

Lecture -06

Dynamic Programming
Monte Carlo
Temporal Difference (TD) learning
n-step TD methods.

Value Function approximation

Tabular methods

- Large state space
 - memory
 - sig. time & data.
- States rarely repeat.

Generalization

Supervised learning:

$$(x_1, y_1)$$

$$(x_2, y_2)$$

$$f: x \rightarrow y$$

→ Linear models

→ Non linear models
(Neuralnet)

$$(x_n, y_n)$$

goal: estimate value fns.

- ① large state spaces.
- ② need for generalization to unseen states.



- ① learning incrementally.
- ② Value FA \rightarrow Nonstationary target fn.

forget: $R_{t+1} + \gamma \underline{\underline{v_{\pi}(s_{t+1})}}$

Value fn. approximation:-

$V_{\pi}(s)$ - value fn.

$$\tilde{V}(s; \omega) \stackrel{\text{weights}}{\approx} V_{\pi}(s)$$

↑
approximation.

$$d \ll |s|$$

$$\omega \in \mathbb{R}^d$$

linear model:

's'

$$x(s) = (\alpha_1(s), \alpha_2(s), \dots, \alpha_d(s))$$

representation of the state

$$\tilde{V}(s; \omega) = \omega^T x(s)$$

$$= \sum_{i=1}^d \omega_i x_i(s)$$

State	Val
s_1	
s_2	
s_3	
\vdots	
s_{100}	

tabular approach
no generalization

State forget update
 $S \mapsto u$

Mobile Grb $S_t \mapsto G_t$

1-step TD $S_t \mapsto R_{t+1} + \gamma \hat{V}(S_{t+1}; w_t)$

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

DP $s \mapsto \mathbb{E}_{\pi} \left[R_{t+1} + \gamma \hat{V}(s_{t+1}; w_t) \right]_{s=s}$

(S_t, G_t)

FA:

~~FA~~

$d \ll |S|$

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

The predict objec.

Mean squared error.

$$\overline{VE}(\omega) = \sum_{s \in S} \mu(s) (V_{\pi}(s) - \hat{V}(s; \omega))^2$$

$$VE(\omega) = \frac{1}{|S|} \sum_{s \in S} (V_\pi(s) - \hat{V}(s; \omega))^2$$

$$\boxed{\omega^* = \underset{\omega}{\operatorname{argmin}} \quad VE(\omega)}$$

$\mu(s)$

mean
standard
error

$$\bar{VE}(\omega) = \sum_{s \in S} \mu(s) (V_\pi(s) - \hat{V}(s; \omega))^2$$

Square root of \bar{VE}



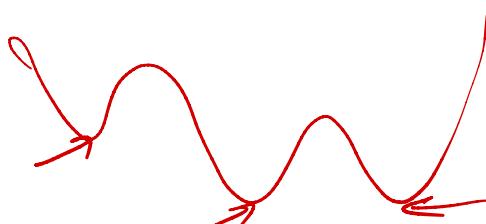
Policy gradient methods

$$\omega^* = \underset{\omega}{\operatorname{argmin}} \quad \bar{VE}(\omega)$$

global opt'
 ω^*

$$\bar{VE}(\omega^*) \leq \bar{VE}(\omega) \text{ for all possible } \omega$$

local optimum



$$\bar{VE}(\omega^*) \leq \bar{VE}(\omega)$$

for all ω in
some neighbor of ω^* .

On-policy prediction with approximation:-

Stochastic gradient descent :-

$$\omega = (\omega_1, \omega_2 \dots \omega_d)^T$$

$$\hat{v}(s; \omega) = \omega^T x(s)$$

$$\omega_0, \omega_1, \omega_2, \omega_3 \dots$$

$$s_t \mapsto v_{\pi}(s_t)$$

SGD:

$$\omega_{t+1} = \omega_t - \frac{1}{2} \alpha \frac{\partial}{\partial \omega} [v_{\pi}(s_t) - \hat{v}(s_t; \omega_t)]$$

$$\omega_{t+1} = \omega_t + \alpha [v_{\pi}(s_t) - \hat{v}(s_t; \omega_t)] \frac{\partial \hat{v}(s_t; \omega)}{\partial \omega}$$

$$\omega_{t+1} = \omega_t + \alpha [U_t - \hat{v}(s_t; \omega_t)] \frac{\partial \hat{v}(s_t; \omega)}{\partial \omega}$$

if U_t is unbiased,

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

$$U_t = G_t$$

Prediction:

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

Input: the policy π to be evaluated

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

$$v(S; w) = w^\top \mathbf{x}(s)$$

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

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

n-step return : $G_{t:T+n}$

$$\text{DP} : \sum_{a, s', r} \text{tfals}_t P(s' | s_t, a) [r + \gamma \hat{V}(s'; w_t)]$$

Step-TD: $R_{t+1} + \gamma \hat{V}(s_{t+1}; w)$

$$\text{Error} = \left(\boxed{R_{t+1} + \gamma \hat{V}(s_{t+1}; w)} - \hat{v}(s_t; w) \right)^2$$

Constant

$$\frac{\partial \text{error}}{\partial w}$$

Semi-gradient :

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'

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

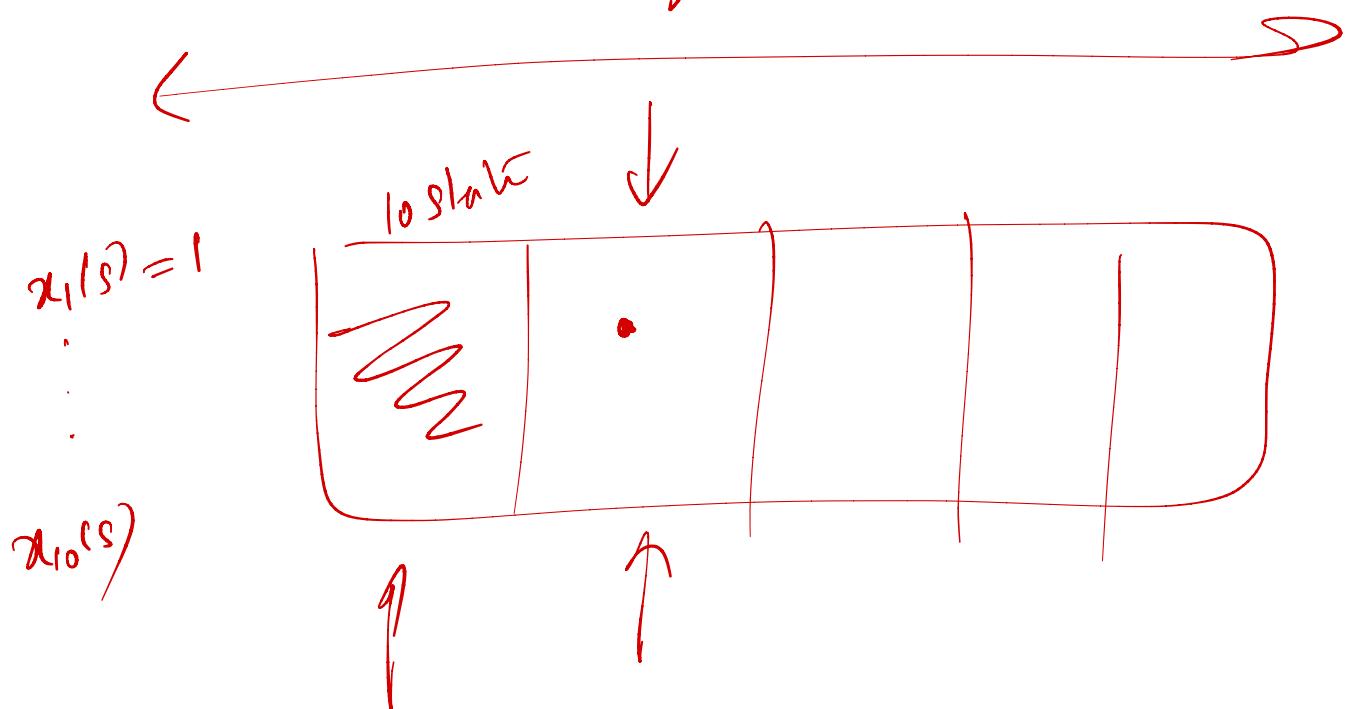
$$S \leftarrow S'$$

 until S is terminal

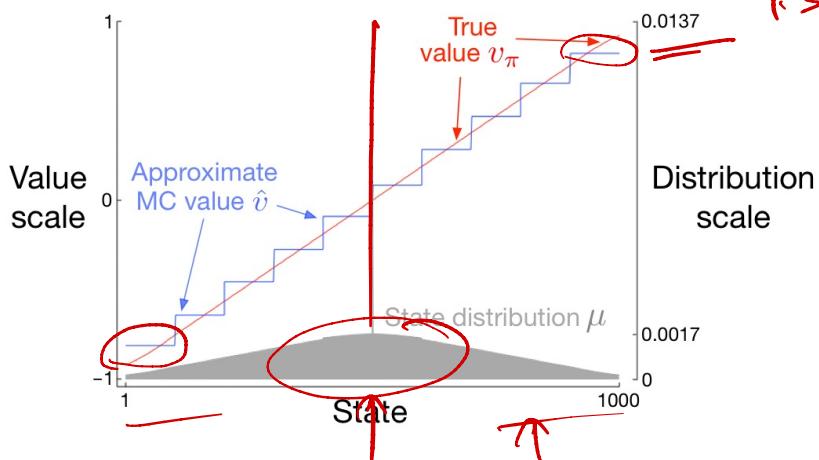
*tabular
methods*

*state
aggregation*

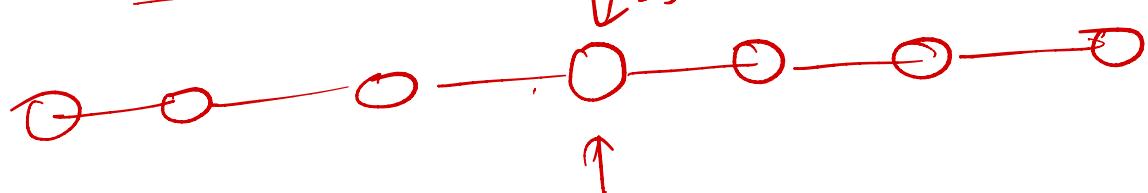
*function
approx*



1.37-1.



(500-state) random walk



-1

+1

Linear methods:

$$x(s) = (n_1(s), n_2(s), \dots, n_d(s))^T$$

$$\hat{v}(s; \omega) = \underbrace{\omega^T x(s)}_{= \sum_i w_i n_i(s)} = \sum_{i=1}^d w_i n_i(s)$$

$$\frac{\partial}{\partial \omega} \hat{v}(s; \omega) = x(s)$$

SGD update:

$$\omega_{t+1} = \omega_t + \alpha [U_t - \hat{v}(s_t; \omega_t)] x(s_t)$$

Semi-gradient TD :

Conting Case :

$$\begin{aligned}\omega_{t+1} &= \omega_t + \alpha \left(\frac{R_{t+1} + \gamma \omega_t^\top x_{t+1}}{b} - \frac{\omega_t^\top x_t}{A} \right) x_t \\ &= \omega_t + \alpha \left(\frac{R_t x_t}{b} - \frac{x_t (x_t - \gamma x_{t+1})^\top \omega_t}{A} \right)\end{aligned}$$

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

where $b = \mathbb{E}[R_t x_t] \in \mathbb{R}^d$

$$A = \mathbb{E}[x_t (x_t - \gamma x_{t+1})^\top] \in \mathbb{R}^{d \times d}$$

at Convergence ω_{TD}

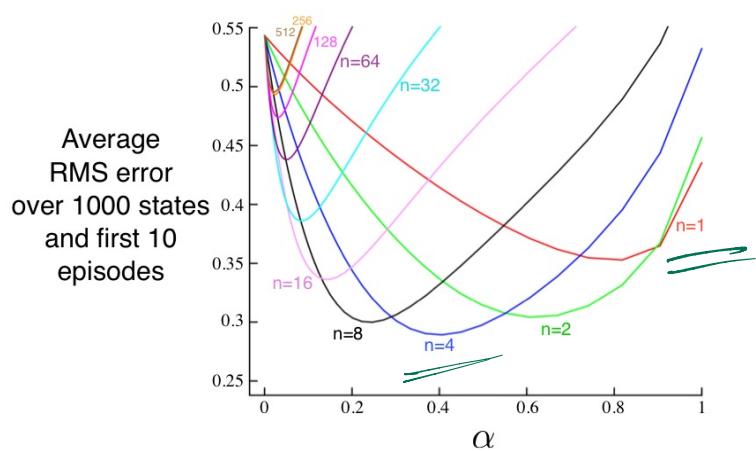
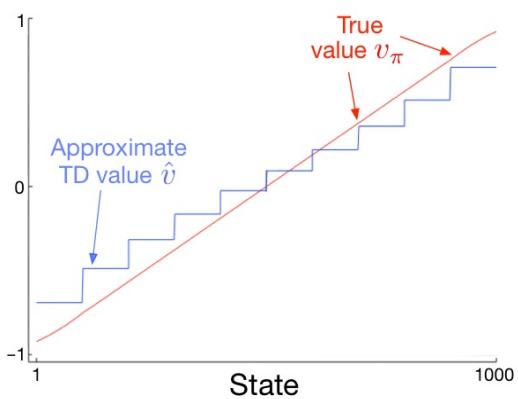
$$b - A\omega_{TD} = 0$$

$$b = A\omega_{TD}$$

$$\omega_{TD} = A^{-1}b$$

TD
fixed
point-

$$\widehat{VE}(\omega_{tb}) \leq \frac{1}{1-\gamma} \min_{\overline{\omega}} \widehat{VE}(\omega)$$



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

$$A = \mathbb{E} [x_t (x_t - \gamma x_{t+1})^T]$$

$$b = \mathbb{E} [R_{t+1} x_t]$$

Compute estimate of A and b .

Least-Squares TD (LSTD)

$$\hat{A}_t = \sum_{k=0}^{t-1} x_k (x_k - \gamma x_{k+1})^T + GI$$

$$\hat{b}_t = \sum_{k=0}^{t-1} R_{k+1} x_k$$

$$w_t = \hat{A}_t^{-1} \hat{b}_t$$

Semi-gradient JD

LSTD.

- iterative algork.
estm. val. → est. A, b.

$O(d)$

Comp. $A_t : O(d^2)$

$A_t^{-1} : O(d^3)$

$\cancel{O(d^2)}$

$$\hat{A}_t^{-1} = \left(\hat{A}_{t-1} + x_{t-1} (x_{t-1} - \gamma x_t)^T \right)^{-1}$$
$$\hat{A}_0 = C\mathbb{I} = \hat{A}_{t-1} - \frac{\hat{A}_{t-1} x_{t-1} (x_{t-1} - \gamma x_t)^T \hat{A}_{t-1}^{-1}}{1 + (x_{t-1} - \gamma x_t)^T \hat{A}_{t-1}^{-1} x_{t-1}}$$

Shermann - Morrison formula.

LSTD for estimating $\hat{v} = \mathbf{w}^\top \mathbf{x}(\cdot) \approx v_\pi$ ($O(d^2)$ version)

Input: feature representation $\mathbf{x} : \mathcal{S}^+ \rightarrow \mathbb{R}^d$ such that $\mathbf{x}(\text{terminal}) = \mathbf{0}$

Algorithm parameter: small $\varepsilon > 0$

$$\widehat{\mathbf{A}}^{-1} \leftarrow \varepsilon^{-1} \mathbf{I}$$

$$\mathbf{b} \leftarrow \mathbf{0}$$

A $d \times d$ matrix

A d -dimensional vector

Loop for each episode:

Initialize S ; $\mathbf{x} \leftarrow \mathbf{x}(S)$

Loop for each step of episode:

Choose and take action $A \sim \pi(\cdot | S)$, observe R, S' ; $\mathbf{x}' \leftarrow \mathbf{x}(S')$

$$\mathbf{v} \leftarrow \widehat{\mathbf{A}}^{-1} (\mathbf{x} - \gamma \mathbf{x}')$$

$$\widehat{\mathbf{A}}^{-1} \leftarrow \widehat{\mathbf{A}}^{-1} - (\mathbf{A}^{-1} \mathbf{x}) \mathbf{v}^\top / (1 + \mathbf{v}^\top \mathbf{x})$$

$$\widehat{\mathbf{b}} \leftarrow \widehat{\mathbf{b}} + R \mathbf{x}$$

$$\mathbf{w} \leftarrow \widehat{\mathbf{A}}^{-1} \widehat{\mathbf{b}}$$

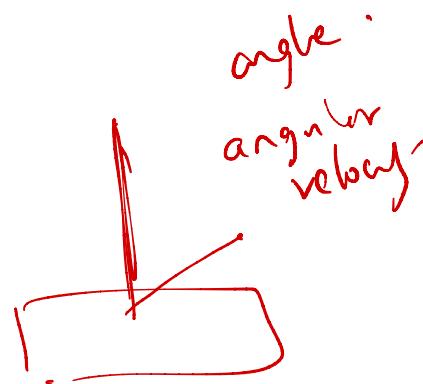
$$S \leftarrow S'; \mathbf{x} \leftarrow \mathbf{x}'$$

until S' is terminal

Representation of states! —

Linear models:

$$\mathbf{w}^\top \mathbf{x}$$



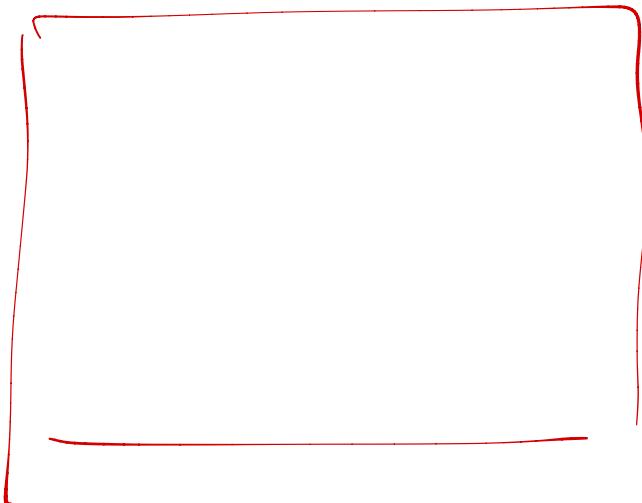
Polynomials

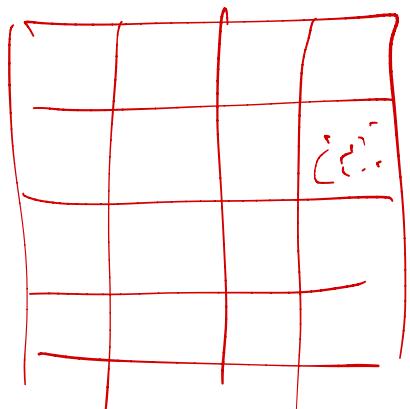
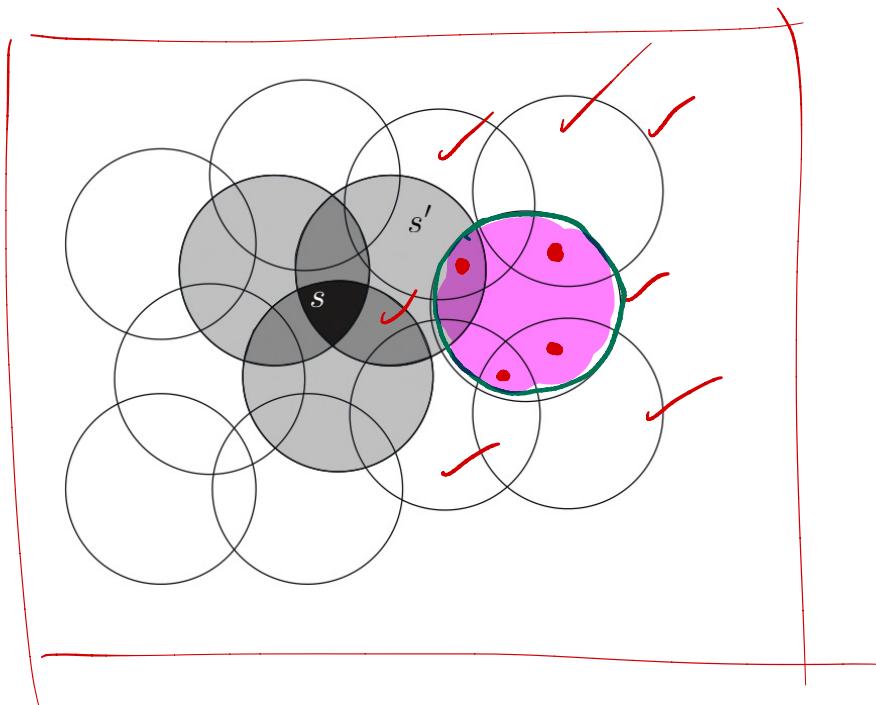
$$(s_1, s_2)$$

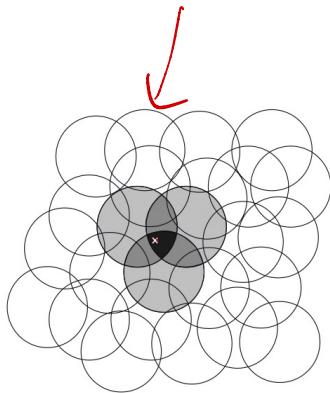
$$(1, s_1, s_2, s_1 s_2)$$

$$(1, s_1, s_2, s_1^2, s_2^2, s_1 s_2)$$

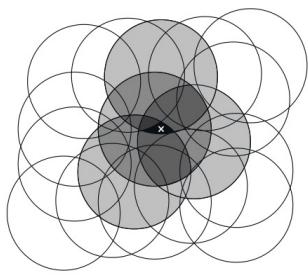
Course Coding:



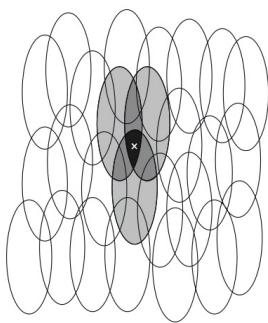


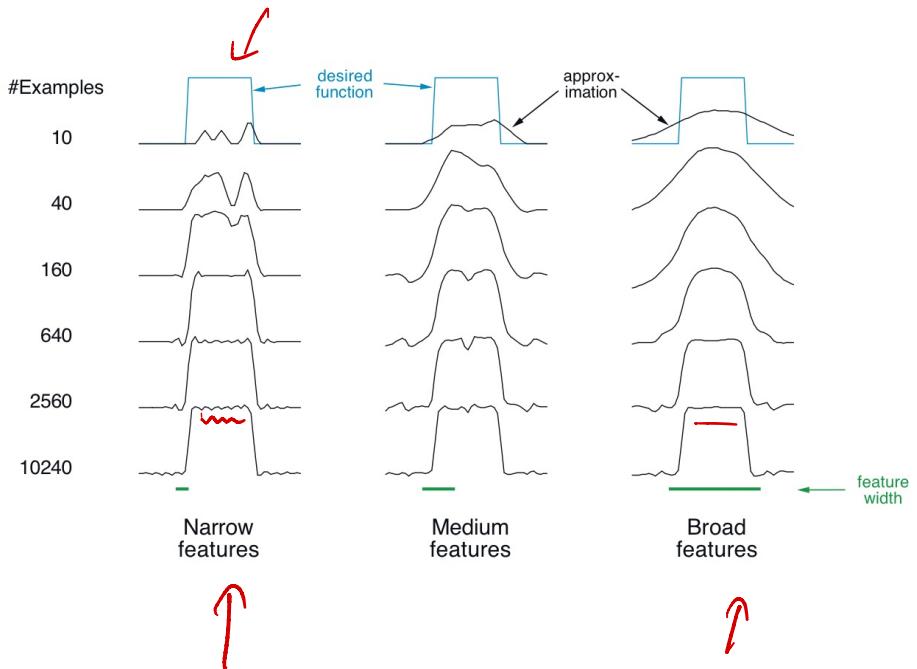


Narrow



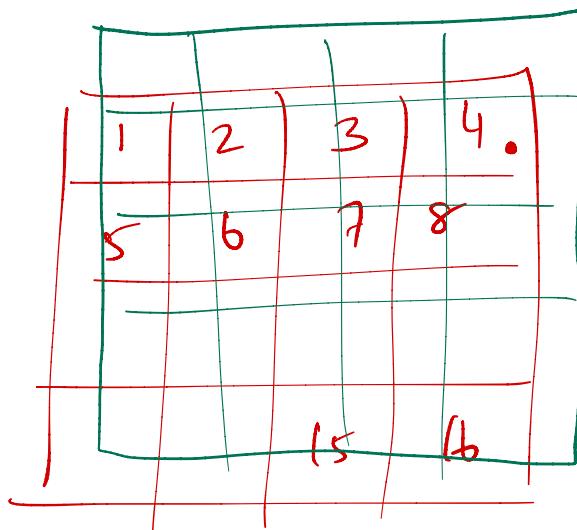
broad

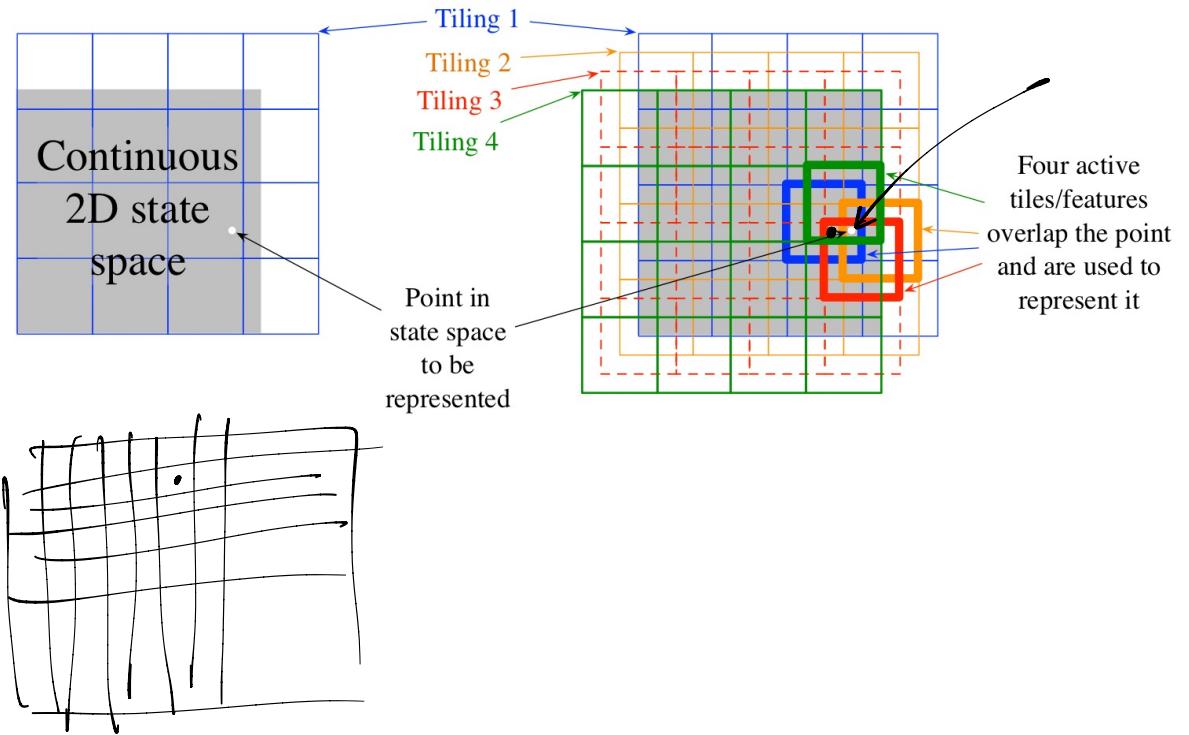


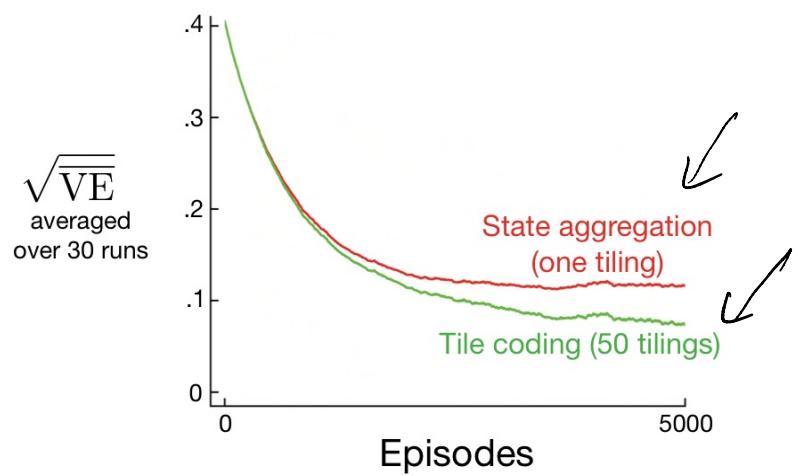


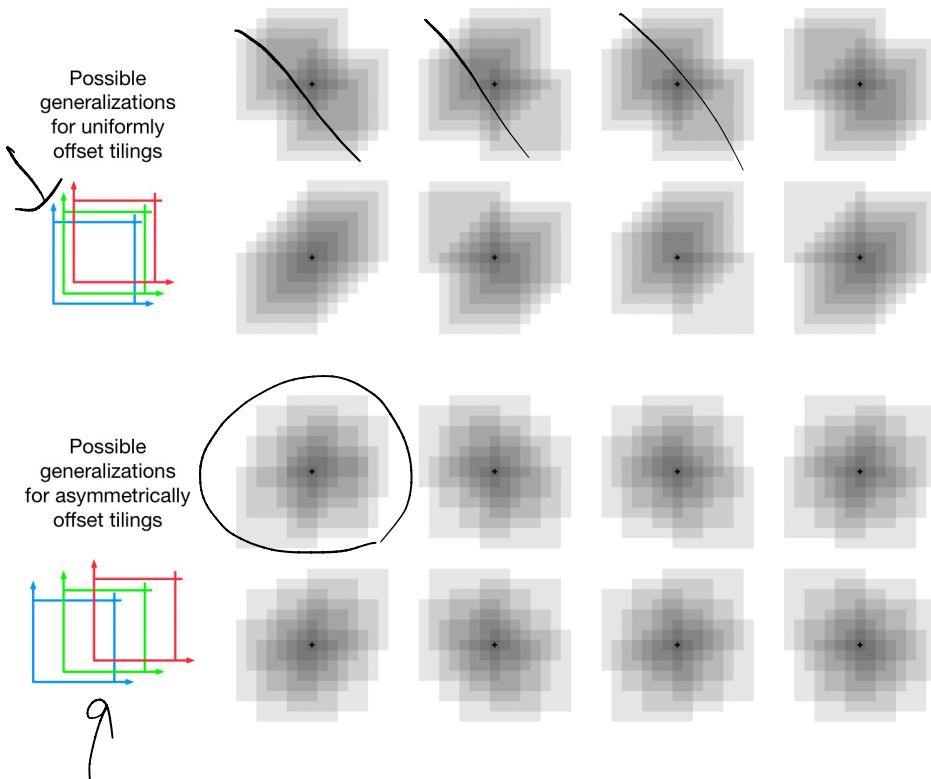
~~(1)(1)~~

Tile-Coding :









On-policy Control with Approximation :-

$$w_{t+1} = w_t + \alpha \left[v_t - \hat{q}(s_t, a_t; w_t) \right] \underbrace{\frac{\partial \hat{q}(s_t, a_t; w_t)}{\partial w}}$$

one-step Sarsa :-

$$w_{t+1} = w_t + \alpha \left[R_{t+1} + \gamma \hat{q}(s_{t+1}, a_{t+1}; w_t) - \hat{q}(s_t, a_t; w_t) \right] \underbrace{\frac{\partial \hat{q}(s_t, a_t; w_t)}{\partial w}}$$



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

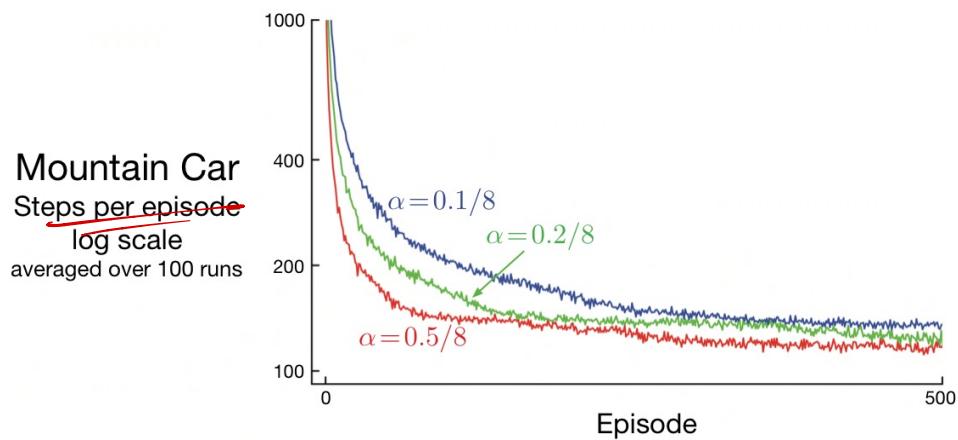
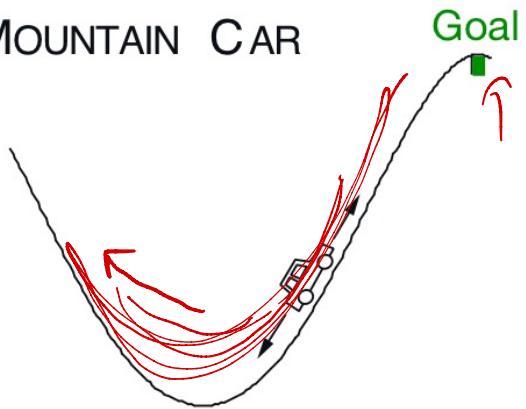
 Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

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

$S \leftarrow S'$

$A \leftarrow A'$

MOUNTAIN CAR



Episodic semi-gradient n -step Sarsa for estimating $\hat{q} \approx q_*$ or q_π

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

Input: a policy π (if estimating q_π)

Algorithm parameters: step size $\alpha > 0$, small $\varepsilon > 0$, a positive integer n

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

All store and access operations (S_t , A_t , and R_t) can take their index mod $n + 1$

Loop for each episode:

 Initialize and store $S_0 \neq$ terminal

 Select and store an action $A_0 \sim \pi(\cdot | S_0)$ or ε -greedy wrt $\hat{q}(S_0, \cdot, \mathbf{w})$

$T \leftarrow \infty$

 Loop for $t = 0, 1, 2, \dots$:

 | If $t < T$, then:

 | Take action A_t

 | Observe and store the next reward as R_{t+1} and the next state as S_{t+1}

 | If S_{t+1} is terminal, then:

 | $T \leftarrow t + 1$

 | else:

 | Select and store $A_{t+1} \sim \pi(\cdot | S_{t+1})$ or ε -greedy wrt $\hat{q}(S_{t+1}, \cdot, \mathbf{w})$

 | $\tau \leftarrow t - n + 1$ (τ is the time whose estimate is being updated)

 | If $\tau \geq 0$:

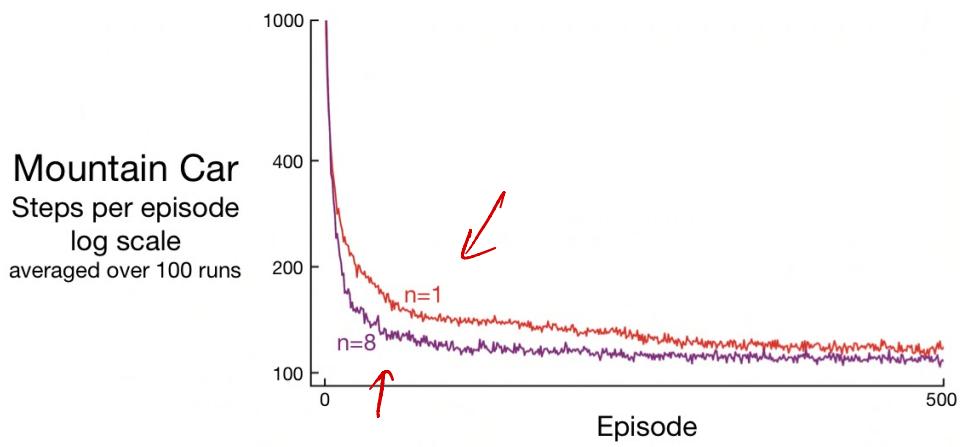
 | $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$

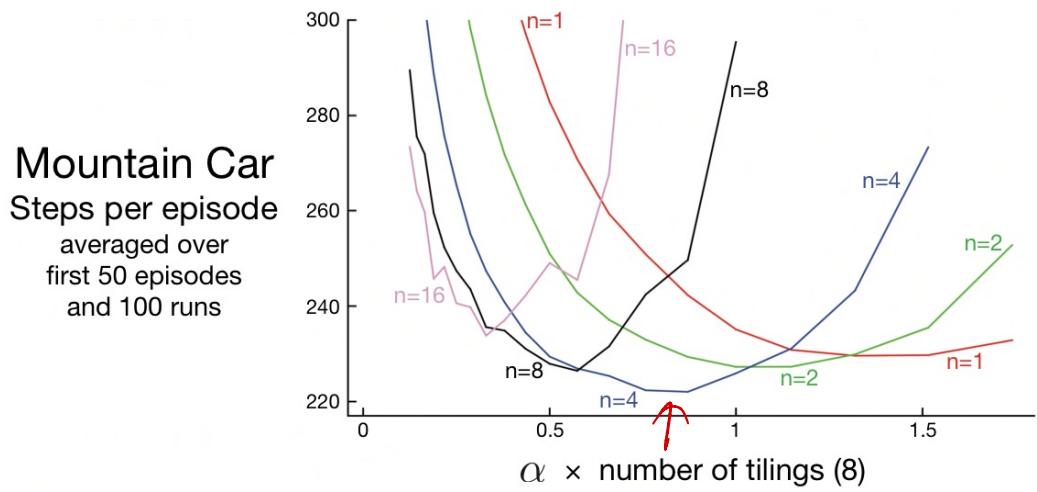
 | If $\tau + n < T$, then $G \leftarrow G + \gamma^n \hat{q}(S_{\tau+n}, A_{\tau+n}, \mathbf{w})$

($G_{\tau:\tau+n}$)

 | $\mathbf{w} \leftarrow \mathbf{w} + \alpha [G - \hat{q}(S_\tau, A_\tau, \mathbf{w})] \nabla \hat{q}(S_\tau, A_\tau, \mathbf{w})$

 | Until $\tau = T - 1$





Value fn. approximation

