

24th Sept

MONTE-CARLO METHODS

(1) Prediction

Given π , Find v^π through interaction with MDP.
Control (Also wants to change and converge to a good policy)

$v^\pi(s)$ Episodic Task Assumption

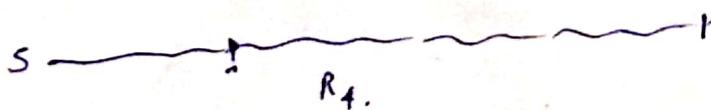
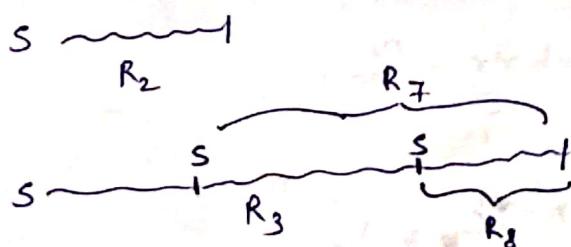
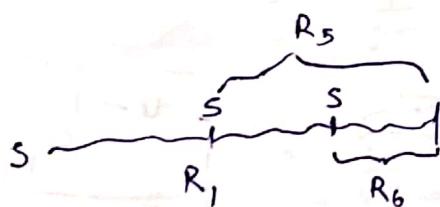
i: Run $s_i^0 \ a_i^0 \ r_i^0 \ s_i^1 \ a_i^1 \ r_i^1 \ s_i^2 \ \dots \ s_i^{q_i}$

i-th "Return": $r_i^0 + \gamma r_i^1 + \gamma^2 r_i^2 + \dots + \gamma^{q_i-1} r_i^{q_i-1}$

$R_1, R_2, R_3, \dots, R_n$

Estimate $v^\pi \approx \frac{1}{n}(R_1 + R_2 + \dots + R_n)$

$v^\pi(s)$



$$v_{\pi_1}^\pi(s) = \frac{R_1 + R_2 + R_3 + R_4}{4}$$

$$\text{Instead, } v_{\pi_2}^\pi(s) = \frac{R_1 + R_2 + R_3 + R_4 + R_5 + R_6 + R_7 + R_8}{8}$$

Is $\lim_{n \rightarrow \infty} v_1^n(s) = v^\pi(s)$? Yes \rightarrow FIRST-VISIT MONTE CARLO
 Is $\lim_{n \rightarrow \infty} v_2^n(s) = v^\pi(s)$? Yes! Proof different from 1st one.
 \rightarrow EVERY-VISIT MONTE CARLO.

$\# n E[v_1^n(s)] = v^\pi(s)$ Unbiased Estimator for v^π

$$E[v_2^n(s)] \neq v^\pi(s) \rightarrow = v^\pi(s) + [b(s \text{ AT } R), n]$$

As long as you don't look at the future, you're okay, but
 if you something like taking last or second-last occurrence of
 s is not correct.

Every-visit probability converges faster.

② Control

Assumption of exploring starts.

Monte Carlo's exploring states

$$Q[s][a] \leftarrow 8$$

$\pi \leftarrow$ Arbitrary policy

$$\text{visits}[s][a] \leftarrow 0$$

For each episode: $e = 1, 2, \dots$

Pick $s \in S$ uniformly at random

Pick $a \in A$ uniformly at random

Start at s , take a , follow π , let return be R .

$$Q[s][a] \leftarrow \frac{Q[s][a] * \text{visits}[s][a] + R}{\text{visits}[s][a] + 1}$$

$$\text{visits}[s][a] \leftarrow +1$$

$\pi \leftarrow$ greedy (Q)

$\pi \leftarrow \underset{a \in A}{\text{argmax}} Q(s, a)$.

(Suppose you pick s using arbit distribution on S) \rightarrow need not converge.

Is $\pi_{c \rightarrow \infty} = \pi^*$? \rightarrow This will converge.

If you converge you'll converge to optimal policy

Problem: (Probably Approximately Correct Formulation)
Continuing Tasks
Exploring starts (Random actions)
 (s, a)
Rewards b/w R_{\min} & R_{\max} .

Prediction: $v^{\pi} \leftarrow \text{Given } \pi$ (Initial) $\rightarrow v = [v^{\pi}(s), v]$ is off

Write an algorithm L

Required: Let L be output v_L as the

We want with probability at least $(1-\delta)$ that

$$\forall s \in S \quad |v_L(s) - v^{\pi}(s)| \leq \epsilon$$

(i) How many episodes?

(ii) Length of each episode?

Try and work out yourselves.

Next: Bootstrapping.

1st Oct

① Estimators

Process with parameter $p \in [0, 1/2]$

$$x \sim \text{Bernoulli}(p)$$

$$y \sim \text{Bernoulli}(2p)$$

Run process $x=1, y=0$

What is p ?

$$\text{Let's take } \hat{p} = \frac{x+y}{2}$$

* Least-squares estimator

$$\hat{p} = \underset{p \in [0, 1/2]}{\operatorname{argmin}} (p-1)^2 + (2p-0)^2$$

$$= \underset{p \in [0, 1/2]}{\operatorname{argmin}} p^2 - 2p + 1 + 4p^2$$

$$p = 0.2$$

* Maximum Likelihood Estimator

$$\hat{p} = \underset{p \in [0, 1/2]}{\operatorname{argmax}} p(1-2p)$$

$$= 0.25$$

$$p = 2p^2$$

② π V^n s

$$s \sim R_1$$

$$s \sim R_2$$

$$\sim R_3$$

$$\sim R_n$$

$$V_n(s) = \frac{R_1 + R_2 + \dots + R_n}{n}$$

$$V_{n+1}(s) = \frac{R_1 + R_2 + \dots + R_{n+1}}{n+1}$$

$$= V_n(s) \left(1 - \frac{1}{n+1}\right) + R_{n+1} \left(\frac{1}{n+1}\right)$$

Result (Robbins and Monro, 1951)

Let $(\alpha_i)_{i=1}^{\infty}$ be a sequence such that $\sum_{i=1}^{\infty} \alpha_i = \infty$ and $\sum_{i=1}^{\infty} \alpha_i^2 < \infty$,

Consider updating $v_{n+1}(s) \leftarrow v_n(s) \{1 - \alpha_{n+1}\} + \frac{\alpha_{n+1} r_{n+1}}{\downarrow}$.

Then, $\lim_{n \rightarrow \infty} v_n(s) = v^{\pi}(s)$. Learning Rate.

Stochastic Approximation - Robbins & Monro.

$$③ s = s^0 \underset{n+1}{\alpha} r^0 \underset{n+1}{\alpha} s^1 \underset{n+1}{\alpha} r^1 \cdots$$

Monte-Carlo Method $\leftarrow v_{n+1}(s) \leftarrow v_n(s) \{1 - \alpha_{n+1}\} + \alpha_{n+1} \{r^0_{n+1} + \gamma r^1_{n+1} + \cdots + \gamma^2 r^2_{n+1} \cdots\}$

Bootstrapping $v_{n+1}(s) \leftarrow v_n(s) \{1 - \alpha_{n+1}\} + \alpha_{n+1} \{r^0_{n+1} + \gamma v_n(s^1_{n+1})\}$

Form of

Full Bootstrapping

$$v^{\text{TD}(0)}_{n+1}(s) \leftarrow v^{\text{TD}(0)}_n(s) + \frac{\alpha_{n+1} \left\{ r^0_{n+1} + \gamma v^{\text{TD}(0)}_n(s^1_{n+1}) - v^{\text{TD}(0)}_n(s) \right\}}{\text{Old Estimate}}$$

Temporal Difference prediction error.

④ TD(0) Algorithm

$v_0 \leftarrow$ Initial guess of v^{π}

Assume that agent is born in state s^0

For $t = 0, 1, 2, \dots$

Take action $a^t = \pi(s^t)$

Obtain r^t, s^{t+1}

$$v_{t+1}(s^t) \leftarrow v_t(s^t) + \alpha_{t+1} \{r^t + \gamma v_t(s^{t+1}) - v_t(s^t)\}$$

~~s^{t+1} s^{t+1} s^{t+1}~~

$$\lim_{t \rightarrow \infty} v_t = v^{\pi}$$

⑤ Batch MC and Batch TD(0)

TB(1) sars sars sars \rightarrow (2), and girlobug stabbing

Assume that episodes loop-back after end.

Least-Squared Error Estimate

$\nwarrow V_{\text{Batch}}^{\text{MC}}$ (Dataset set of size N)

Will they converge to
the same thing or
different things?

Value f^N of the MDP to

Value of f^N of the MDP, to
have most-likely generated this
data.

Mauritius 190

~~markedly emaciated~~ ~~and~~ ~~thin~~

Mitospila (OKT)

• $\text{d} \cdot \text{app}$ \rightarrow $\text{d} \cdot \text{app} \text{ dapp}$ \rightarrow dapp

It's time we moved in longer term now.

1998-99
SCHOOL YEAR
BOSTON PUBLIC SCHOOLS

$(\frac{1}{2}f_2) \text{ for } \alpha \neq 0$ and $(\frac{1}{2}) \text{ for } \alpha = 0$.

① Multi-Step Return s^t

4th Oct

notes - h

$$G_{t:t+1} = r^t + \gamma v_t(s^{t+1})$$

$$G_{t:t+2} = r^t + \gamma r^{t+1} + \gamma^2 v_{t+1}(s^{t+2})$$

$$G_{t:t+n} = r^t + \gamma r^{t+1} + \dots + \gamma^n v_{t+n-1}(s^{t+n})$$

$$G_t = \lim_{n \rightarrow \infty} G_{t:t+n}$$

$$v_{t+1}(s^t) \leftarrow v_t(s^t) + \alpha_{t+1} \{ G_{t:t+n} - v_t(s^t) \}$$

$$v_{t+1}(s^t) \leftarrow v_t(s^t) + \alpha_{t+1} \left\{ \frac{G_{t:t+1} + G_{t:t+5}}{2} - v_t(s^t) \right\}$$

The minus sign might be a problem. Need not work. $\leftarrow \times \left\{ 1.5 G_{t:t+2} - 0.5 G_{t:t+3} \right\}$
 $\cdot \cdot \cdot \leftarrow \left\{ G_{t:t+1} + G_{t:t+5} \right\}$

$$\begin{array}{c} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{array} \left\{ \begin{array}{l} G_{t:t+1} \\ G_{t:t+2} \\ G_{t:t+3} \\ \vdots \\ G_{t:t+n} \end{array} \right.$$

Any convex combination of weights would work, coefficients non-negative, sum to 1. $(w_1 v_t + w_2 v_{t+1} + \dots + w_n v_{t+n})$

② λ -return

$$\lambda \in [0, 1]$$

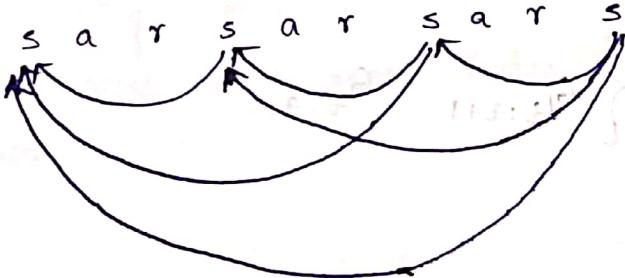
λ -return

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_{t:T}$$

where T is the last time-step in the episode.

$$V_{t+1}(s^t) \leftarrow V_t(s^t) + \alpha_{t+1} \{ G_{t:t+n} - V_t(s^t) \}$$

③ TD(λ) Algorithm



Initialise V

Repeat for each episode

$e \leftarrow 0$ // $e : s \rightarrow \mathbb{R}$, eligibility trace vector

Be born in state s

Repeat for each step of episode

Take action $a = \pi(s)$, obtain r and s'

$$s \leftarrow r + \gamma V(s') - V(s)$$

Analogous to Every-Vist Monte-Carlo: $e(s) \leftarrow e(s) + 1$

of $e(s) \leftarrow 1$
Analogous to First-Vist Monte-Carlo.

For all \bar{s} :

Assuming constant:

$$1. V(\bar{s}) \leftarrow V(\bar{s}) + \alpha \downarrow e(\bar{s})$$

$$2. e(\bar{s}) \leftarrow \gamma \lambda e(\bar{s})$$

$$s \leftarrow s'$$

④ Control

$$1^* \quad a^t = \pi_t(s^t)$$

$$Q_{t+1}(s^t, a^t) \leftarrow Q_t(s^t, a^t) + \alpha_{t+1} \{ r^t + \gamma Q_t(s^{t+1}, a^{t+1}) - Q_t(s^t, a^t) \}$$

$$\pi_t = \pi \Rightarrow \lim_{t \rightarrow \infty} Q_t = Q^\pi$$

$$\pi_t = \epsilon_t - \text{greedy } (Q_{t-1}) \Rightarrow \lim_{t \rightarrow \infty} Q_t = Q^* \quad \begin{matrix} \text{Sarsa} \\ \downarrow \\ \text{on-policy} \end{matrix} \quad \begin{matrix} \text{Handwritten} \\ \text{1996.} \end{matrix}$$

If you follow
on-policy fixed policy,
you'll converge
to Q^π

$$2^* \quad Q_{t+1}(s^t, a^t) \leftarrow Q_t(s^t, a^t) + \alpha_{t+1} \{ r^t + \gamma \max_a Q_t(s^{t+1}, a) - Q_t(s^t, a^t) \}$$

$$\pi_t = \pi \Rightarrow \lim_{t \rightarrow \infty} Q_t = Q^* \quad \begin{matrix} \text{assume } \pi \text{ takes all} \\ s, a \text{ with positive prob.} \end{matrix}$$

$$\pi_t = \epsilon_t - \text{greedy } (Q_{t-1}) \Rightarrow \lim_{t \rightarrow \infty} Q_t = Q^* \quad \text{Not a.}$$

Q-Learning.

Q-Learning is the 1st
Algorithm which has this
guarantee.

Watkins, P.h.D thesis.

off-policy -

Importance? Model-Free!

$\Theta(\|s\| \|A\|)$ space.

3* Expected-Sarsa (Recent)

Only change.

$$Q_{t+1}(s^t, a^t) \leftarrow Q_t(s^t, a^t) + \alpha_{t+1} \left\{ r^t + \gamma \sum_a \pi_t(s^{t+1}, a) Q_t(s^{t+1}, a) \right\}$$

Some behaviours as Sarsa,
less variance than Sarsa because of expectation.

~~same thing as Sarsa~~

Wall off may be

very hard to do in real world
because it may

$\pi_t(s, a)$

$$(Q_t(s^t, a^t))^{\beta} \times \pi_t(s^{t+1}, a^{t+1})^{\beta} + \gamma \{ r^t + (Q_t(s^{t+1}, a^{t+1}))^{\beta} \}^{\beta} \rightarrow (Q_t(s^t, a^t))^{\beta}$$

No wall off means $\beta = 1$
Very interesting the $\beta < 1$ case will be right

but difficult

- good effect for initial values

the initial value is very important, drop value

the more the better the value

drop between old and new

and when the new value is good

drop (the initial) + drop

11th Oct

⑥ TD(λ)

s // visited

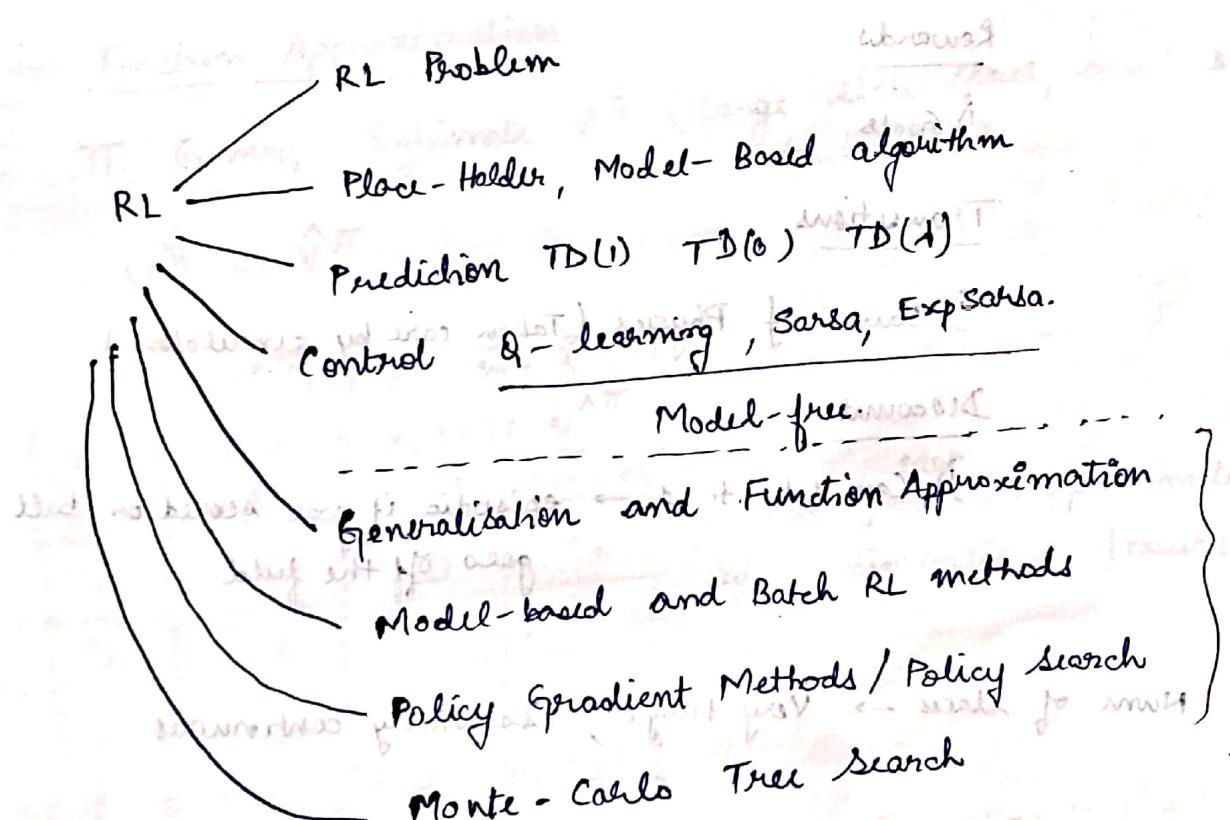
$$e(s) \leftarrow \begin{cases} e(s) + 1 & // \text{Accumulating trace} \rightarrow \text{Every-visit MC} \\ 1 & // \text{Replacing trace} \rightarrow \text{First-visit MC} \end{cases}$$

For $s \in S$: // All states Valid for certain types of α .

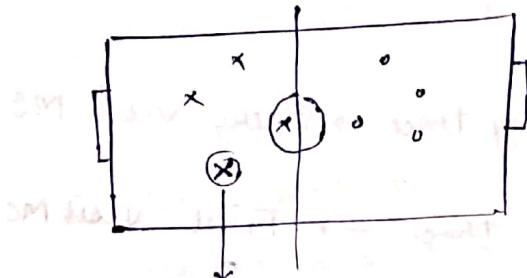
$$v(s) \leftarrow v(s) + \alpha_2 s e(s)$$

$$e(s) \leftarrow \gamma \lambda e(s)$$

① Overview



② Illustration : Soccer



Let us assume that one player

is in our control.

State :

1) Positions of players and football

2) Velocities

3) Score ?

4) Stamina.

Rewards

1) Goals.

Transitions

2) Laws of Physics (Taken care by simulator).

Discounting

3) Can set to 1 → episodic if goal scored or ball goes off the field.

Num of states → Very large, essentially continuous.

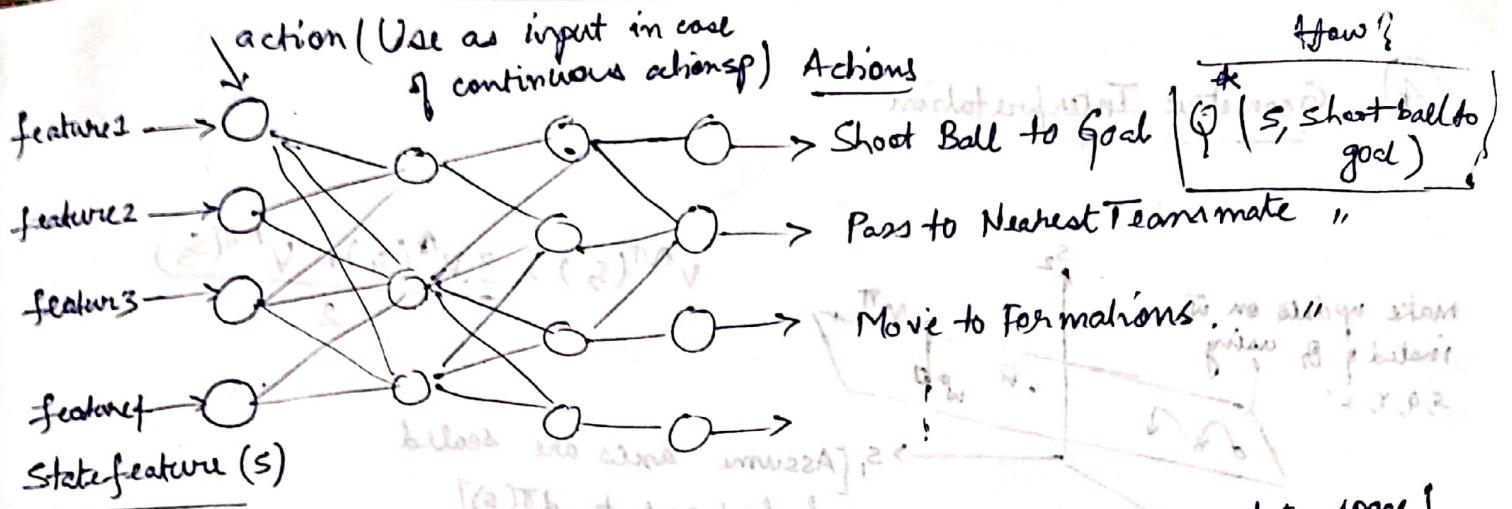
⇒ Extract Features :

1) Which Team Attacking?

2) Position of nearest teammates?

3) Distance to goal?

4) Angle to goal?



Output $\neq Q^*$ exactly? Why? \rightarrow 1) Not looking at entire state space!
Only looking at extracted features.

Optimality x

Convergence?

③ Linear Function Approximation

π Given, estimate V^π (Large state space, can't store in table).

$$V^\pi \approx \hat{V}^\pi$$

$$V^\pi(s) = \bar{w} \cdot \bar{x}(s)$$

V^π	s	$x_1(s)$	$x_2(s)$	V^π
7	s_1	2	-1	$2w_1 - w_2$
2	s_2	4	0	w_1
-4	s_3	2	3	$2w_1 + 3w_2$

(3-state features approximated by 2-dimensional features.)

Which is the "best" (w_1, w_2) ?

$$(w_1^*, w_2^*) = \underset{(w_1, w_2)}{\operatorname{argmin}} \left([7 - (2w_1 - w_2)]^2 + [2 - w_1]^2 + [-4 - 2w_1 + 3w_2]^2 \right)$$

$$\bar{w}^* = \underset{(\bar{w})}{\operatorname{argmin}} \text{MSVE}(\bar{w})$$

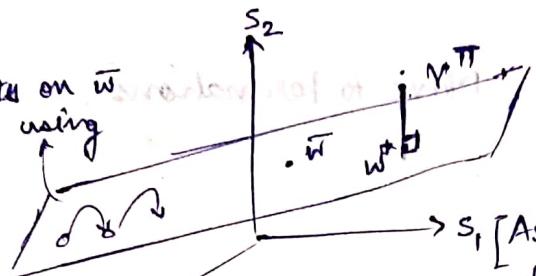
$$\text{MSVE}(\bar{w}) = \frac{1}{3} \sum_{s \in S} d\pi(s) \left\{ V^\pi(s) - \frac{\hat{V}^\pi(s)}{\bar{w} \cdot \bar{x}(s)} \right\}^2$$

(4) Geometric Interpretation

Statement: $\pi(s_2)$ is the average of $\pi(s_1)$ and $\pi(s_3)$

$$V^{\pi}(s_2) = \frac{3V^{\pi}(s_1) + V^{\pi}(s_3)}{2}$$

make update on \bar{w} using
initial θ using
 s, a, r, s'



s_1 [Assume axes are scaled proportional to $d\pi(s)$].

→ weight scale → need to convert that ($1 \rightarrow 3$ plus 3 vectors up to the next level has to proceed plus

\times plateaued

↓ experienced

restomix again, without com-

pute times (loop does ages) π_V starts at 3, moves to π_V

$$\pi_V \approx \pi_V$$

$$(2) \bar{x} \cdot \bar{w} \approx (2)\pi_V$$

intermix again?

(initial) dimensions as pi

$$\frac{d\pi}{dt}$$

Oct - 14th

prabal 217

① Stochastic Gradient Descent

$\bar{w}_0 \leftarrow \text{Initialisation}$

Follow π

For $t = 0, 1, 2, \dots$

$$\bar{w}_{t+1} \leftarrow w_t - \alpha_{t+1} \nabla_{\bar{w}} \left\{ \frac{[(v^\pi(s^t) - \hat{v}^\pi(\bar{x}(s^t), \bar{w}_t))]^2}{2} \right\}$$

$$\leftarrow w_t + \alpha_{t+1} \left\{ v^\pi(s^t) - \hat{v}^\pi(\bar{x}(s^t), \bar{w}_t) \right\} \times \nabla_{\bar{w}} \hat{v}^\pi(\bar{x}(s^t), \bar{w}_t)$$

If linear,

$$\leftarrow w_t + \alpha_{t+1} (v^\pi(s^t) - \bar{x}(s^t) \cdot \bar{w}_t) \cdot \bar{x}(s^t)$$

Option - 1 : Use $\gamma^t + \gamma \gamma^{t+1} + \gamma^2 \gamma^{t+2} \dots$
 True Gradient Descent

Linear TD(1)

Option 2 : Use $\gamma^t + \gamma \hat{v}^\pi(\bar{x}(s^{t+1}), \bar{w}_t')$
 Linear TD(0), semi Gradient Descent.

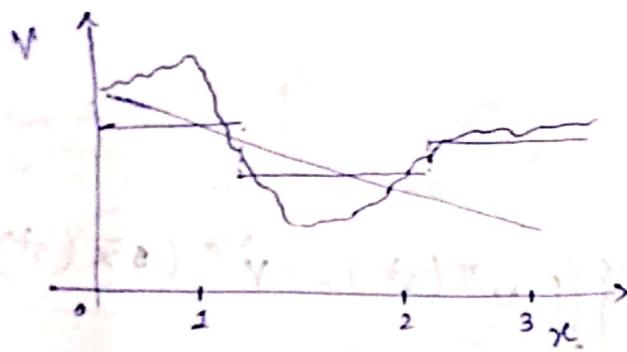
Option - 3 : Use λ - returns
 Linear TD(1).

Eligibility Trace : One parameter for each element of \bar{w} .

$$\text{For TD}(\lambda), \text{MSVE}(\bar{w}_\infty) \leq \frac{1-\gamma\lambda}{1-\gamma} \text{MSVE}(\bar{w}^*)$$

$\tau\text{TD}(\lambda) \rightarrow$ Convergence not proved, but tbf divergence also not proved

② Tile Coding



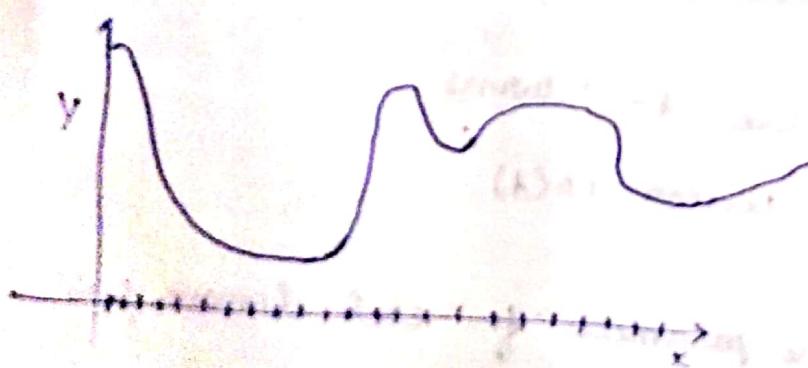
$$\hat{V} = mx + c$$

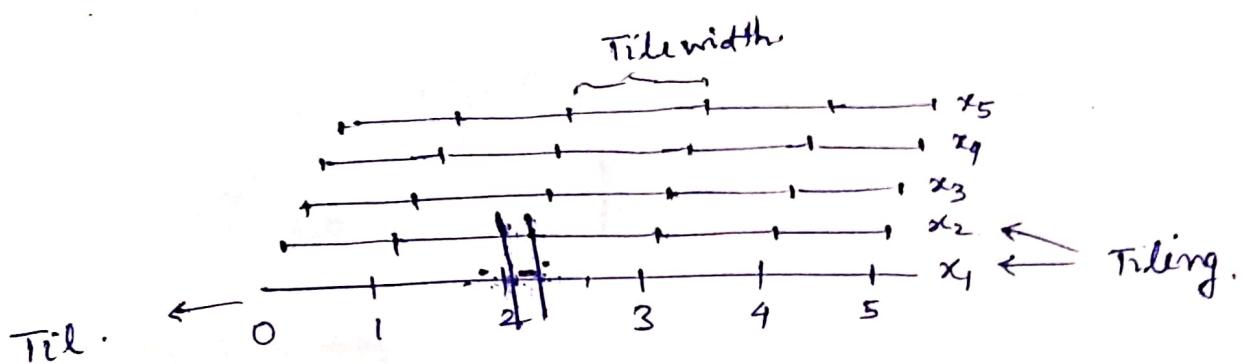
$$x_1 = \begin{cases} 1 & 0 \leq x < 1 \\ 0 & \text{o/w} \end{cases}$$

$$x_2 = \begin{cases} 1 & 1 \leq x < 2 \\ 0 & \text{o/w} \end{cases}$$

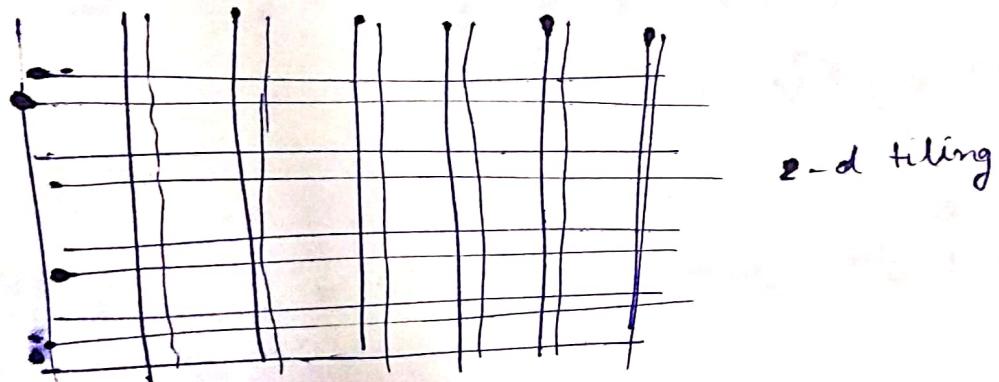
$$x_3 = \begin{cases} 1 & 2 \leq x < 3 \\ 0 & \text{o/w} \end{cases}$$

$$\hat{V}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3$$





$$\text{Resolution} = \frac{\text{Tile width}}{\text{Number of Tilings}}$$



$$\rightarrow \pi^t v^{\pi} TD(\lambda) \bar{w}$$

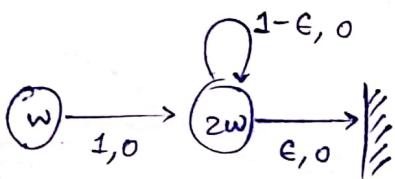
π^t // Control

$$s a \times s' a' (0)$$

$$\bar{w}_{t+1}^a \leftarrow \bar{w}_t^a + \alpha_{t+1} \left\{ r^t + \bar{w}_t^{a'} \times \bar{x}(s') - \bar{w}_t^a \cdot \bar{x}(s) \right\} \alpha$$

① Tsiatskis and Vomkoy's Counterexample

Oct - 18



$$w_0 \rightarrow w_1 \rightarrow w_2 \dots$$

For each state, a "better" estimate of the value

$$\text{is } \mathbb{E}_\pi [r + \gamma \hat{V}^\pi(\bar{x}(s), w)]$$

$$\text{Set } w_{k+1} \leftarrow \underset{w \in \mathbb{R}}{\operatorname{argmin}} \left\{ \mathbb{E}_\pi [r + \gamma \hat{V}^\pi(\bar{x}(s_{\text{next}}, w_k))] - \hat{V}^\pi(\bar{x}(s), w_k) \right\}$$

Function Approx.

Bootstrapping
off-policy update

$$= \underset{w \in \mathbb{R}}{\operatorname{argmin}} \left\{ (2\gamma w_k - w_k)^2 + [(1-\epsilon)2\gamma w_k + \epsilon 0 - 2w_k]^2 \right\}$$

Root-Cause

$$\begin{aligned} \text{Not visiting states according to } d_\pi(s) \} &= \underset{w \in \mathbb{R}}{\operatorname{argmin}} \left\{ w^2 + 4\gamma^2 w_k^2 - 4\gamma w w_k + 4w_k^2 + (1-\epsilon)^2 4\gamma^2 w_k^2 - 8\gamma w w_k (1-\epsilon) \right\} \\ &= \underset{w \in \mathbb{R}}{\operatorname{argmin}} \left\{ 5w^2 + 4\gamma^2 w_k^2 (1+(1-\epsilon)^2) + -4\gamma w w_k (1+2-2\epsilon) \right\} \\ &= \underset{w \in \mathbb{R}}{\operatorname{argmin}} \left\{ 5w^2 - 4\gamma w w_k (3-2\epsilon) + \text{const.} \right\} \end{aligned}$$

$$\Rightarrow 10w = 4\gamma w_k (3-2\epsilon)$$

$$\Rightarrow w = \frac{2}{5} \gamma w_k (3-2\epsilon)$$

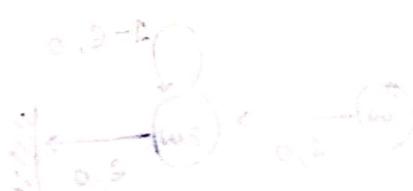
$$\Rightarrow \boxed{w_{k+1} = \frac{2}{5} \gamma w_k (3-2\epsilon)}$$

May Diverge

(2) Approximate Policy Improvement

$$Q^\pi / \theta^{\star\pi} \quad a_1 \quad a_2$$

$$s_1 \quad 3/2.5 \quad 2/3$$



s_2

$$\text{under } \frac{1}{2} \text{ stabilize } \frac{6}{5} \text{ word } \rightarrow \text{state value } v_2$$

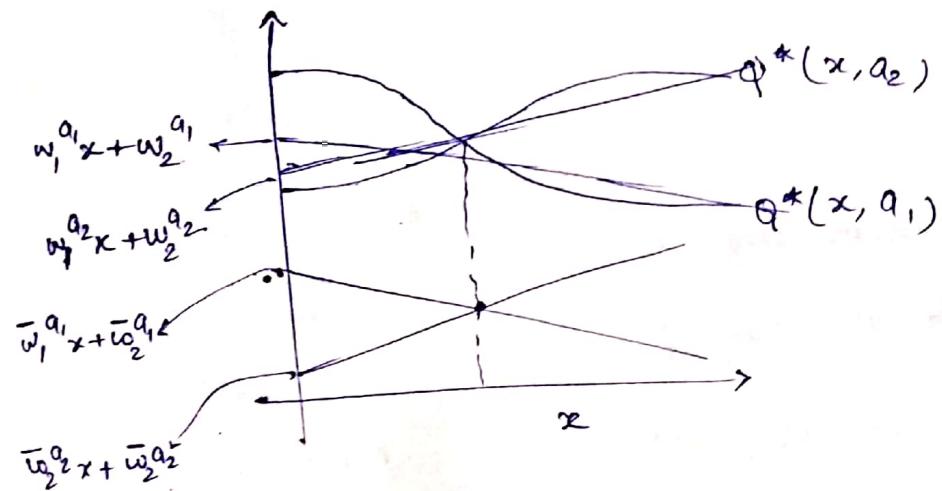
$$\text{Greedy } (Q^\pi) = (a_1, a_2)$$

$$\left\{ \begin{array}{l} \text{Greedy } (Q^{\star\pi}) = (a_2, a_2) \end{array} \right.$$

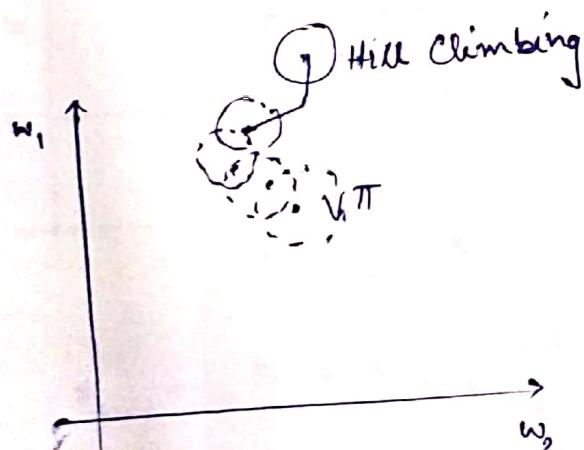
(3) Summary of results

Method	Tabular		Linear FA	Non-Linear FA
	Converges	Stable		
TD(0)	C, exact	+ fast	C, bounded error	NK
TD(λ)	C, exact	+ fast	C, bounded error	NK
TD(1)	C, exact	+ fast	C, best	C, Local Optimal
SARSA(0)	C, opt	+ slow	Chattering	NK
SARSA(λ)	NK	+ slow	NK	NK
SARSA(1)	NK	= slow	NK	NK
Q-Learning(0)	C, opt	+ slow	NK	NK

④ Policy search methods



$$Q = w_1 x + w_2$$



: Evolutionary
Algorithm

: Particle Swarm
Optimisation

: Grid-Search

① Introduction

π_θ

Policy

If $\theta < 6$:

a_1 gets reward of 5

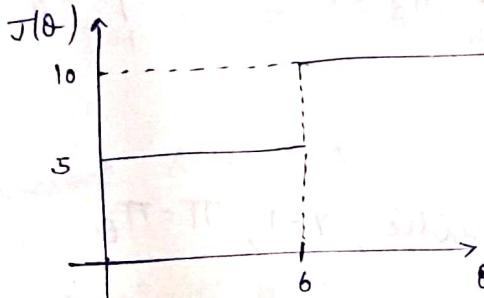
Take a_1 ,

a_2 gets reward of 10

Else

Take a_2 .

$J(\theta)$: reward achieved by π_θ

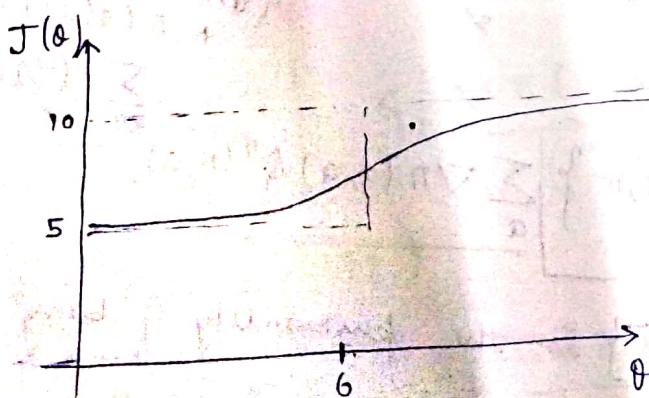


$$\theta_{t+1} \leftarrow \theta_t + \alpha_t \nabla_{\theta} J(\theta) \Big|_{\theta=\theta_t}$$

$$\pi(s, a_1) = \frac{1}{1 + e^{-(6-\theta)}} \quad \left. \begin{array}{l} \text{probabilities of} \\ \text{taking action } a_1, a_2. \end{array} \right\}$$

$$\pi(s, a_2) = 1 - \pi(s, a_1)$$

$$J(\theta) = \pi(s, a_1) \times 5 + \pi(s, a_2) \times 10$$



Start state s_0

$$J(\theta) = V^{\pi_\theta}(s_0)$$

common choice of π_θ

$$\pi_\theta(s, a) = \frac{e^{-\theta \cdot \bar{x}(s, a)}}{\sum_b e^{-\theta \cdot \bar{x}(s, b)}}$$

$$\nabla_\theta \pi_\theta(s, a) = \left\{ \bar{x}(s, a) - \sum_b \pi_\theta(s, b) \bar{x}(s, b) \right\} \pi_\theta(s, a).$$

② Policy Gradient Theorem

Assumptions: $\nabla w.r.t \hat{\theta}$, Episodic, $\gamma=1$, $\pi=\pi_\theta$

$$\begin{aligned} \nabla V^\pi(s) &= \nabla \left[\sum_a \pi(s, a) Q^\pi(s, a) \right] \\ &= \sum_a \left[\nabla \pi(s, a) Q^\pi(s, a) + \right. \\ &\quad \left. \pi(s, a) \nabla \left\{ \sum_{s'} T(s, a, s') (R(s, a, s') + V^\pi(s')) \right\} \right] \\ &= \sum_a \left[\nabla \pi(s, a) Q^\pi(s, a) + \right. \\ &\quad \left. \pi(s, a) \left(\sum_{a'} T(s, a, s') \nabla V^\pi(s') \right) \right] \\ &= \sum_a \left[\nabla \pi(s, a) Q^\pi(s, a) + \pi(s, a) \sum_{s'} T(s, a, s') \left(\sum_{a'} \nabla \pi(s', a') Q^\pi(s', a') \right. \right. \\ &\quad \left. \left. + \pi(s', a') \times \sum_{s''} T(s', a', s'') \nabla V^\pi(s'') \right) \right] \\ &= \sum_{x \in S} \underbrace{\left| \sum_{K=0}^{\infty} P\{s \rightarrow x, K, \pi\} \right|}_{\sum_a \nabla \pi(x, a) Q^\pi(x, a)} \end{aligned}$$

where $P\{s \rightarrow x, K, \pi\}$ is the probability of being
in x after K steps following π starting at s .

$$\nabla J^\pi(s) \propto \sum_x \mu^{\pi,s}(x) \sum_a \nabla \pi(x, a) Q^\pi(x, a)$$

where $\mu^{\pi,s}(x)$ is the stationary of being in x starting from s and following π .

$$\boxed{\nabla J(\theta) \propto \sum_{s \in S} \mu^{\pi,s}(s) \sum_{a \in A} \nabla \pi(s, a) Q^\pi(s, a)}$$

Proportionality constant is expected length of the episode.

Policy gradient theorem

$$\theta_{t+1} \leftarrow \theta_t + \alpha_t \nabla J(\theta)$$

(3) Policy Gradient Algorithm

$$\pi_\theta : s_0 = s^0, a^0, r^0, s^1, a^1, r^1, \dots, s^t, a^t, r^t, \dots$$

$$\begin{aligned} \nabla J(\theta) &\propto \sum_s \underline{\mu^{\pi,s}(s)} \sum_a \nabla \pi(s, a) Q^\pi(s, a) \\ &= E_\pi \left[\sum_a \nabla \pi(s^t, a) Q^\pi(s^t, a) \right] \\ &= E_\pi \left[\sum_a \pi(s^t, a) \left\{ Q^\pi(s^t, a) \frac{\nabla \pi(s^t, a)}{\pi(s^t, a)} \right\} \right] \\ &= E_\pi \left[Q^\pi(s^t, a^t) \cdot \frac{\nabla \pi(s^t, a^t)}{\pi(s^t, a^t)} \right] \\ &= E_\pi \left[Q^\pi(s^t, a^t) \frac{\nabla \pi(s^t, a^t)}{\pi(s^t, a^t)} \right] \end{aligned}$$

REINFORCE (Williams, '92)

$$\begin{aligned} \theta_{t+1} &\leftarrow \theta_t + \alpha_{t+1} g_t \frac{\nabla \pi(s^t, a^t)}{\pi(s^t, a^t)} \\ &= \theta_t + \alpha_{t+1} g_t \nabla \log \pi(s^t, a^t) \end{aligned}$$