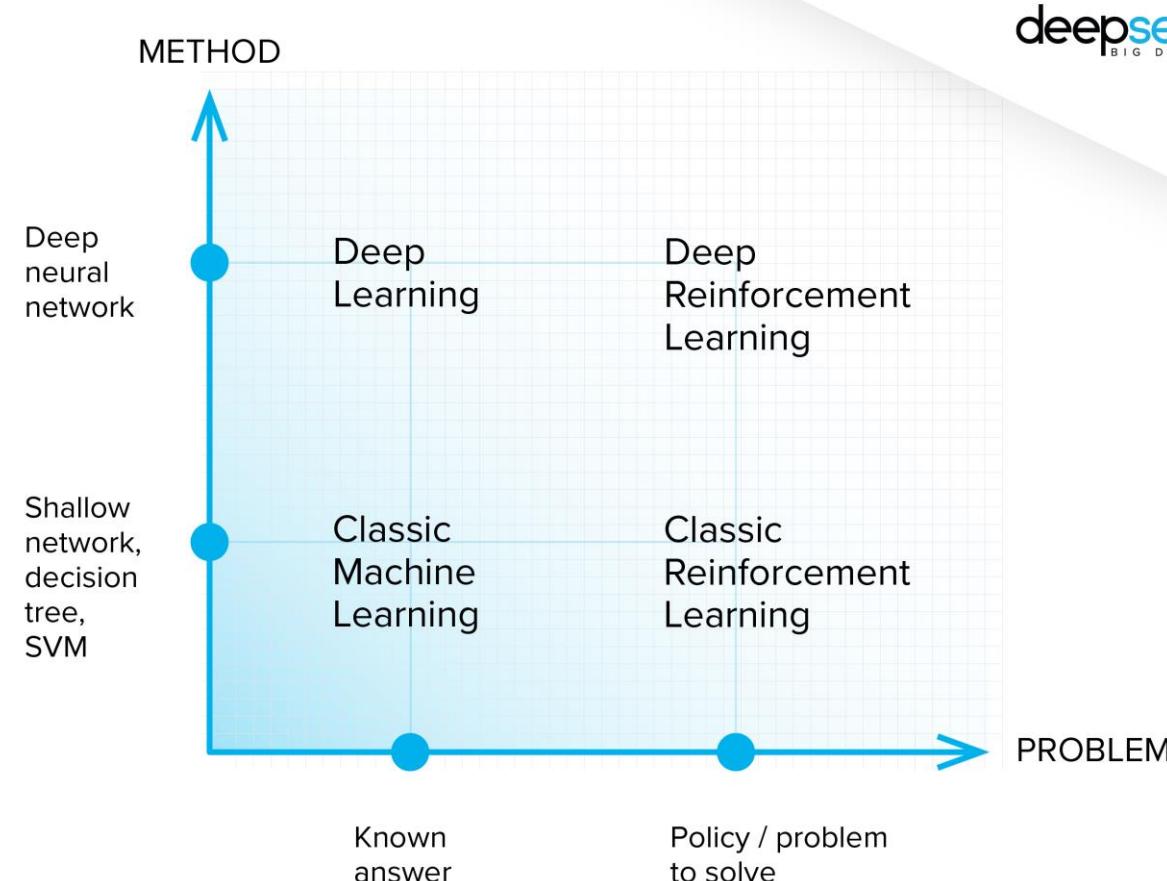
Function approximation setting

- **Goal:** Approximation of the entire function
- **Function approximation** is an instance of **supervised learning**
 - Statistical curve fitting
 - Decision tree
 - Artificial neural networks

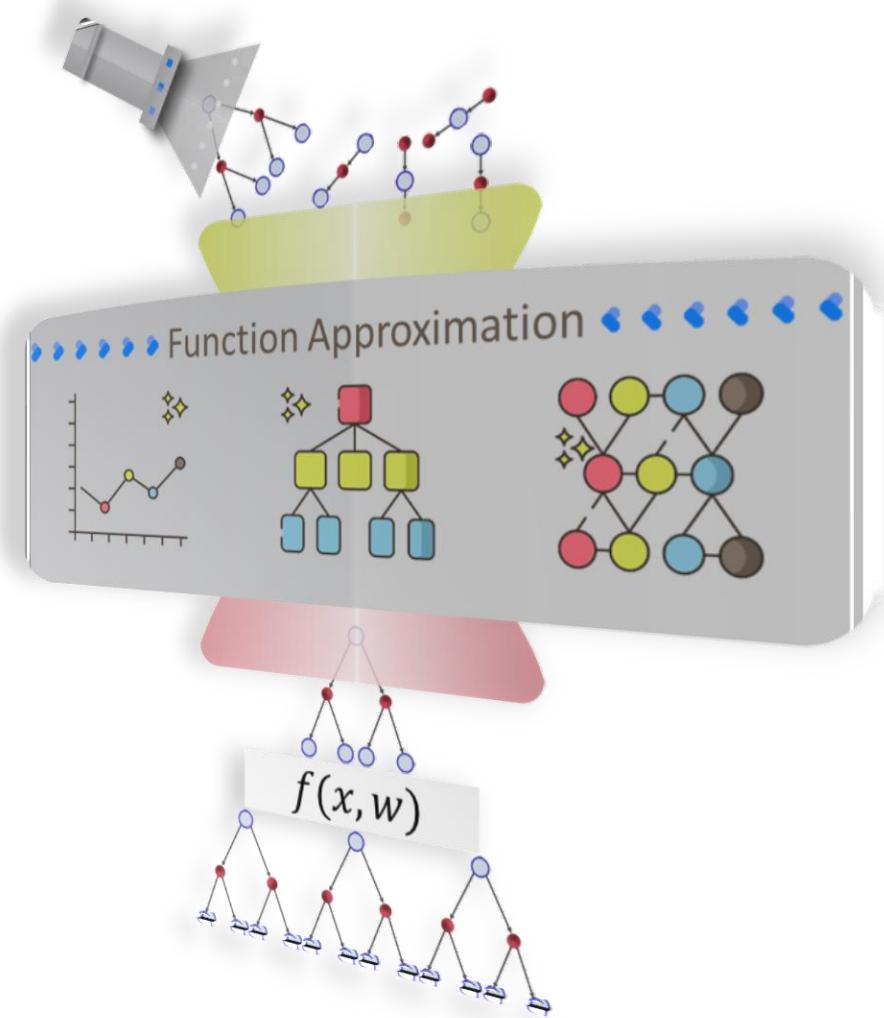


Introduction

Outline



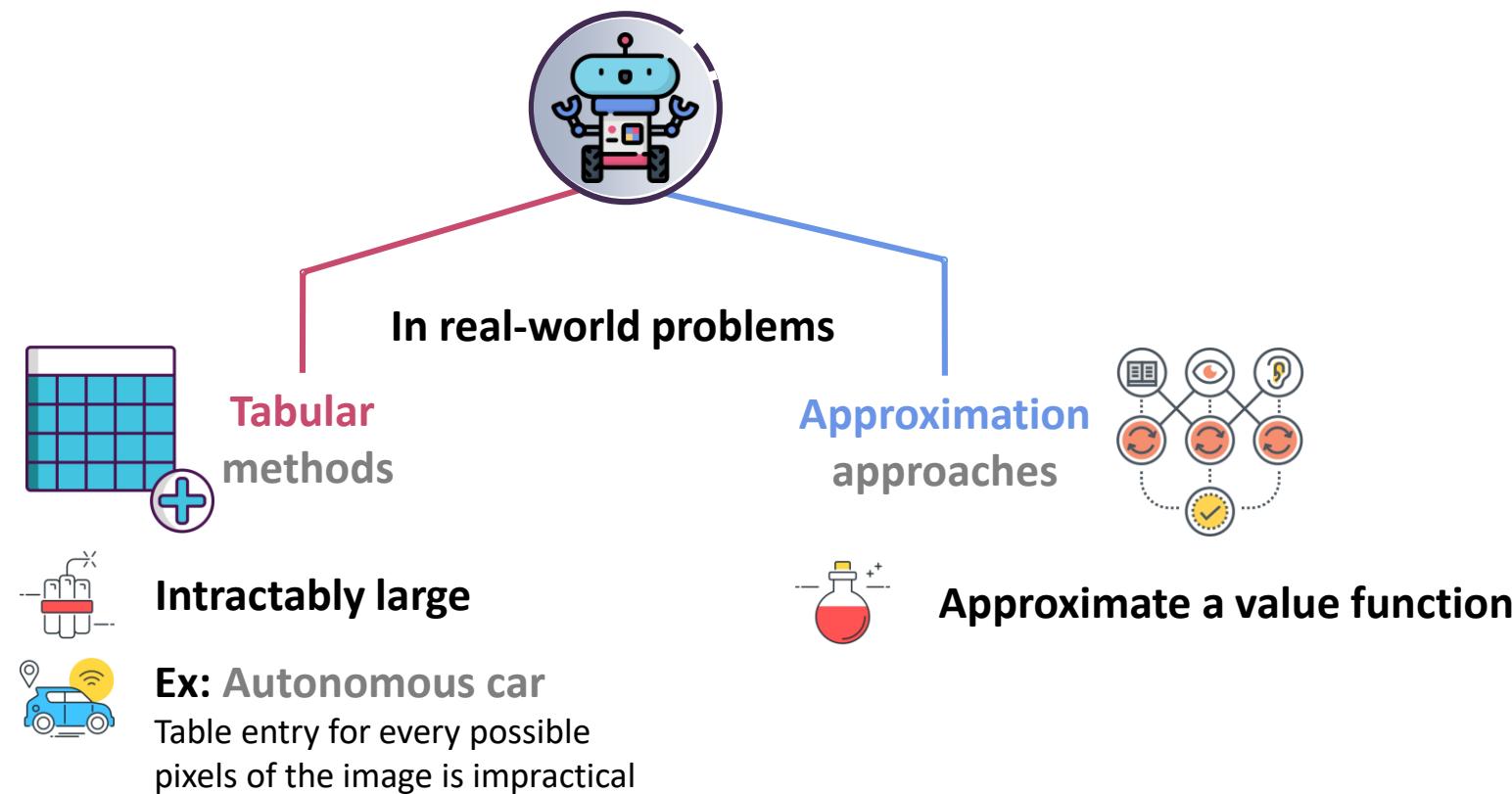
Value Function Approximation



Value Function Approximation

→ Tabular methods vs. Approximation approaches

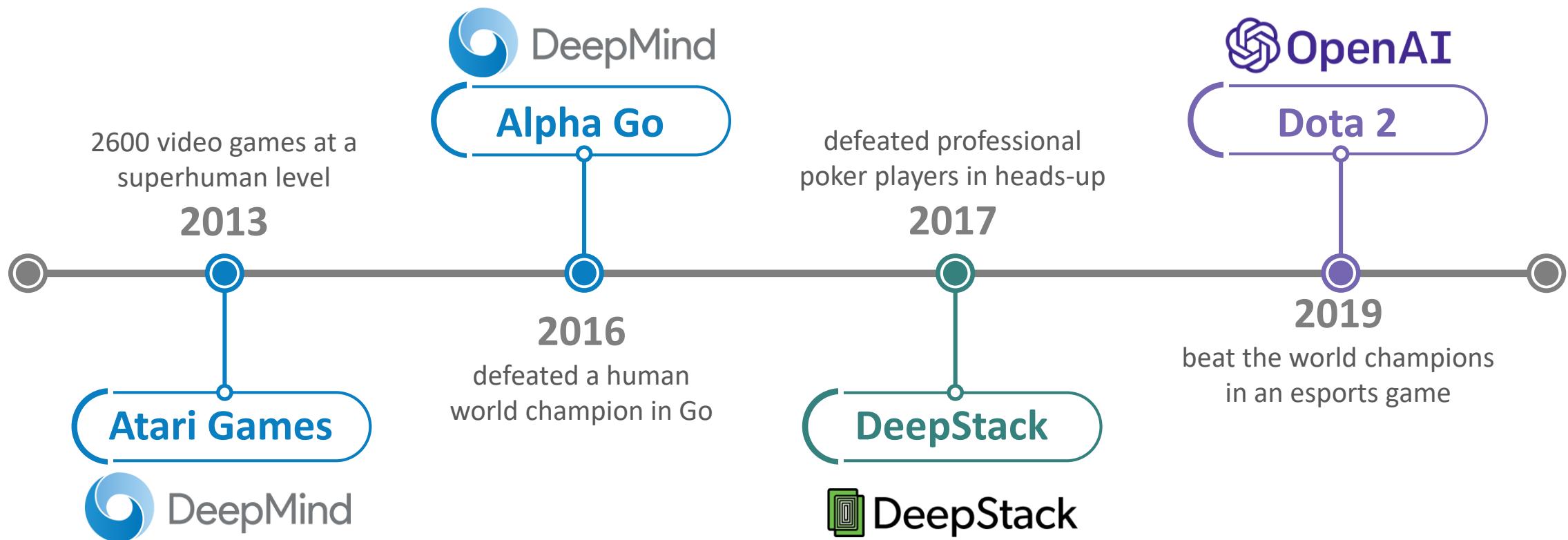
- Tabular approaches are inherently limited to *fairly low-dimensional* problems and suffer from *lack of scalability*



Value Function Approximation

Trendline in games and applications

- Tabular case is a special case of linear value function approximation



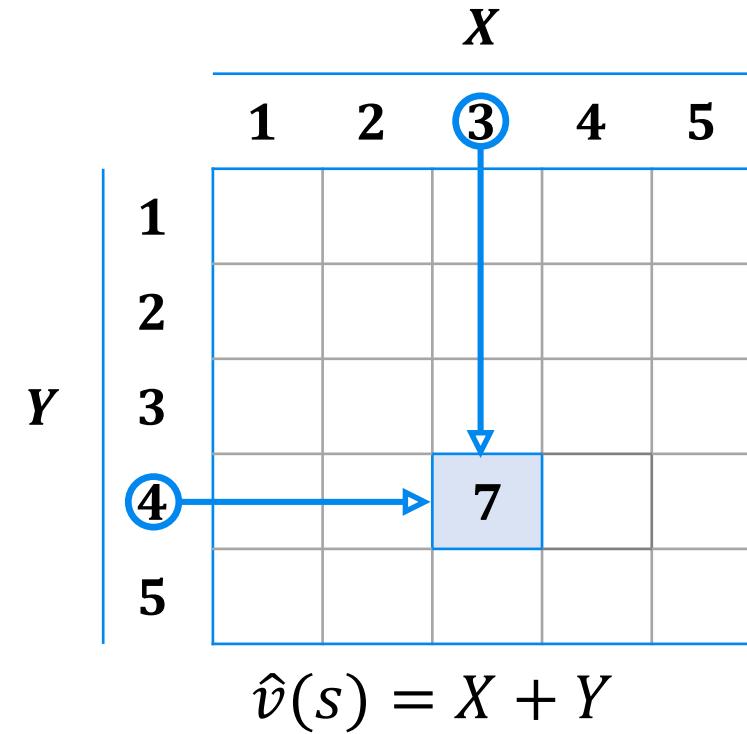
Value Function Approximation

Value function representation

Tabular
representation

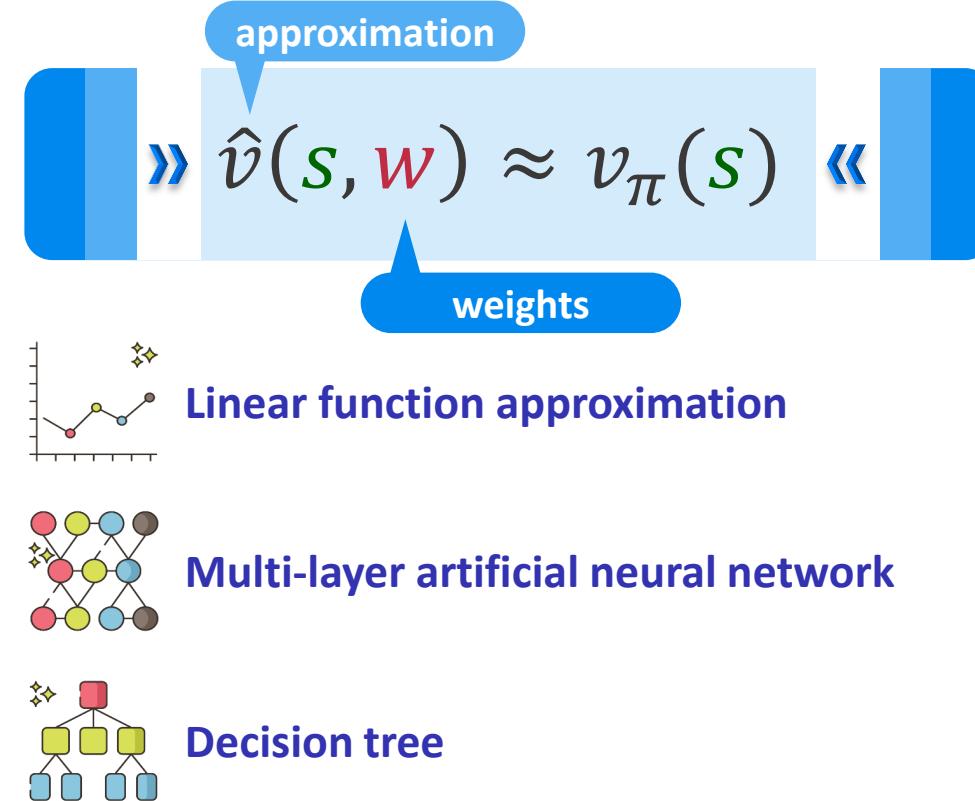
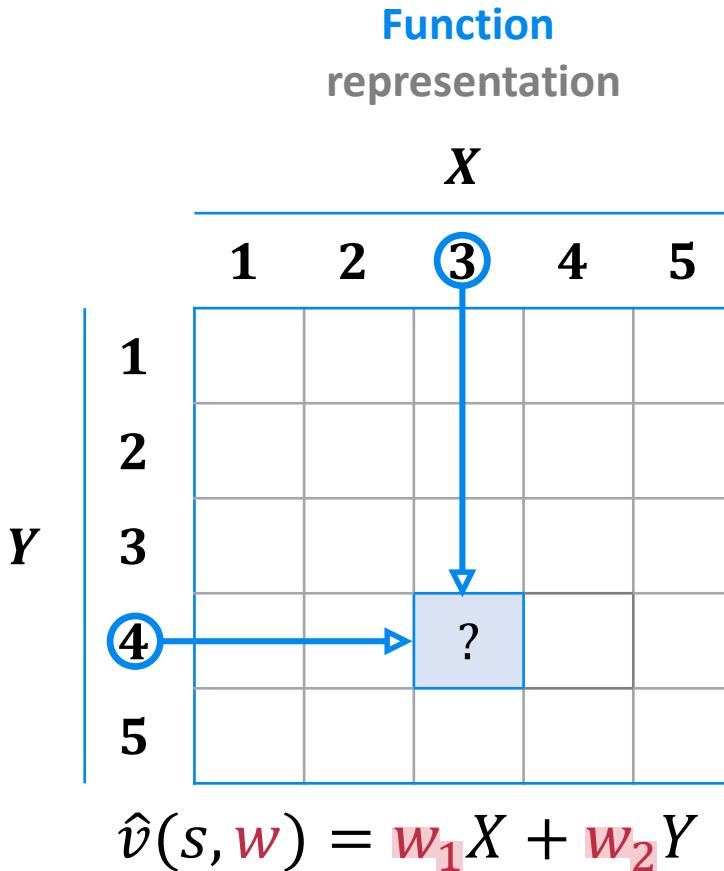
State	Value
s_1	20
s_2	0
s_3	-10
s_4	8
...	...
s_{25}	6

Function
representation



Value Function Approximation

Parametrizing the Value Function



Value Function Approximation

Characteristics

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

- ✓ Typically, the number of weights (the dimensionality of w) is much less than the number of states ($d \ll |\mathcal{S}|$)
- ✓ Changing one weight changes the estimated value of many states.
 - » Such *generalization* makes the learning potentially more powerful but also potentially more difficult to manage and understand.
- ✓ Applicable to partially observable problems

Value Function Approximation

Advantages and challenges

Advantages:

Flexible for high dimensional problems:

- » Parameterized value function is flexible enough to represent every function mapping to overcome the high dimensional problems (possibly continuous space - $\mathcal{S} \times \mathcal{A}$ to \mathbb{R})

Faster:

- » Allows the model to learn in reasonable time and space

Less memory requirements:

- » Instead of recording all state (or state-action) values in a tabular format, function approximation provides estimation using weights.
- » Weights parametrize the function, and the only values that needs to be stored.

Challenges

Setting weights

- » Our learning outcomes will modify the weights, instead of the individual state values

Tabular representation

State	Value
S_1	20
S_2	0
S_3	-10
S_4	8
...	...
S_{25}	6

Function representation

$$\hat{v}(s, w) = w_1 X + w_2 Y$$

Definition

Linear Function Approximation is a linear combination of features in order to approximate the value function.

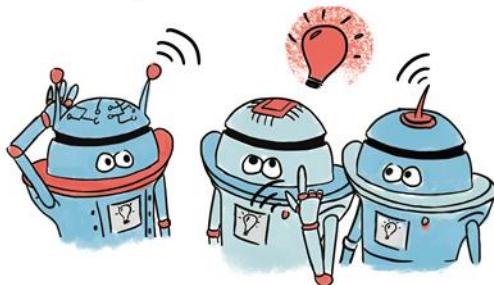
- Linear combination means that the parameters are linearly combined.

$$\widehat{v}(s, w) \doteq \sum w_i x_i(s)$$
$$= \langle w, x(s) \rangle$$

weights Features

- The features can be *the position of a robot*, angular position and *speed of an inverted pendulum*, *configurations of the stones* in a Go game, etc.

PAIR, THINK, SHARE



Tabular value functions vs Linear functions

- In the Tabular methods, state and action spaces are small enough for the approximate value functions to be represented as arrays, or tables.
- Function approximation methods such as linear functions, takes examples from a desired function (e.g., a value function) and attempts to generalize from them to construct an approximation of the entire function.

Question

- Can we consider tabular case as a *special case* of linear value function approximation?

Tabular representation

State	Value
s_1	w_1
s_2	w_2
\vdots	\vdots
s_{10}	w_{10}
\vdots	\vdots
s_{25}	w_{25}

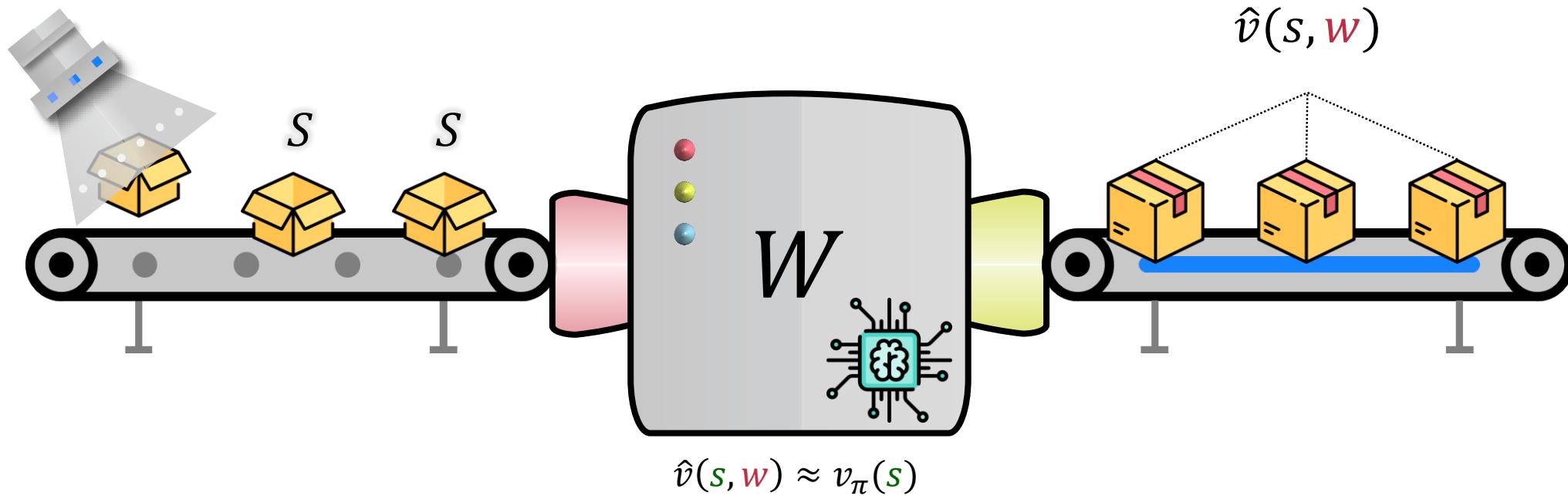
$$x(s_i) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Function representation

$$\hat{v}(s_i, \mathbf{w}) \doteq \langle \mathbf{w}, x(s_i) \rangle = w_i$$

Value Function Approximation

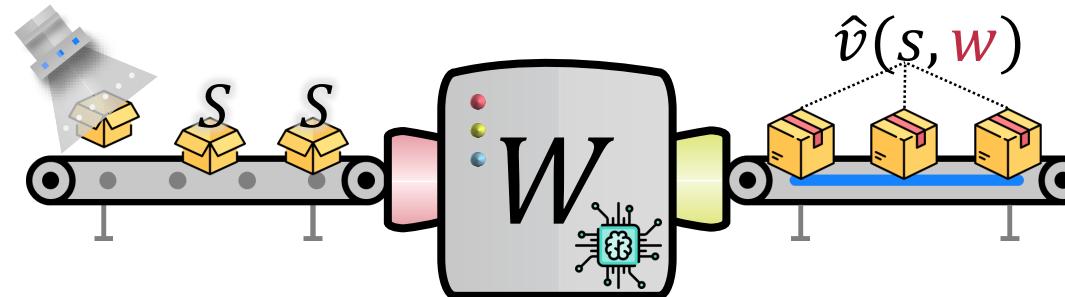
Function Approximation methods: General setting



Value Function Approximation

Applicability of supervised Learning methods

- Estimating a function from samples is addressed by the *supervised learning* paradigm



➢ ➢ ➢ ➢ ➢ ➢ RL Updates ➢ ➢ ➢ ➢ ➢ ➢

Monte Carlo

$$S_t \mapsto G_t$$

TD(0)

$$S_t \mapsto R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t)$$

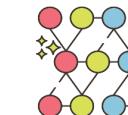
DP

$$S \mapsto \mathbb{E}[R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t) | S_t = s]$$

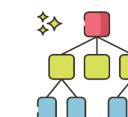
➢ ➢ ➢ ➢ Function Approximation ➢ ➢ ➢ ➢



Statistical curve fitting



Artificial neural networks

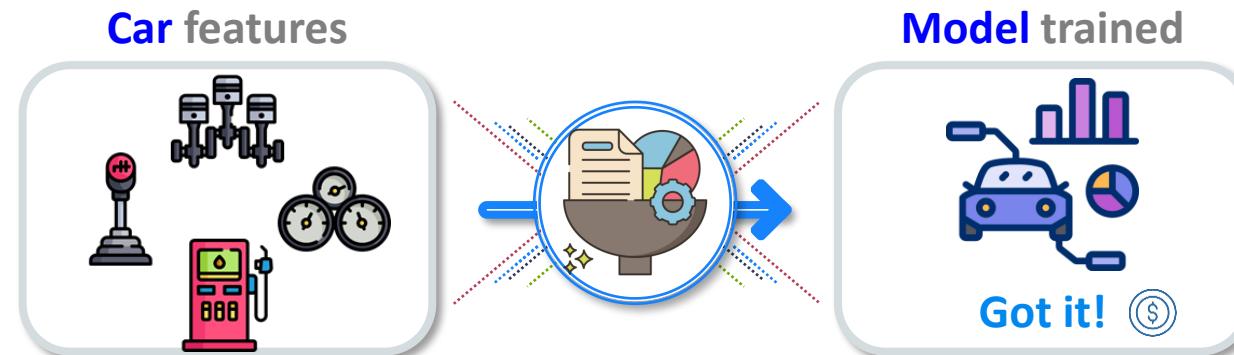


Decision tree

Value Function Approximation

Challenges with using supervised Learning methods (1)

- Estimating a function from samples is addressed by the *supervised learning* paradigm
 - Ex: Car price estimation



- In reinforcement learning, the (state-action) values are *not* directly observed
 - Indeed, values are never directly observed, just rewards which define them.
 - This renders the underlying estimation problem more difficult.

Value Function Approximation

Challenges with using supervised Learning methods (2)

- The most sophisticated artificial neural network and statistical methods all assume a *static* training set over which multiple passes are made
 - Ex: Car price estimation



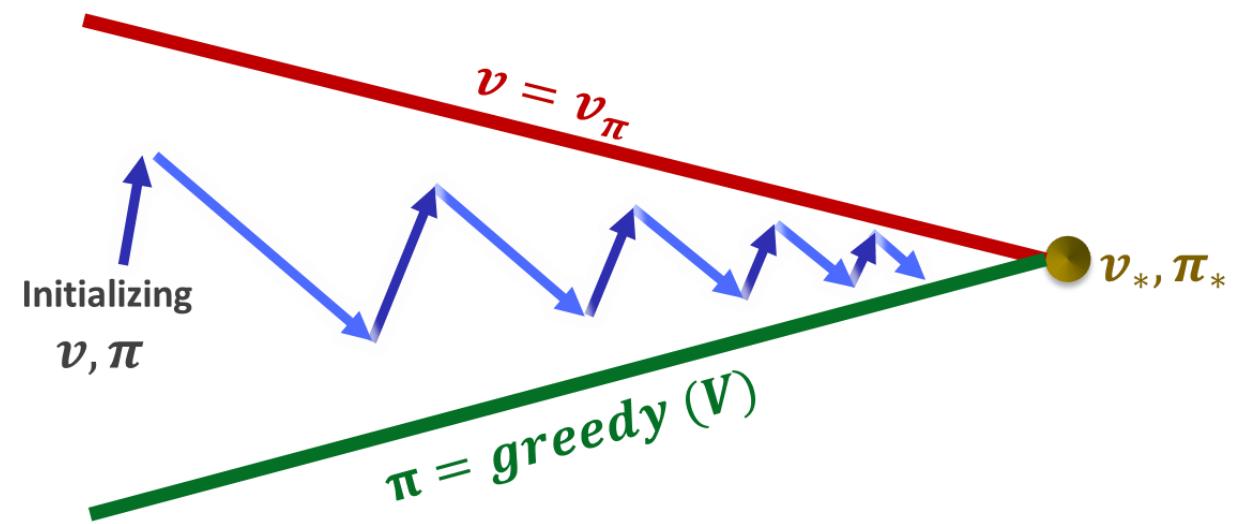
- In reinforcement learning, however, it is important that learning be able to occur online
 - While the agent interacts with its environment or with a model of its environment.
- Supervised learning methods are typically not designed for changing targets

Value Function Approximation

Challenges with using supervised Learning methods (3)

- Handle nonstationary:
 - In addition, reinforcement learning generally requires function approximation methods able to handle **nonstationary** target functions (target functions that change over time)

- Ex: Control methods based on GPI



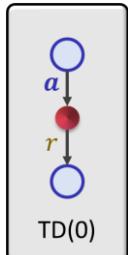
- Methods that cannot easily handle such nonstationarity are less suitable for reinforcement learning.

Value Function Approximation

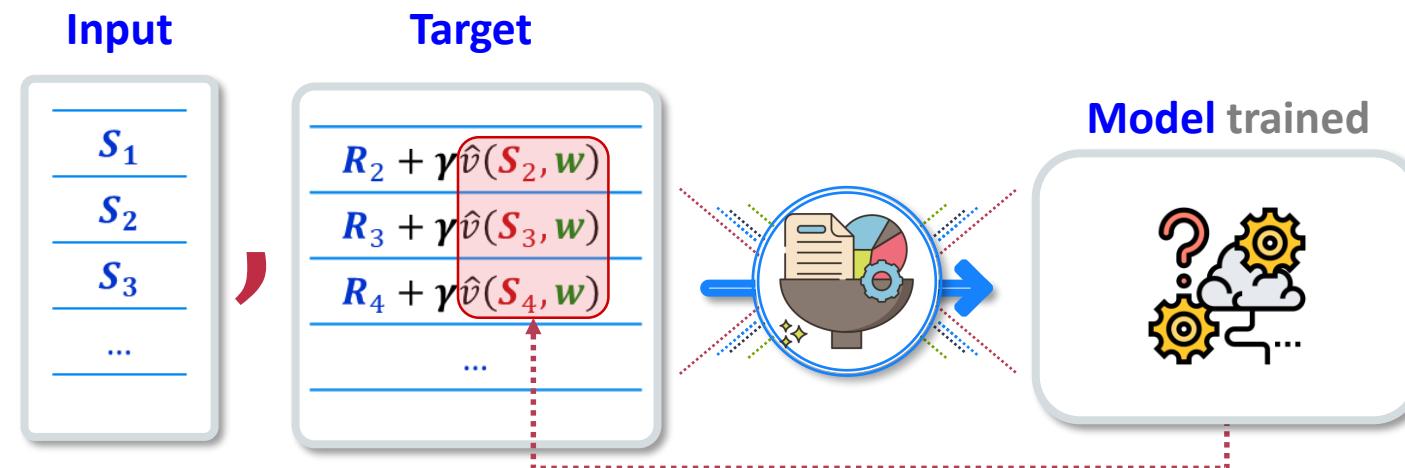
Challenges with using supervised Learning methods (4)

- In bootstrapping methods (i.e TD), predictions are used as targets during the course of learning.
 - Agent learns state values based on trained estimates for **subsequent** states

$$V(S_1) \leftarrow V(S_1) + \alpha[R_2 + \gamma V(S_2) - V(S_1)]$$



- In function approximation, the value of the **subsequent** state is also an approximation by itself.

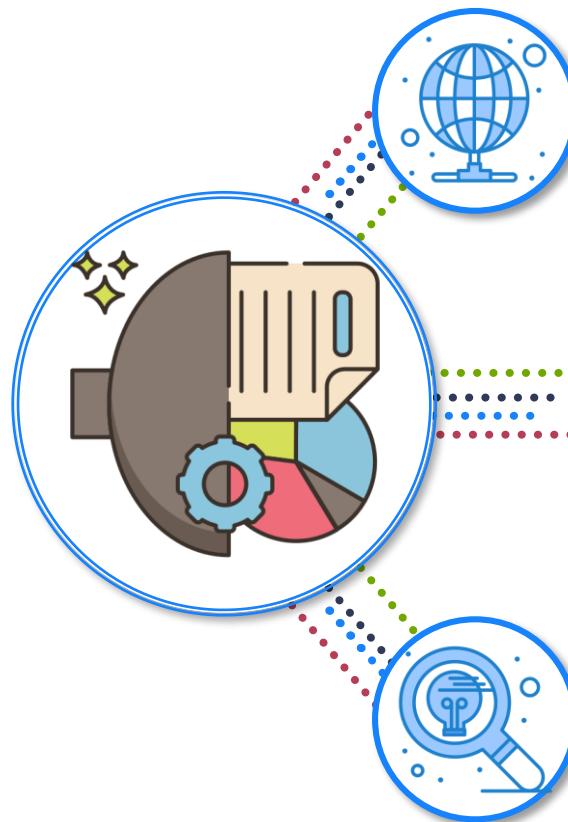


- Supervised learning methods are typically not designed for targets computed from the agent's own estimates.

Value Function Approximation

→ Summary: Appropriate supervised learnings methods for RL

- **Suitable function approximation methods for RL should be:**



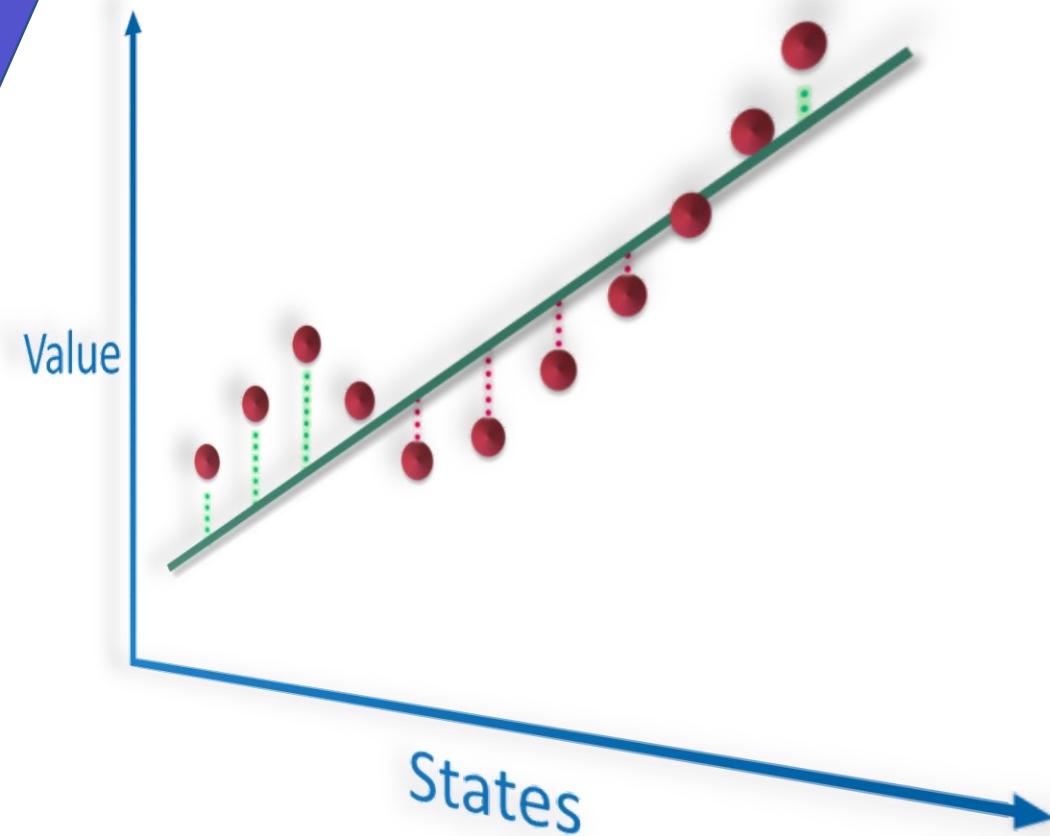
Adaptable to online learning

Capable of handling nonstationary target functions

Compatible with bootstrapping nature of some RL methods

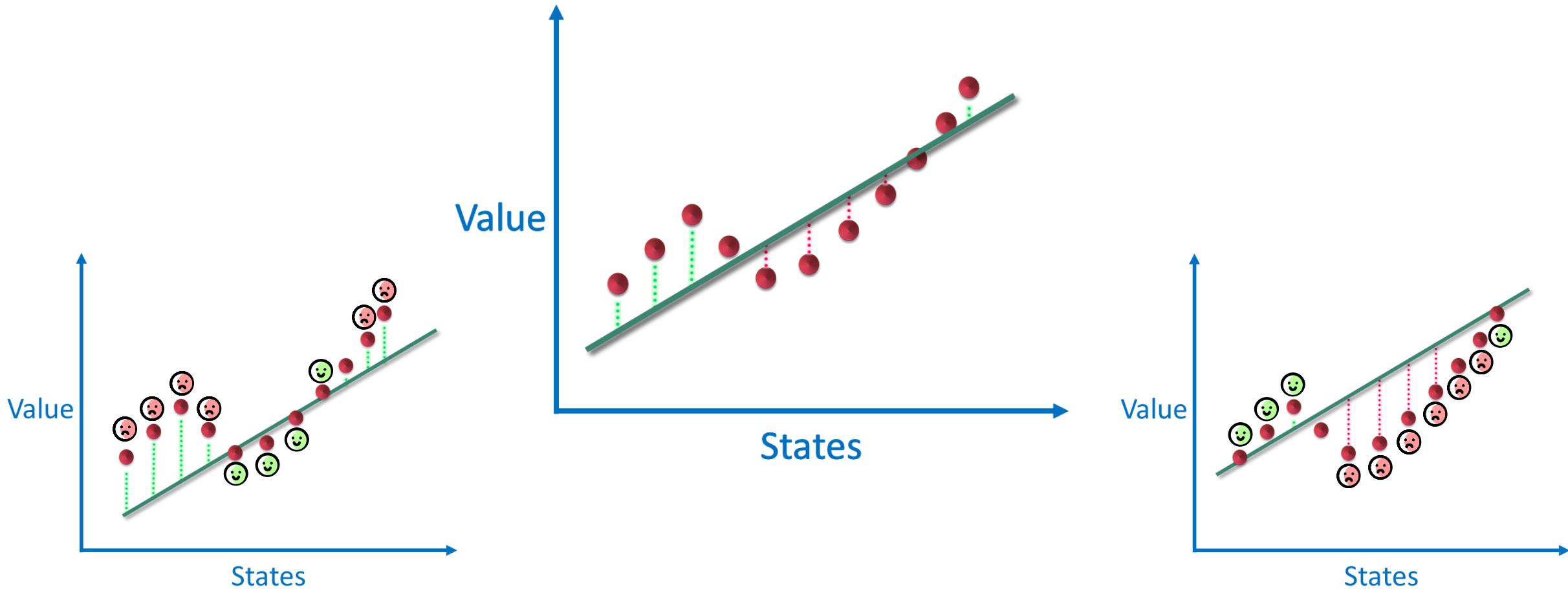
Value Error Objective

(\overline{VE})



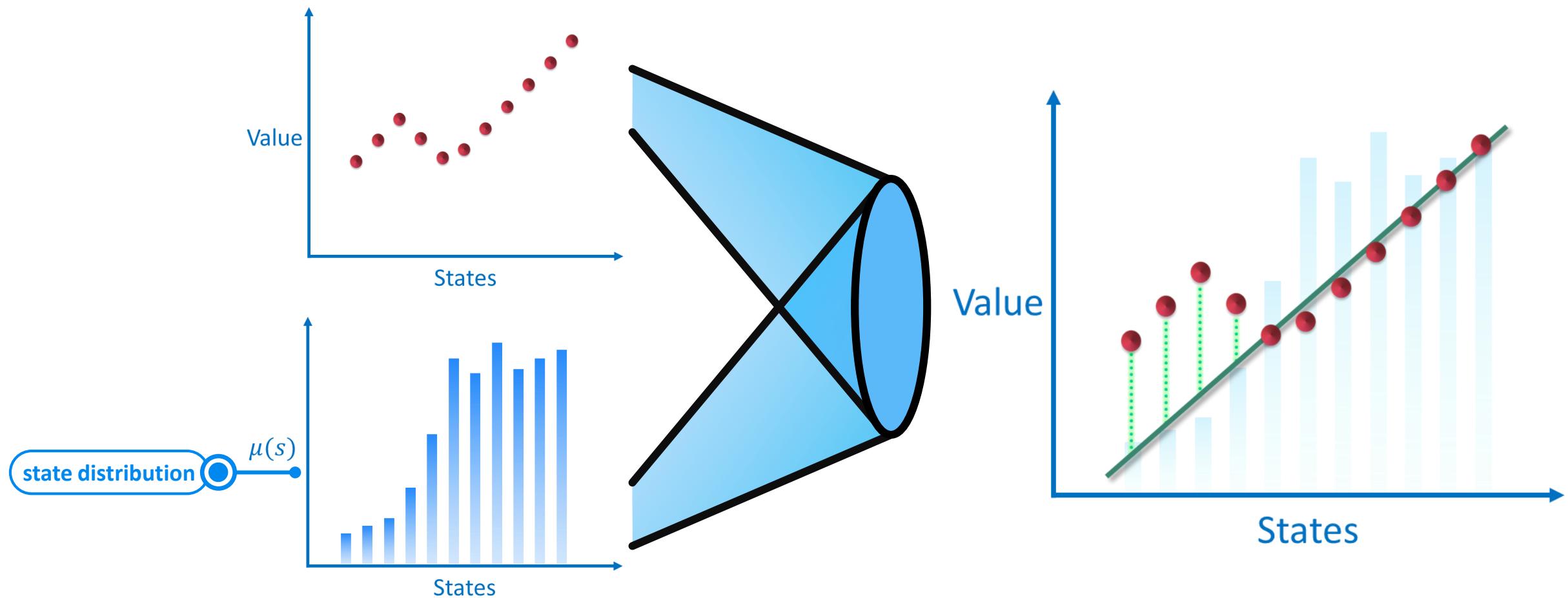
Value Error Objective

Graphical representation



Value Error Objective

Weighting approximation function



Value Error Objective

↳ Mean Squared Value Error

Definition

Mean Squared Value Error, denoted \overline{VE} :

$$\overline{VE}(w) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, w)]^2$$

true value approximate value
State dist.

- $\mu(s)$ representing how much we care about the error in each state s
 - » $\mu(s) \geq 0$
 - » $\sum_s \mu(s) = 1$

Value Error Objective

↳ The on-policy distribution

Definition

- **On-policy distribution** represents how much we care about the error in each state s , and often is chosen to be the *fraction of time spent in s* .
- In episodic tasks, *on-policy distribution* is defined as:

normalized to sum to 1

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

» $\eta(s)$ denote the number of time steps spent, on average, in state s in a single episode.

Prob. of starting from s

Transitions to s form \bar{s}

$$\eta(s) = h(s) + \sum_s h(s) \sum_s \pi(a|\bar{s}) p(s|\bar{s}, a), \quad \text{for all } s \in \mathcal{S}$$

» $h(s)$ denote the probability that an episode begins in each state s

Value Error Objective

↳ Is \overline{VE} the right performance objective for RL?

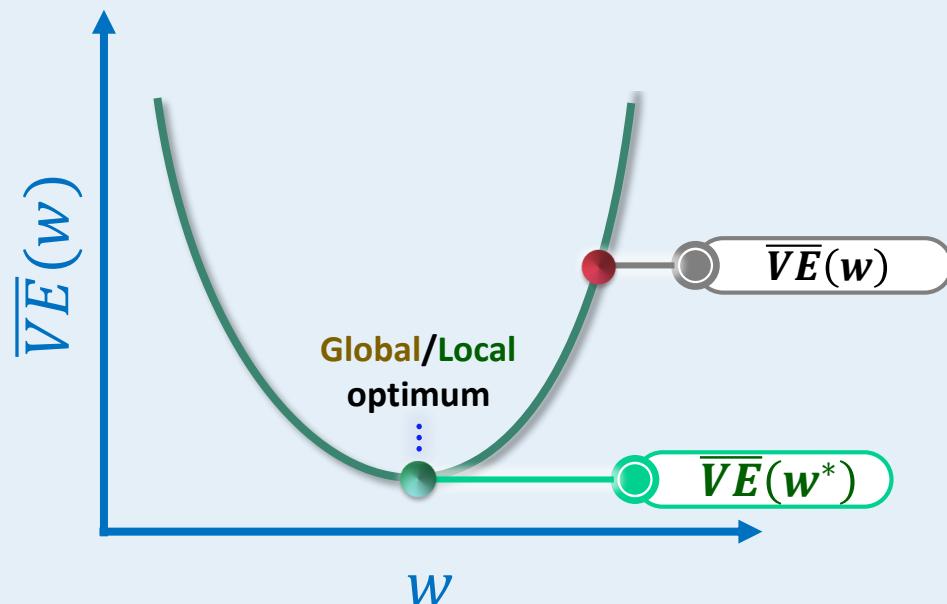
$$\overline{VE}(w) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, w)]^2$$

- **Ultimate purpose—the reason we are learning a value function—is to find a better policy.**
 - The best value function for this purpose is not necessarily the best for minimizing \overline{VE}
 - Nevertheless, it is not yet clear what a more useful alternative goal for value prediction might be
- **BTW, lets focus on other aspects of \overline{VE}** 😊

Value Error Objective

Finding the global optimum

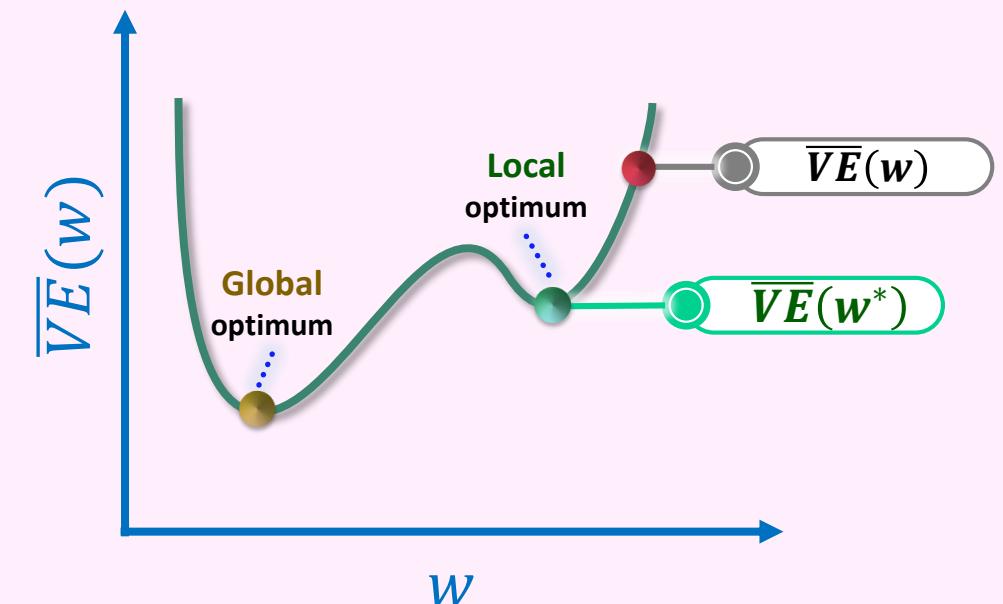
Simple function approximators



$$\overline{VE}(w^*) \leq \overline{VE}(w)$$

for all possible w

Complex function approximators

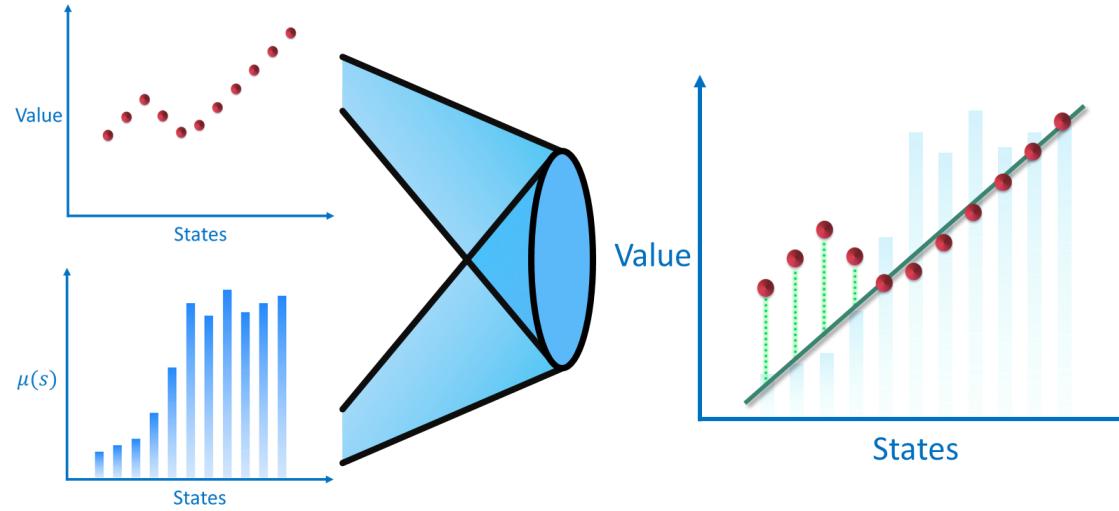


$$\overline{VE}(w^*) \leq \overline{VE}(w)$$

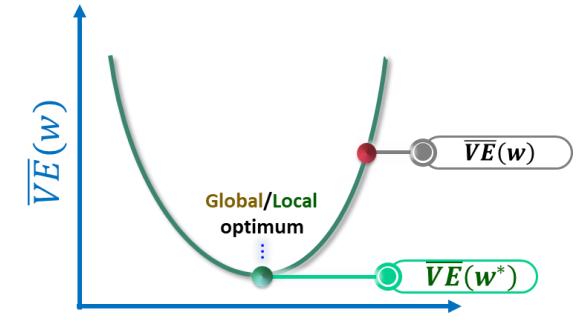
for all w in some neighborhood of w^*

Value Error Objective

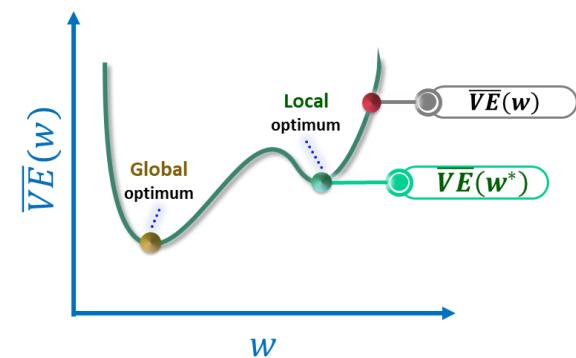
Summary



$$\overline{VE}(w) \doteq \sum_{s \in \mathcal{S}} \mu(s) [v_\pi(s) - \hat{v}(s, w)]^2$$

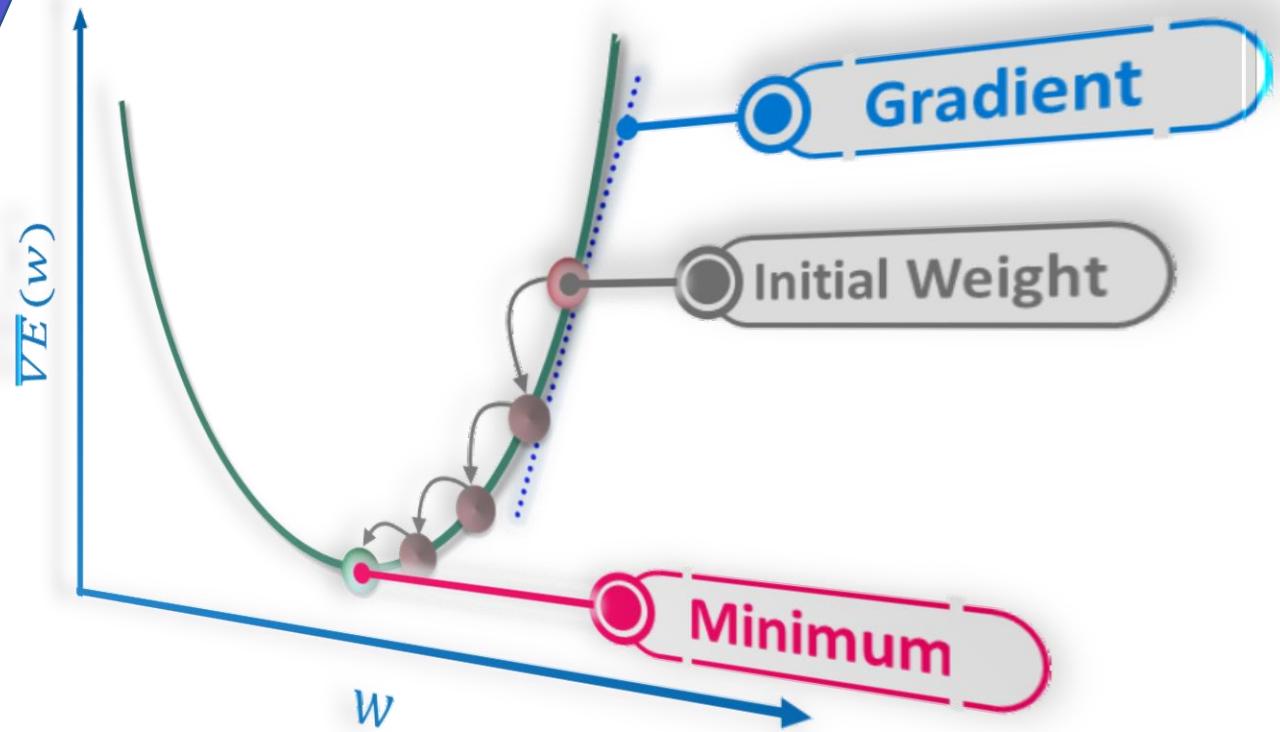


$$\overline{VE}(w^*) \leq \overline{VE}(w) \text{ for all possible } w$$



$$\overline{VE}(w^*) \leq \overline{VE}(w) \text{ for all } w \text{ in some neighborhood of } w^*$$

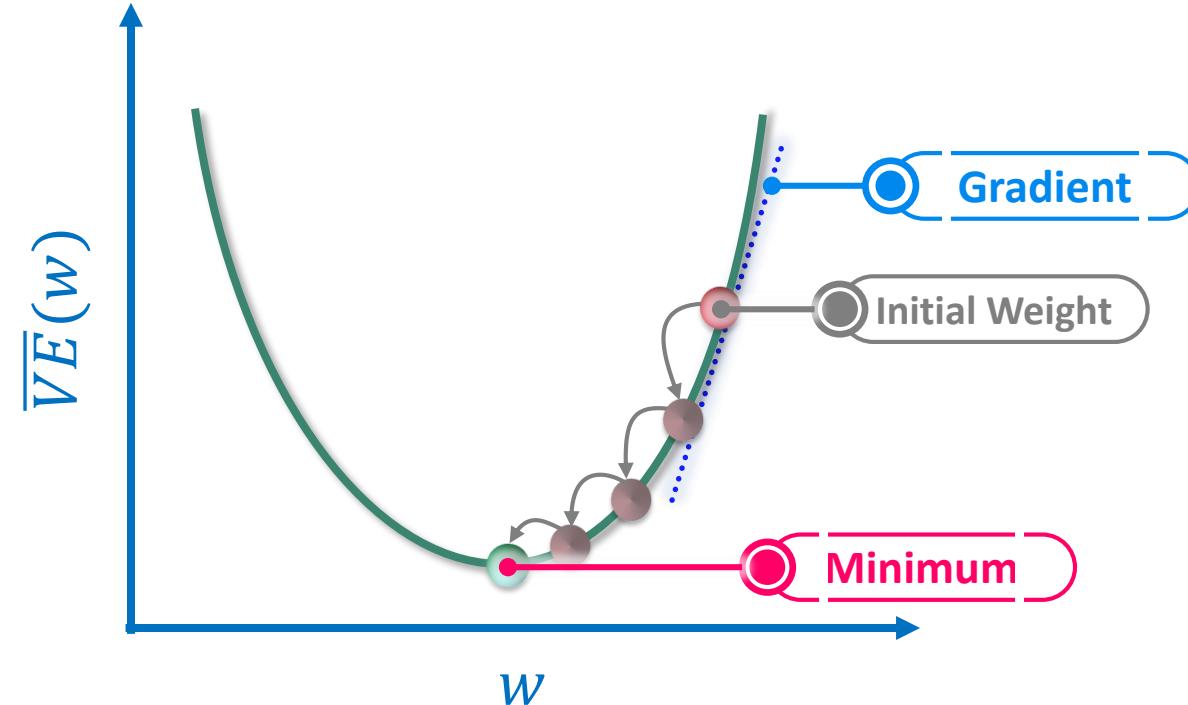
Gradient Descent (GD)



Gradient Descent

Graphical representation

- Stochastic gradient descent (SGD) is one of the most popular and used optimizers in Data Science.



- SGD methods are among the most widely used of all function approximation methods
- SGDs are particularly well suited to online reinforcement learning.

Gradient Descent

↳ Mathematical representation

$$\gg \overline{VE}(w) \doteq \sum_{s \in S} \mu(s) [\mathbf{v}_\pi(s) - \hat{\mathbf{v}}(s, w)]^2 \ll$$

$w \doteq (w_1, w_2, \dots, w_d)^T$

- How to apply gradient decent

$$w \doteq (w_1, w_2, \dots, w_d)^T$$

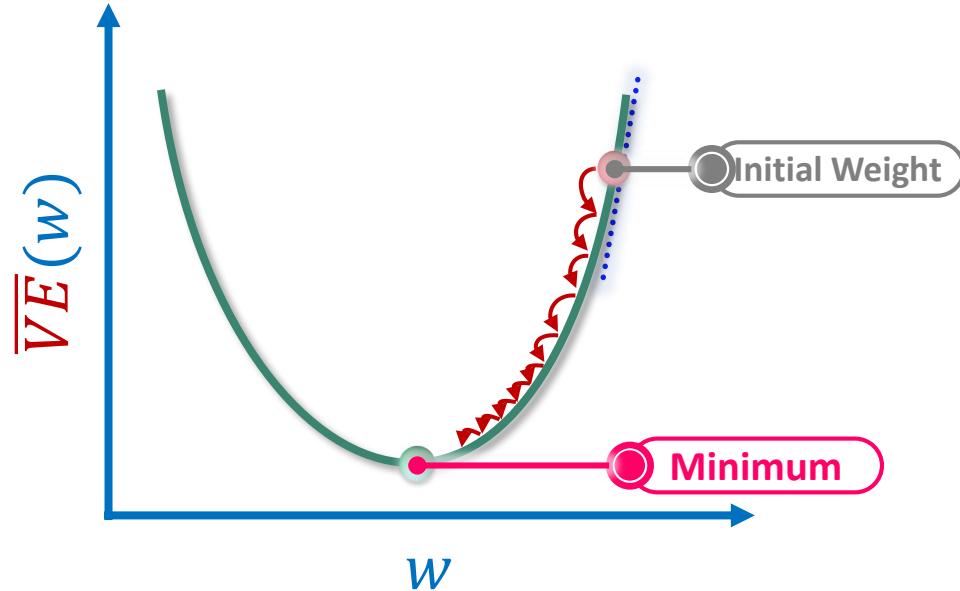
↓

$$\nabla f(w) \doteq \left(\frac{\partial f(w)}{\partial w_1}, \frac{\partial f(w)}{\partial w_2}, \dots, \frac{\partial f(w)}{\partial w_d} \right)^T$$

Gradient Descent

Importance of step size

$$w_{t+1} \doteq w_t + \alpha \nabla \overline{VE}(w_t)$$

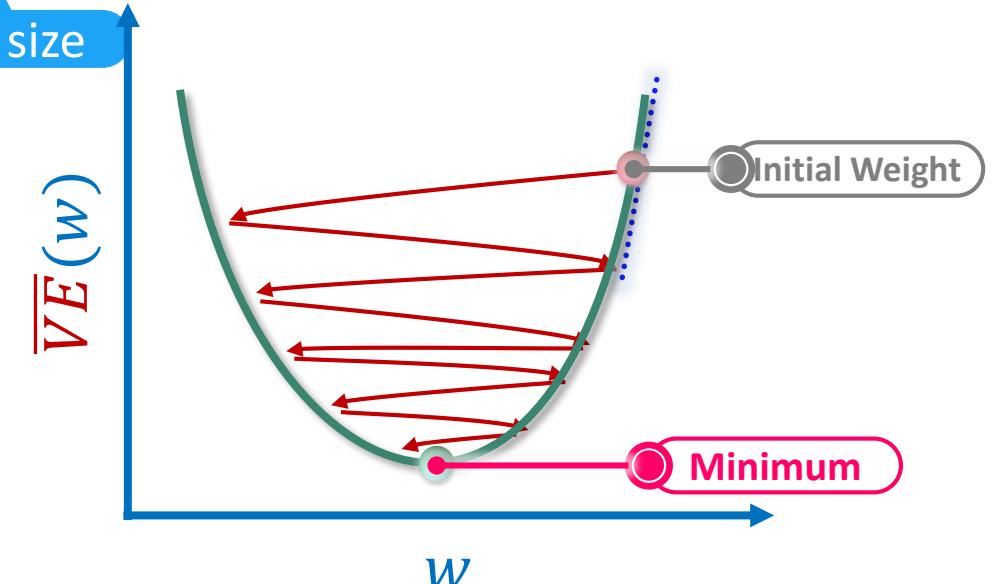


Small step size:

- ▶ gradient decent can be slow

$$\textcircled{1} \quad \sum_{n=1}^{\infty} \alpha_n(a) = \infty$$

and



Large step size:

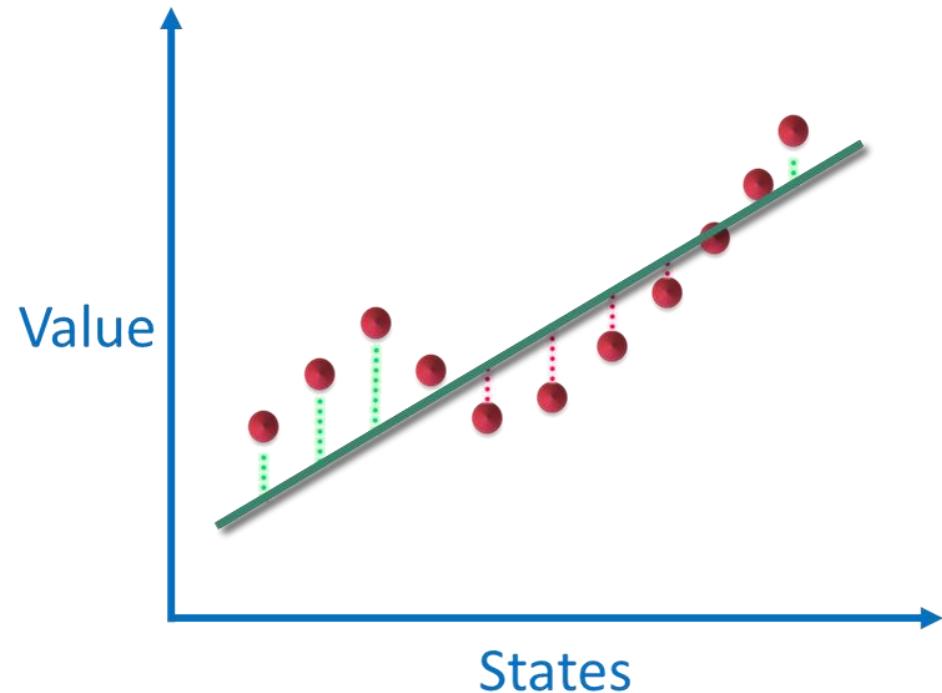
- ▶ gradient decent can overshoot the minimum

$$\textcircled{2} \quad \sum_{n=1}^{\infty} \alpha_n^2(a) < \infty$$

Gradient Descent

Linear function approximation

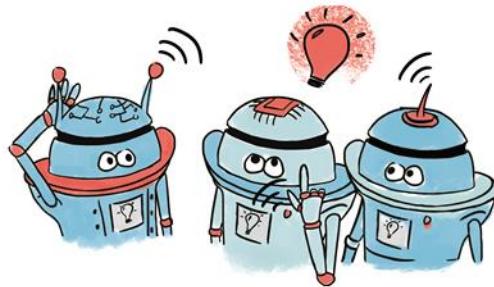
- Linear value approximation is just an inner product of the weights with the feature vector for the state.



$$\hat{v}(s_i, \mathbf{w}) \doteq \sum \mathbf{w}_i x_i(s)$$

$$\frac{\partial \hat{v}(s_i, \mathbf{w})}{\partial w_i} = x_i$$

$$\Delta \hat{v}(s_i, \mathbf{w}) = x(s)$$



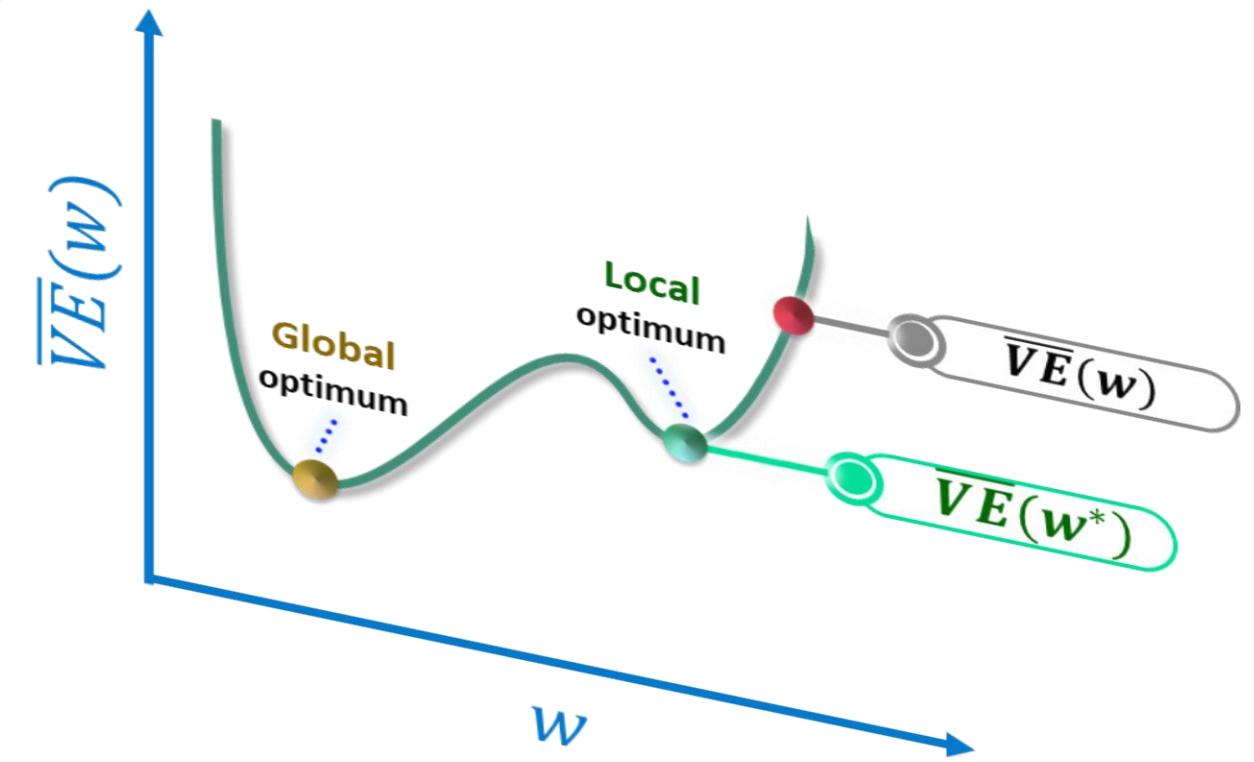
Global Minimum

- ▶ In some cases, gradient descent is guaranteed to converge to the global minimum, which is the best possible setting of the weights for the objective.
 - i.e. the mean squared value error with linear function approximation
- ▶ For more complex function approximate errors such as a neural network, a stationary point may not be a global minimum.

Question

- ▶ In those cases, can we still use gradient descent as an optimization approach? Why?

Stochastic Gradient Descent for (\overline{VE})



Stochastic Gradient Descent

↳ Mathematical Calculation: Gradient Descent of \overline{VE}

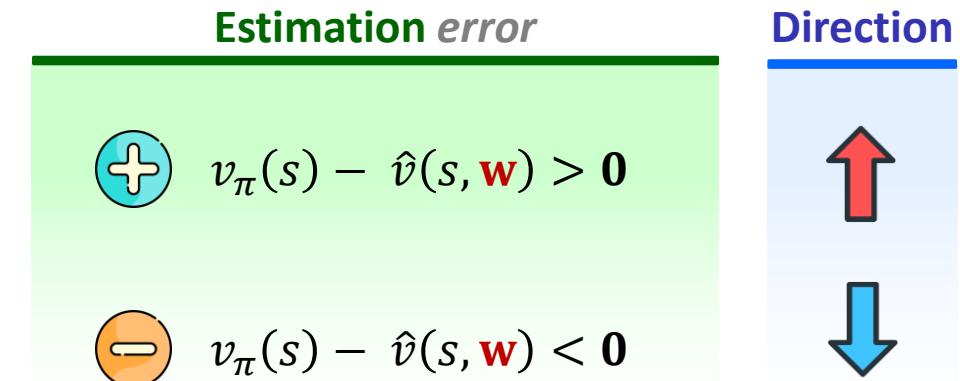
- let us assume that states appear in examples with the same distribution, $\mu(s)$

$$\overline{VE}(\mathbf{w}) \doteq \sum_{s \in \mathcal{S}} [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2$$

- Gradient of \overline{VE}

$$\begin{aligned}\triangledown \overline{VE}(\mathbf{w}) &\doteq \triangledown \sum_{s \in \mathcal{S}} [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 \\ &= \sum_{s \in \mathcal{S}} \triangledown [v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 \\ &= - \sum_{s \in \mathcal{S}} 2[v_\pi(s) - \hat{v}(s, \mathbf{w})] \triangledown \hat{v}(s, \mathbf{w})\end{aligned}$$

Estimation error **Direction**



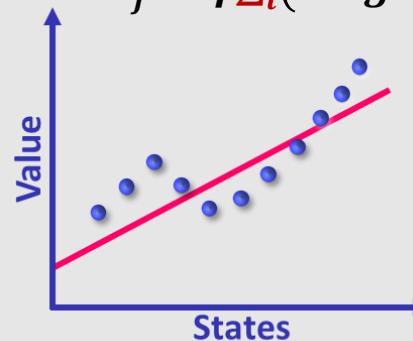
The gradient of the value function approximation depend on the particular parameterized function we are using

Stochastic Gradient Descent

Gradient Descent vs Stochastic Gradient Descent

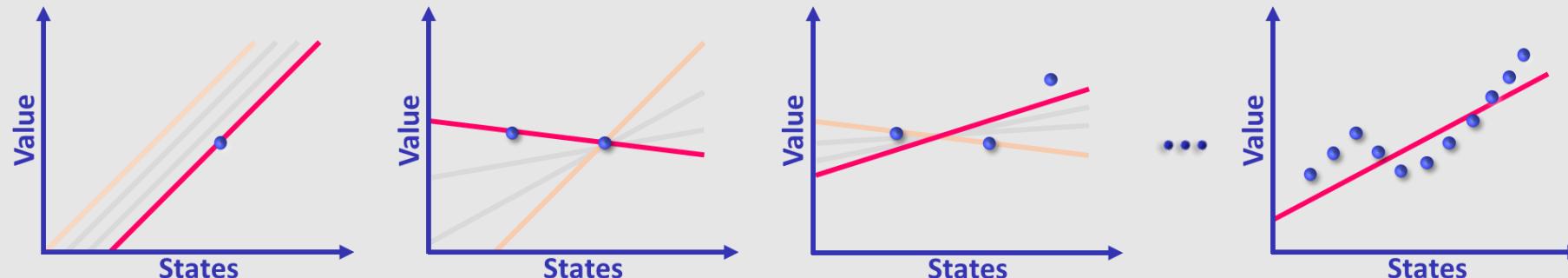
Gradient Descent

$$w_i \doteq w + \Delta w_j, \text{ where: } \Delta w_j = \eta \sum_i (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$$



Stochastic Gradient Descent

$$w_i \doteq w + \Delta w_j, \text{ where: } \Delta w_j = \eta (\text{target}^{(i)} - \text{output}^{(i)}) x_j^{(i)}$$



Stochastic Gradient Descent

↳ Why Stochastic Gradient Descent? (1)

- **Stochastic gradient descent (SGD) is useful when we work with large datasets.**
 - In RL, Value Error objective requires summing over all states, which generally is not feasible

$$\sum_{s \in \mathcal{S}} \mu(s) 2[\nu_\pi(s) - \hat{\nu}(s, \mathbf{w})] \nabla \hat{\nu}(s, \mathbf{w})$$

- Empirically measuring the higher order partial derivatives is computationally too expensive. Therefore, in case of a large size vector \mathbf{w} , gradient decent is not practical. Instead, we use SGD.
- $\mu(s)$ is not available
- There is generally no \mathbf{w} that gets all the states, or even all the examples, exactly correct.

Stochastic Gradient Descent

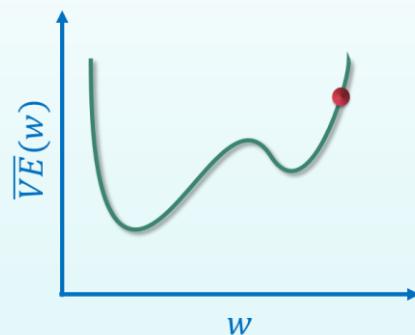
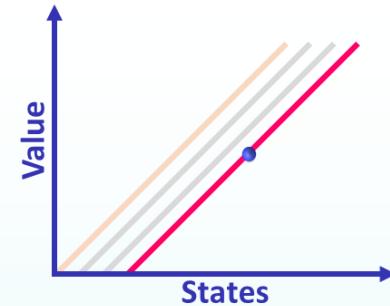
↳ Why Stochastic Gradient Descent? (2)

- **In SGD, we update the weights after each training sample**
 - In RL:
 - » states are selected stochastically (i,e, using a policy)
 - » learning is online, as we observe a new state, we can update weights using SGD
 - » also, we must generalize to all the other states that have not appeared in examples.
- **Also, SGD is efficient when there are redundancy in the data**
 - In RL many states share similar characteristics, and therefore, can be aggregated together.

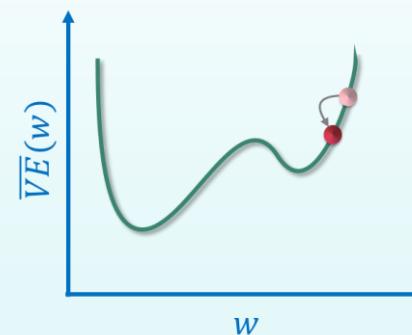
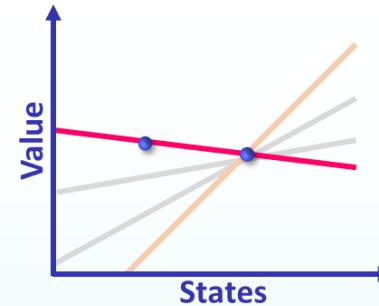
Stochastic Gradient Descent

Graphical representation: SGD for $(\bar{V}E)$

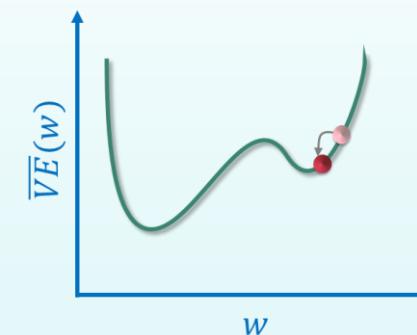
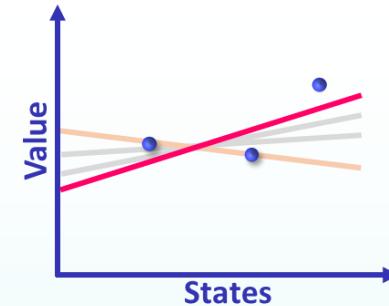
$$S_1, v_{\pi}(S_1)$$
$$w_2 \doteq w_1 + \alpha[v_{\pi}(S_1) - \hat{v}(S_1, w_1)] \nabla \hat{v}(S_1, w_1)$$

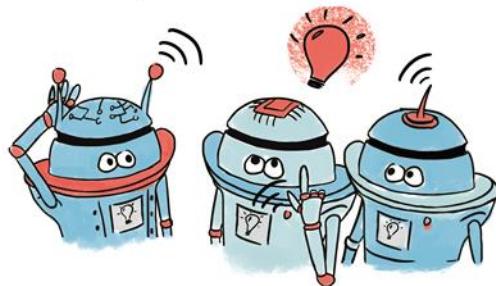


$$S_2, v_{\pi}(S_2)$$
$$w_3 \doteq w_2 + \alpha[v_{\pi}(S_2) - \hat{v}(S_2, w_2)] \nabla \hat{v}(S_2, w_2)$$



$$S_3, v_{\pi}(S_3)$$
$$w_4 \doteq w_3 + \alpha[v_{\pi}(S_3) - \hat{v}(S_3, w_3)] \nabla \hat{v}(S_3, w_3)$$



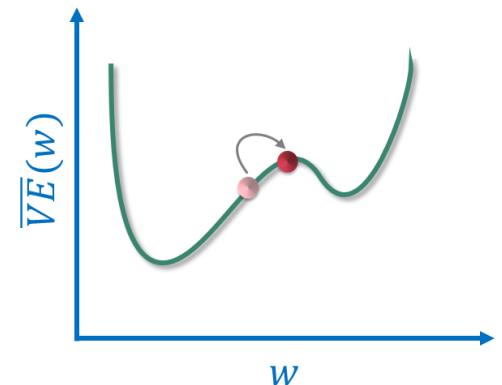


$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [\mathbf{v}_\pi(\mathbf{s}_t) - \hat{\mathbf{v}}(\mathbf{s}_t, \mathbf{w}_t)] \nabla \hat{\mathbf{v}}(\mathbf{s}_t, \mathbf{w}_t)$$

Estimation error
Direction

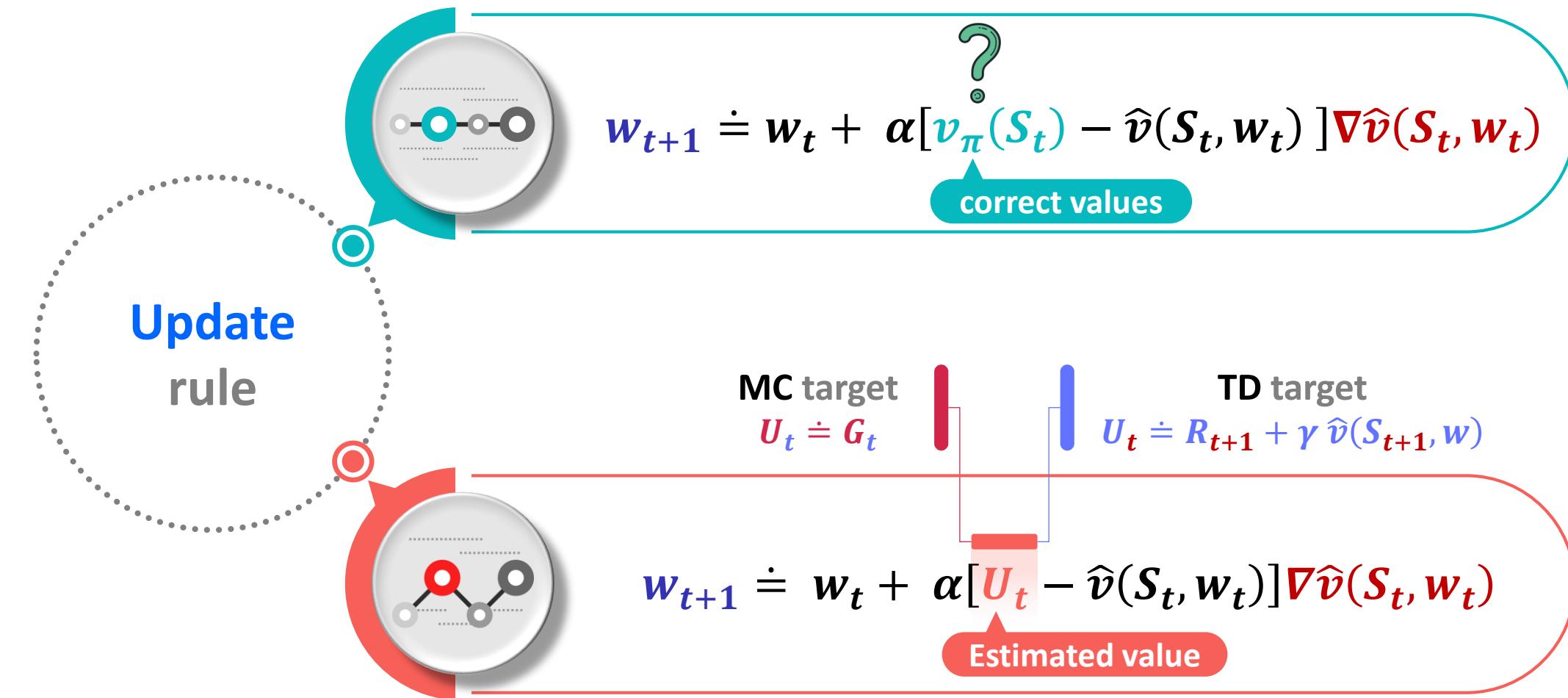
● Stochastic gradient descent

- ▶ Gradient descent methods are called “stochastic” when the update is done on only a single example, which might have been selected stochastically. Over many examples, making small steps, the overall effect is to minimize an average performance measure such as the VE.
 - ▶ It may not be immediately apparent why SGD takes only a small step in the direction of the gradient.
- 1 Could we not move all the way in this direction and completely eliminate the error on the example?**
- 2 In SGD, do we necessarily reduce the error in every single example?**

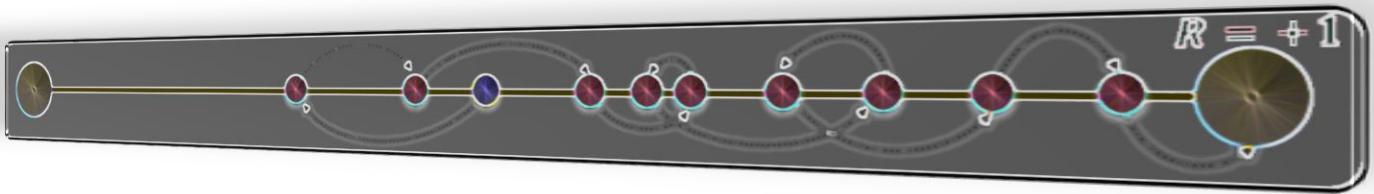


Stochastic Gradient Descent

Update rule



Gradient Monte Carlo and TD



Gradient Monte Carlo and TD

Update Rule with MC estimation

$$v_{\pi}(S_t) \quad U_t$$


- Recall from MCM

Monte Carlo Methods

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

» μ «

\bar{x}



- State-value prediction with MCM



$$w_{t+1} \doteq w_t + \alpha [G_t - \hat{v}(S_t, w_t)] \nabla \hat{v}(S_t, w_t)$$

Gradient Monte Carlo and TD

Pseudocode: Gradient Monte Carlo Algorithm

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

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

1 Initialization

2 Initialize weights

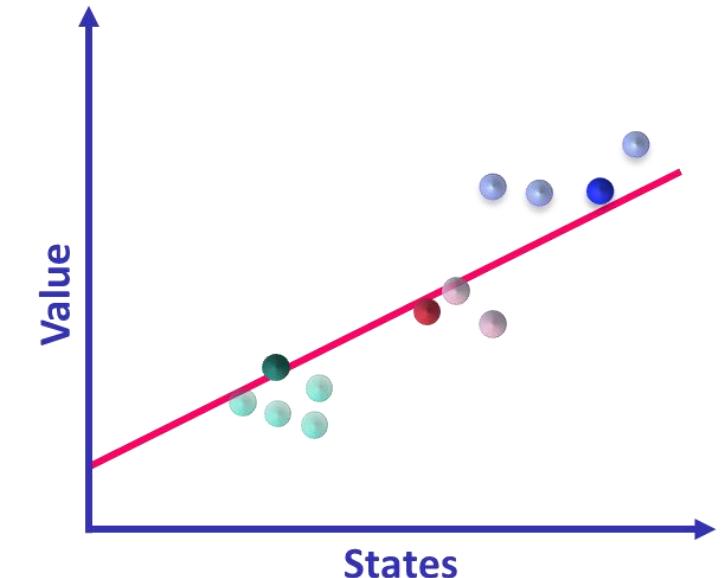
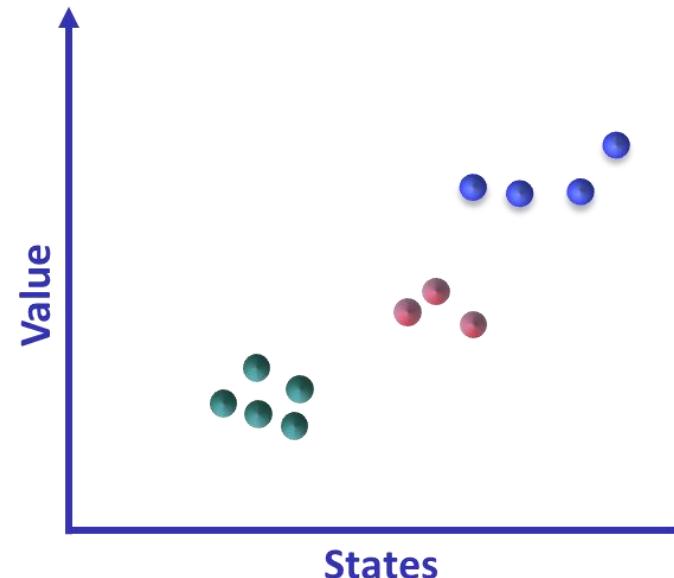
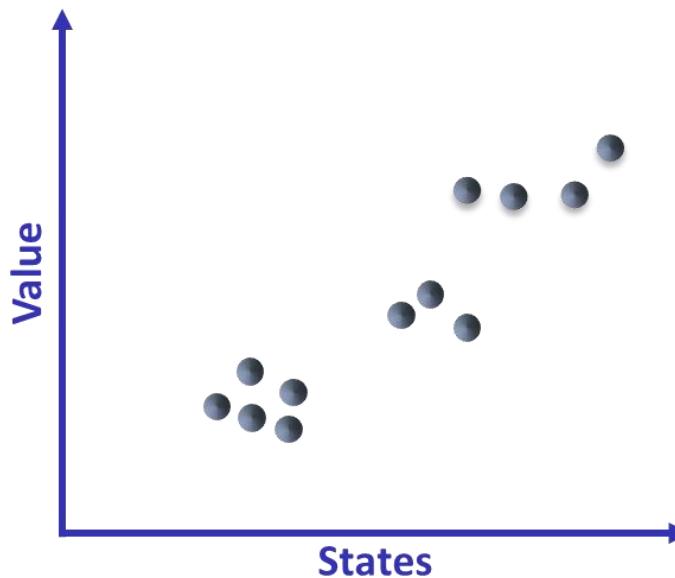
3 Monte Carlo loop
Selecting stochastically

4 SGD
Update weights

Gradient Monte Carlo and TD

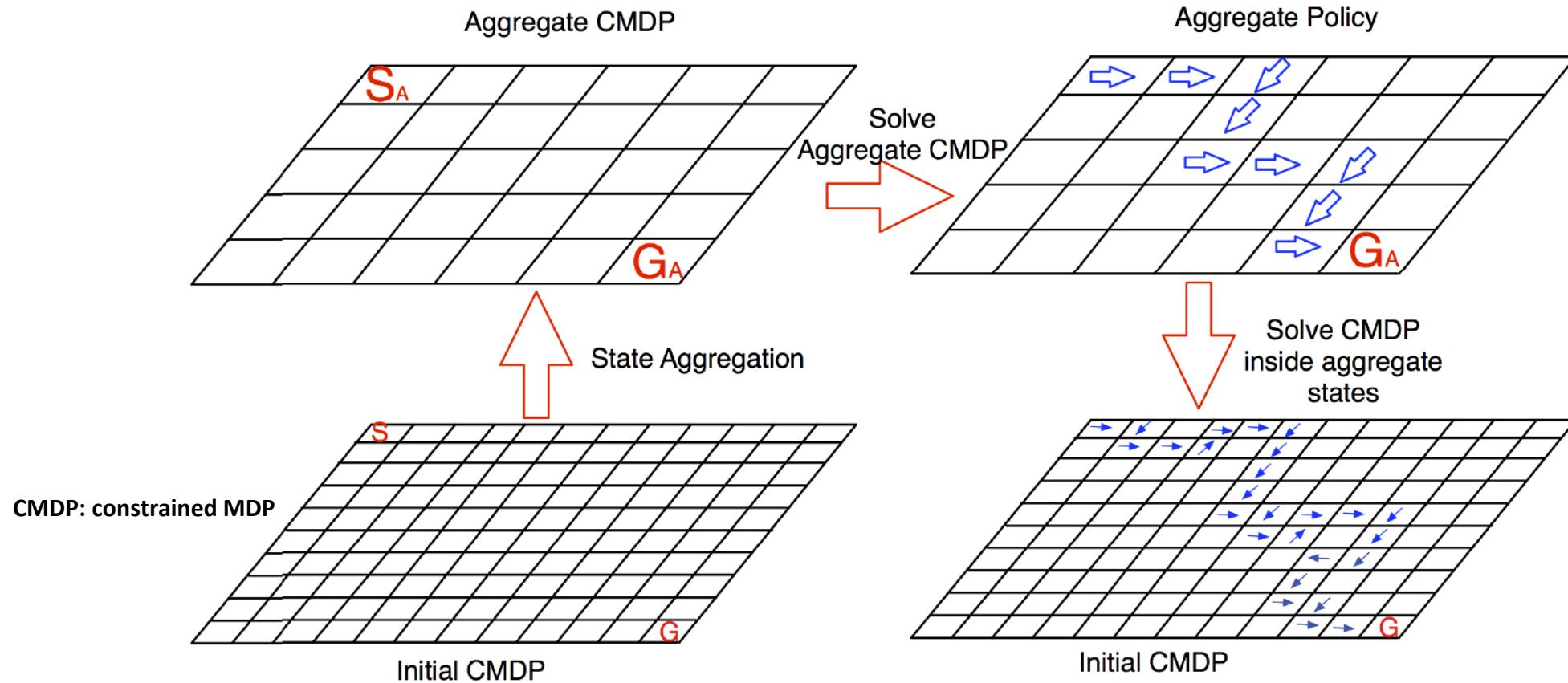
State aggregation

- State aggregation is a simple form of generalizing function approximation in which states are grouped together
 - The value of a group (cluster) generalizes to all states
 - Therefore, there is one estimated value (one component of the weight vector w) for each group.



Gradient Monte Carlo and TD

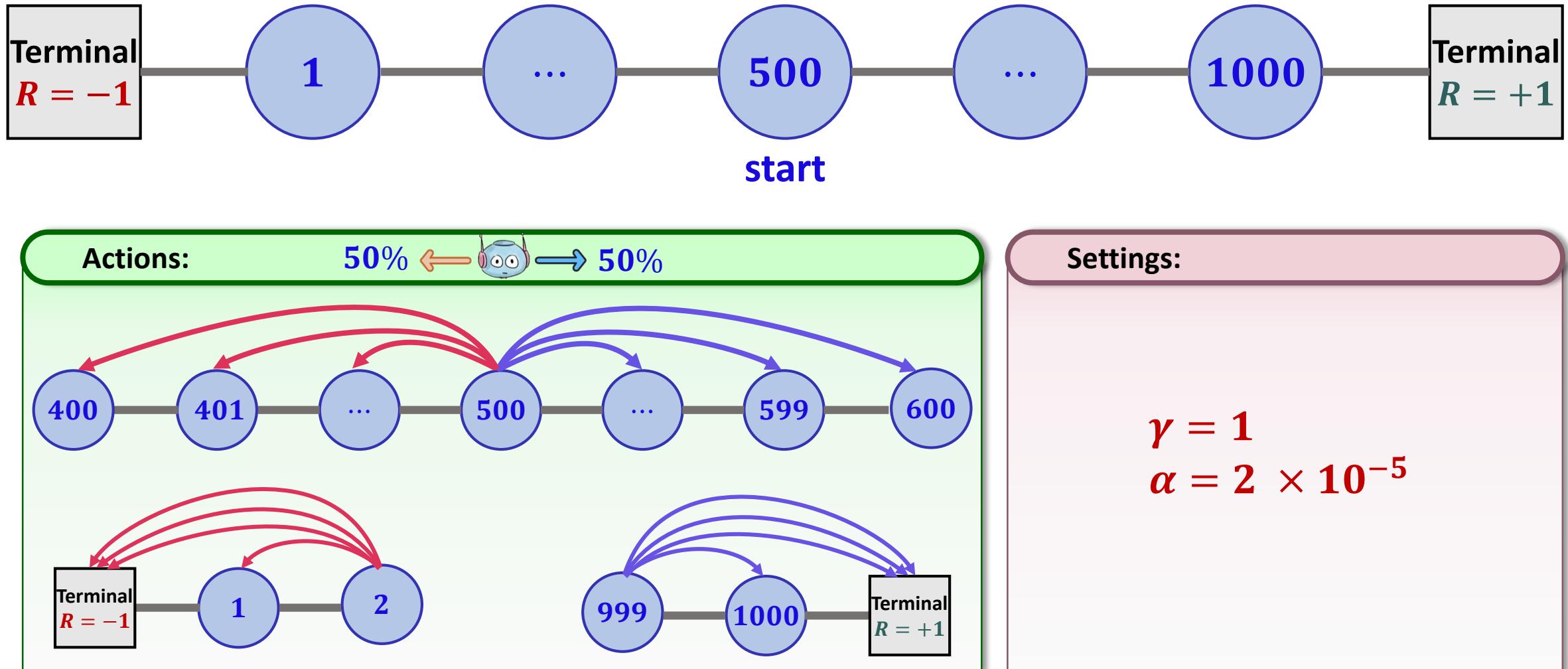
State Aggregation in a GridWorld setting



Feyzabadi, S., & Carpin, S. (2014, August). Risk-aware path planning using hierarchical constrained markov decision processes. In 2014 IEEE International Conference on Automation Science and Engineering (CASE) (pp. 297-303). IEEE.

Gradient Monte Carlo and TD

Example: 1000-state Random Walk



Gradient Monte Carlo and TD

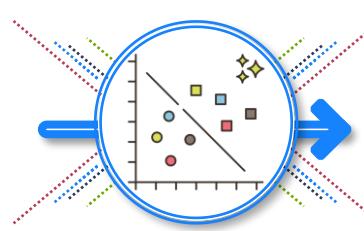
State aggregation

Original states

State	Value
s_1	0
s_2	0
s_3	0
s_4	0
s_5	0
s_6	0
s_7	0
s_8	0

Aggregated states

State	Group	Value
s_1		
s_2		
s_3	Group 1	0
s_4		
s_5		
s_6		
s_7	Group 2	0
s_8		



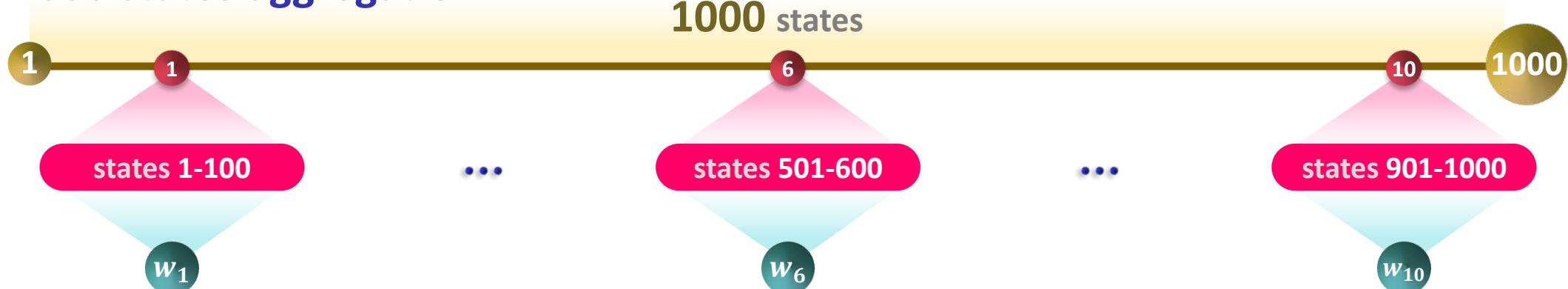
$$x(s) = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad \hat{v}(s, w) = w_1$$

$$x(s) = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \hat{v}(s, w) = w_2$$

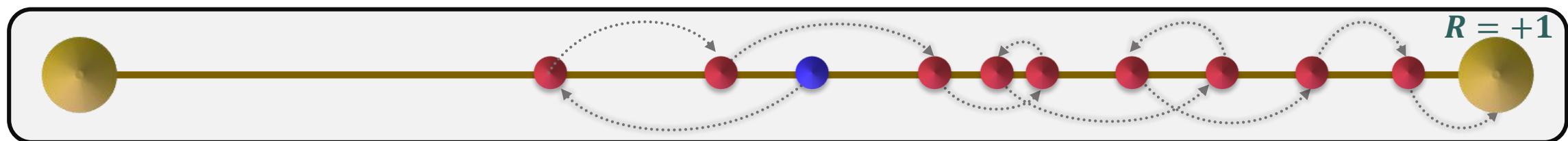
Gradient Monte Carlo and TD

→ State aggregation for 1000-state Random Walk

- **1000 states aggregation**



- **Monte Carlo trajectory**



Visited states 500, 405, 476, ..., 888, 957

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

Update group values $w_5, w_5, w_5, \dots, w_9, w_{10}$ (Terminal state)

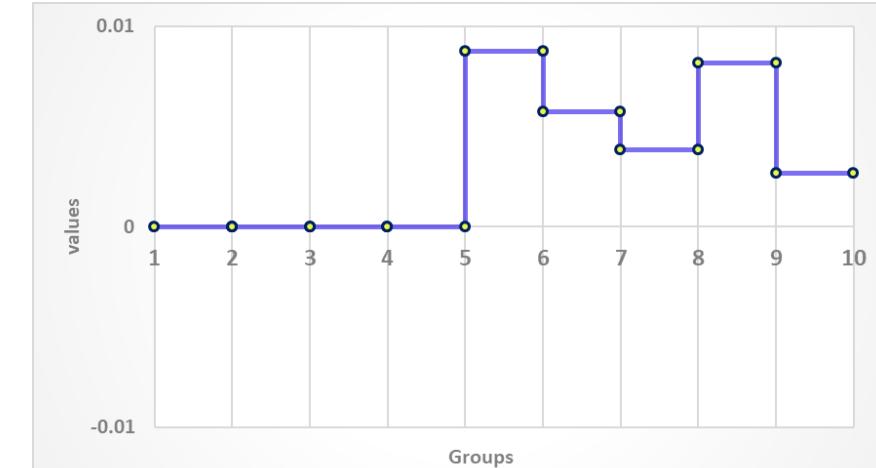
Gradient Monte Carlo and TD

↳ State aggregation for 1000-state Random Walk

Visited states 500, 405, 476, ..., 888, 957

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

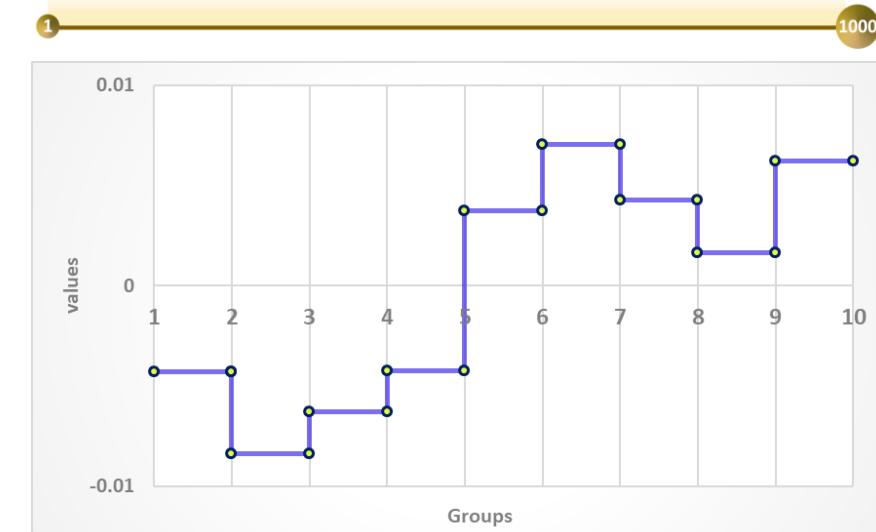
Update group values $w_5, w_5, w_5, \dots, w_9, w_{10}$



Visited states 500, 555, 499, ..., 220, 77

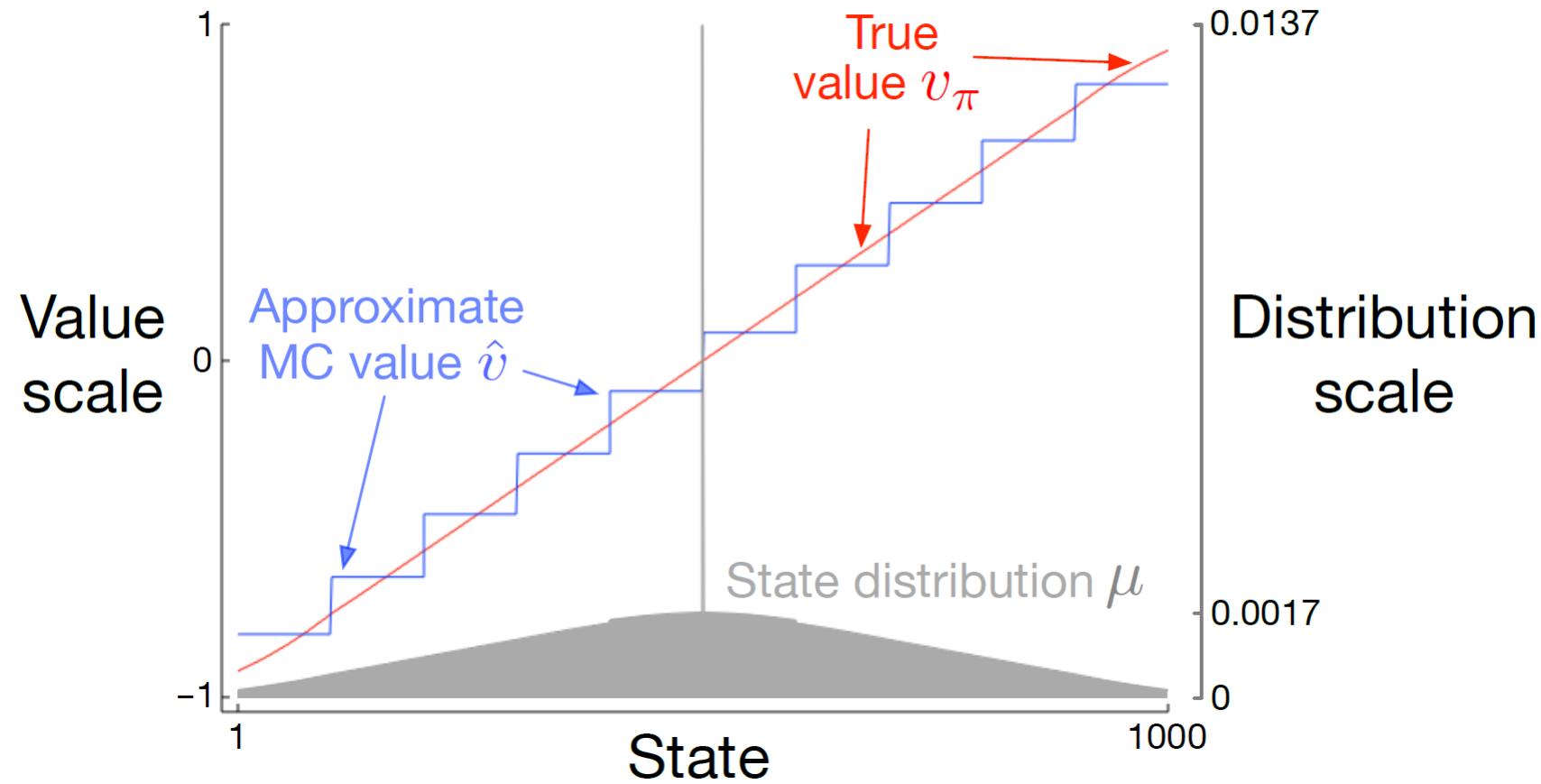
Return -1, -1, -1, ..., -1, -1

Update group values $w_5, w_6, w_5, \dots, w_3, w_1$

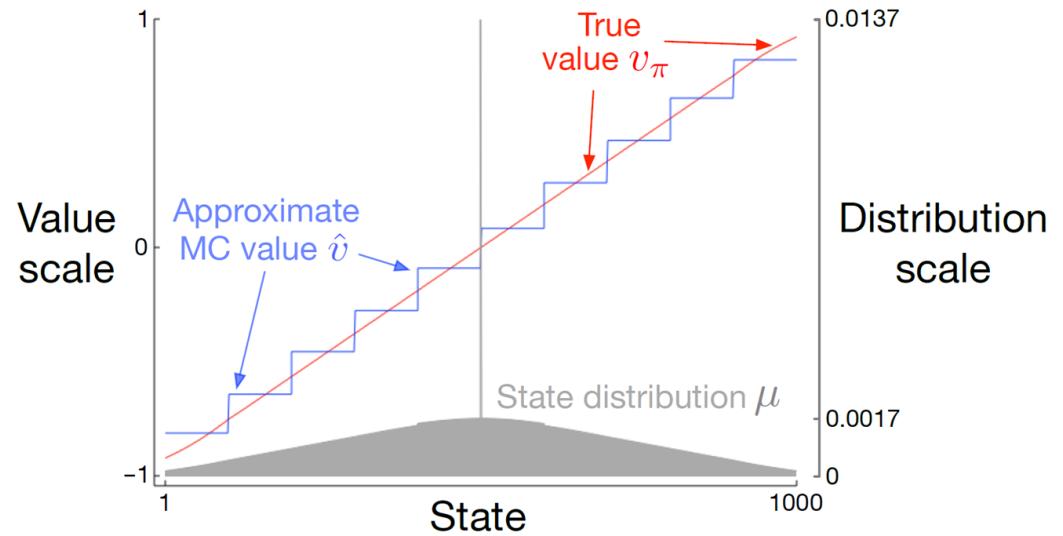
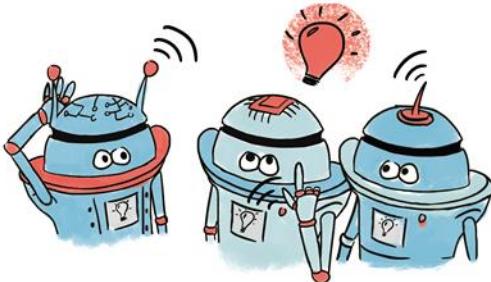


Gradient Monte Carlo and TD

Results: 1000-state Random Walk



PAIR, THINK, SHARE



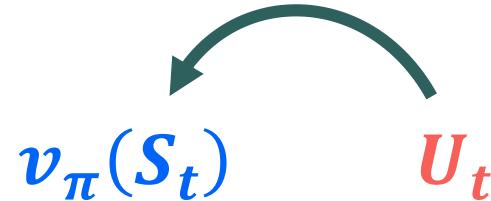
1000-state Random Walk

- ▶ Each step in the plot corresponds to our group of states that share the same approximate value.
- ▶ Some of the details of the approximate values are best appreciated by reference to the state distribution μ for this task, shown in the lower portion of the figure with a right-side scale.

● Explain why the red line does not pass directly through the center of all the blue steps?

Gradient Monte Carlo and TD

Gradient Monte Carlo and TD



- Recall from TD

Temporal Difference

$$V(S_t) \doteq V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

» $V_\pi(S_t)$ «

$V_\pi(S_{t+1})$

- State-value prediction with TD



$$w_{t+1} \doteq w_t + \alpha[R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t) - \hat{v}(S_t, w_t)] \nabla \hat{v}(S_t, w_t)$$

Stochastic Gradient Descent

Pseudocode: Semi-gradient TD(0)

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

- 1 Input: the policy π to be evaluated
- 1 Input: a differentiable function $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$ such that $\hat{v}(\text{terminal}, \cdot) = 0$
- 1 Algorithm parameter: step size $\alpha > 0$
- 2 Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)
- 3 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

- 1 Initialization
- 2 Initialize weights
- 3 Initialize weights
- 4 Semi-gradient TD update

Stochastic Gradient Descent

Convergence of function approximation

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [\mathbf{U}_t - \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)] \nabla \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)$$

Convergence of stochastic approximation

- If U_t is an unbiased estimate, that is, if $\mathbb{E}[U_t | S_t = s] = v_\pi(s_t)$ for each t , then w_t is guaranteed to converge to a local optimum under the usual stochastic approximation conditions for decreasing α .



Monte Carlo Method

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [\mathbf{G}_t - \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)] \nabla \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)$$

- Unbiased estimate:** the general SGD method converges to a locally optimal approximation to $v_\pi(s_t)$
- Therefore, it is guaranteed to find a locally optimal solution



Temporal Difference

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [\mathbf{R}_{t+1} + \gamma \widehat{\mathbf{v}}(\mathbf{S}_{t+1}, \mathbf{w}_t) - \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)] \nabla \widehat{\mathbf{v}}(\mathbf{S}_t, \mathbf{w}_t)$$

- Biased estimate:** in bootstrapping methods like DP and TD, **target value** depends on the current value of the weight vector w_t
- Bootstrapping methods are not in fact instances of true gradient descent (**semi-gradient methods**)

Linear Methods



Linear Methods

Linear methods approximate

- Linear methods approximate the state-value function by the inner product between w and $x(s)$:

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

- $\hat{v}(s, w)$, is a linear function of the weight vector, w .

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

- » In this case the approximate value function is said to be *linear in the weights*, or simply *linear*.
 - » For linear methods, features are basis functions because they form a linear basis for the set of approximate functions.

Linear Methods

↳ Tabular TD as a special case of Linear TD

Tabular representation

State	Values
S_1	v_1
S_2	v_2
S_3	v_3
:	:
S_d	v_d

Function Approx. representation

State	Values
S_1	$\hat{v}(S_1, w)$
S_2	$\hat{v}(S_2, w)$
:	:
S_{10}	$\hat{v}(S_{10}, w)$
:	
S_d	$\hat{v}(S_d, w)$

$$\hat{v}(S_i, \mathbf{w}) \doteq \langle \mathbf{w}, \mathbf{x}(S_i) \rangle = w_i$$

$$w_1 \quad w_2 \quad \dots \quad \mathbf{w}_{10} \quad \dots \quad w_d \quad \bullet \quad \begin{matrix} x_1 = 0 \\ x_2 = 0 \\ \vdots \\ x_{10} = 1 \\ \vdots \\ x_d = 0 \end{matrix} = w_{10}$$



Linear Methods

↳ SGD updates and its convergence

- The gradient of the approximate value function with respect to w in this case is:

$$\nabla \hat{v}(s, w) = x(s)$$

- Thus, in the linear case the general *SGD update rule* is:

$$w_{t+1} \doteq w_t + \alpha [U_t - \hat{v}(S_t, w_t)] x(S_t)$$

- Because it is so simple, the linear SGD case is one of the most favorable for mathematical analysis.

- Convergence results

- In particular, in the linear case there is only *one optimum* (or, in degenerate cases, one *set of equally good optima*)
- Thus, any method that is guaranteed to converge to or near a local optimum is automatically guaranteed to converge to or near the global optimum.
 - » i.e. gradient Monte Carlo algorithm



Linear Methods

The True Objective for TD: TD fixed point (1)

- The semi-gradient TD(0) algorithm converges under linear function approximation but this does not follow from general results on SGD.

- In TD(0), the update at each time t is:

$$\begin{aligned} w_{t+1} &\doteq w_t + \alpha [R_{t+1} - \gamma w_t^T x_{t+1} - w_t^T x_t] x_t \\ &= w_t + \alpha [R_{t+1} x_t - x_t (x_t - \gamma x_{t+1})^T w_t] \end{aligned}$$

shorthand notation
 $x_t : x(S_t)$

- The update rule can be written as:

$$\Delta w = \alpha [R_{t+1} x_t - x_t (x_t - \gamma x_{t+1})^T w_t]$$

noise expected update

Linear Methods

The True Objective for TD: TD fixed point (2)

$$\Delta w = \alpha [R_{t+1}x_t - x_t(x_t - \gamma x_{t+1})^T w_t]$$

- The expected weight update vector can be written

$$\mathbb{E}[\Delta w] = \alpha(b - Aw_t), \text{ where:}$$

noise $b \doteq \mathbb{E}[R_{t+1}x_t] \in \mathbb{R}^d$

expected update $A \doteq \mathbb{E}[x_t(x_t - \gamma x_{t+1})^T] \in \mathbb{R}^d \times \mathbb{R}^d$

$$\mathbb{E}[\Delta w_{TD}] = \alpha(b - Aw_{TD})$$

- By setting $\mathbb{E}[\Delta w_{TD}] = 0$

$$w_{TD} = A^{-1}b$$

TD fixed point

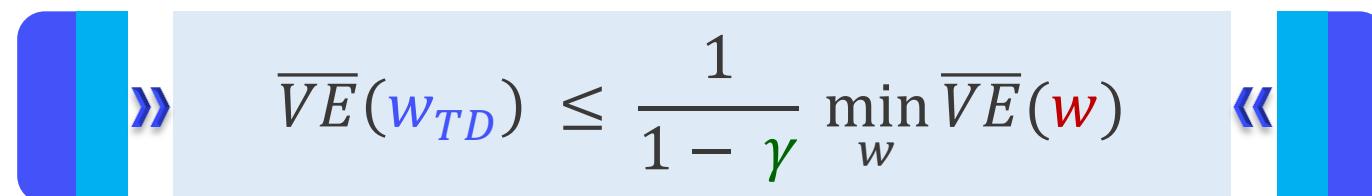
Linear Methods

↳ Value Error of TD fixed point

- **Review Linear TD**

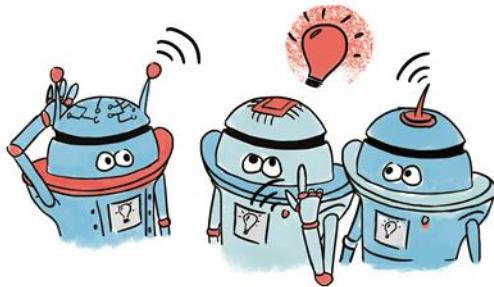
- Biased estimation
- Convergence to TD fixed point (check SB book for proof)
- Linear TD approximates values based on Bellman equation and minimize an objective function called *Projected Bellman error*.

- At the TD fixed point, it has also been proven (in the continuing case) that the \overline{VE} is within a bounded expansion of the lowest possible error:


$$\overline{VE}(\mathbf{w}_{TD}) \leq \frac{1}{1 - \gamma} \min_w \overline{VE}(\mathbf{w})$$



PAIR, THINK, SHARE



$$\overline{VE}(\mathbf{w}_{TD}) \leq \frac{1}{1 - \gamma} \min_{\mathbf{w}} \overline{VE}(\mathbf{w})$$

TD fixed point and the Minimum of \overline{VE}

- ▶ Linear TD does not minimize the MSE objective, but in fact minimizes ***Projected Bellman error***
 - ▶ *However, TD fixed point is very close to the minimum value error solution*
- 1 **What is the impact of Gamma (γ) on the difference between the TD fixed point and the minimum value error solution**
 - 2 **Does this bound depend on the quality of the features? What happened if we can perfectly represent the value function**
 - 3 **In general, why isn't the TD fixed point equal to the minimum value error solution?**

Summary

