

Problem: arbitrarily large state spaces



Solution: find a good approximate solution

Using limited computational resources  
to learn a better policy

"generalization"

function approximation

Approximate  $V_\pi$  from experience generated  
Using a known policy  $\pi$

parameterized function:

$$\hat{V}(s, w) \approx V_\pi(s)$$

$w \in \mathbb{R}^d$  weight vector.

Linear model

DNN

Decision Tree

$$d \ll |\mathcal{S}|$$

$\left\{ \begin{array}{l} d: \text{the number of weights} \\ |\mathcal{S}|: \text{the number of states} \end{array} \right.$

RL

update to an estimated value function that  
shift its value at particular states toward a

"back-up value" or "update target"  
for that state

$$\boxed{S \longrightarrow v}$$

$$M_{\pi} : S_t \longrightarrow G_t$$

$$\text{TD}(0) : S_t \longrightarrow R_{t+1} + \gamma \hat{v}(S_{t+1}, w)$$

$$n\text{-step TD} : S_t \longrightarrow G_{t:t+n}$$

$$\text{DP} : S \longrightarrow \mathbb{E}_{\pi} [R_{t+1} + \gamma \hat{v}(S_{t+1}, w_t) | S_t = S]$$

Prediction Objective  $\overline{VE}$

Mean Squared Value Error ( $\overline{VE}$ )

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

$\mu(s)$  : fraction of time spent in  $s$   
(on-policy distribution)

$$\mu(s) \geq 0, \quad \sum_s \mu(s) = 1$$



$$\eta(s) = h(s) + \sum_{\bar{s}} \eta(\bar{s}) \sum_a \pi(a|\bar{s}) p(s|\bar{s}, a)$$

$$\mu(s) = \frac{\eta(s)}{\sum_{s'} \eta(s')}$$

$\overline{VE}$  : find a global optimum

Stochastic gradient and semi-gradient Method:

Stochastic gradient descent (SGD):

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

$\hat{v}(s, w)$  : differentiable function of  $w$  for all  $s$

$$S_t \rightarrow V_\pi(S_t)$$

$$\begin{aligned} w_{t+1} &\doteq w_t - \frac{1}{2} \alpha \nabla [V_\pi(S_{t+1}) - \hat{v}(S_t, w_t)]^2 \\ &= w_t + \alpha [V_\pi(S_t) - \hat{v}(S_t, w_t)] \nabla \hat{v}(S_t, w_t) \end{aligned}$$

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

target output:  $U_t \in \mathbb{R}$

$$W_{t+1} \doteq W_t + \alpha [U_t - \hat{v}(S_t, W_t)] \nabla \hat{v}(S_t, W_t)$$

$U_t$ : unbiased estimate

$$\mathbb{E}[U_t | S_t = s] = V_{\pi}(S_t)$$

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

Gradient MC:

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

Semi-gradient TD(0):

$$W \leftarrow W + \alpha [R + \gamma \hat{v}(s', W) - \hat{v}(s, W)] \nabla \hat{v}(s, W)$$



# Linear Models

$$X(s) \doteq (X_1(s), X_2(s), \dots, X_d(s))^T$$

$$\hat{V}(s, w) \doteq w^T X(s) \doteq \sum_{i=1}^d w_i X_i(s)$$

(linear in the weights)

$X(s)$ : feature vector

$$\nabla \hat{V}(s, w) = X(s)$$

$$W_{t+1} \doteq W_t + \alpha [U_t - \hat{V}(S_t, W_t)] X(S_t)$$

$$\doteq W_t + \alpha (R_{t+1} + \gamma W_t^T X_{t+1} - W_t^T X_t) X_t$$

$$= W_t + \alpha (R_{t+1} X_t - X_t (X_t - \gamma X_{t+1})^T W_t)$$

$$\mathbb{E}[W_{t+1} | W_t] = W_t + \alpha (b - A W_t)$$

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

$$A \doteq \mathbb{E}[X_t (X_t - \gamma X_{t+1})^T] \in \mathbb{R}^d \times \mathbb{R}^d$$

$$b - A W_{T0} = 0$$

$$b = A w_{TD}$$

$$w_{TD} = A^{-1}b \quad (\text{TD fixed point})$$

ANN: Artificial Neural Networks

Semi-linear units

Weighted sum  $\rightarrow$  activation function  
(nonlinear function)

S-shaped (Sigmoid)

logistic function  $f(x) = \frac{1}{1 + e^{-x}}$

rectifier nonlinearity  $f(x) = \max(0, x)$

"modify the objective function to discourage complexity of the approximation"

"The dropout method efficiently approximates this combination by multiplying each outgoing weight of a unit by the probability that that unit



was retained during training.

Method:  $\left\{ \begin{array}{l} \text{gradient descent} \\ \text{Semi-gradient descent} \end{array} \right.$

Linear Methods

Nonlinear Function Approximation: ANN

Least-Squares TD

Memory-based Function Approximation

Kernel-based Function Approximation