

Start Reading Sutton and Barto!

24<sup>th</sup> Sept

## MONTE-CARLO METHODS

### (1) Prediction

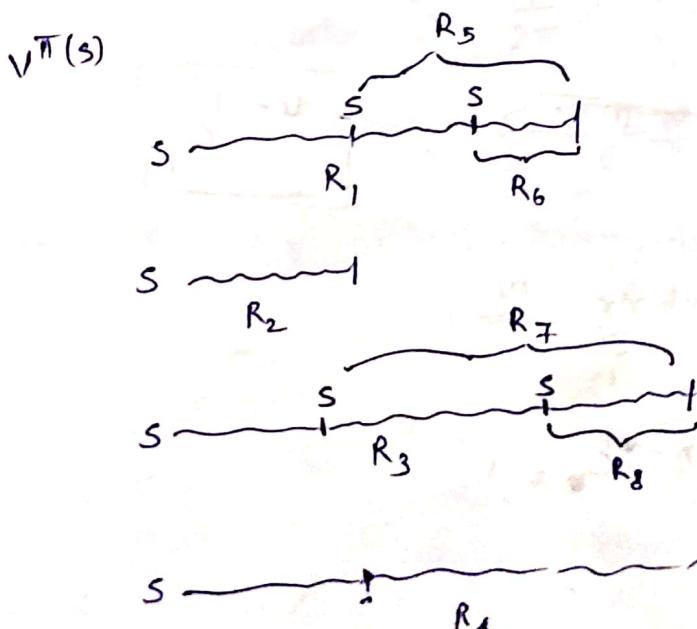
Given  $\pi$ , Find  $v^\pi$  through interaction with MDP.  
Control (Also wants to change and converge to a good policy)

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



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

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

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 1<sup>st</sup> one.  
 $\rightarrow$  EVERY-VISIT MONTE CARLO.

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

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

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  $\pi$ , 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 ( $s \xrightarrow{a}$ ) (Random actions)

Rewards b/w  $R_{\min}$  &  $R_{\max}$ .

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

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

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

$$p = 2p^2$$

## ② $\pi$ $V^n$ $s$

$$s \sim R_1$$

$$s \sim R_2$$

$$\sim R_3$$

$$\sim R_n$$

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

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

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

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

$\downarrow$   
Form of  
Full Bootstrapping

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

Temporal Difference prediction error.

#### ④ TD(0) Algorithm

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

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<sup>t+1</sup> s<sup>t+1</sup> s<sup>t+1</sup>~~

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

## 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

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

Will they converge to  
the same thing or  
different things?

Value f<sup>N</sup> of the MDP to

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

start 23 km

~~marked by a small white mark~~

• the best way to keep dark  $\rightarrow$  it  
is not in water, but around a

$\{f_n\}_{n=1}^{\infty}$  is not uniformly bounded.

$$(\mathbb{H}^2) \otimes V + (\mathbb{H}^2) \otimes (\mathbb{H}^2) \rightarrow (\mathbb{H}^2) \otimes V \xrightarrow{\cong} V.$$

# 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) $\lambda$ -return

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

$\lambda$ -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( $\lambda$ ) 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  $s'$

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

Analogous to Every-Vist Monte-Carlo

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

For all  $\bar{s}$ :

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

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

Assuming constant.

#### 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  $Q^\pi$

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

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

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

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  
in some cases  $\beta = 0$  and  $\beta > 1$  will be wrong

but difficult

- good effect for both with variance of

the  $\beta$  and  $\beta < 1$  will have higher variance  
but  $\beta > 1$  will have lower variance

higher variance can make

more exploration - more variance

worse (the variance)

11<sup>th</sup> Oct

## ⑥ TD(λ)

s // visited

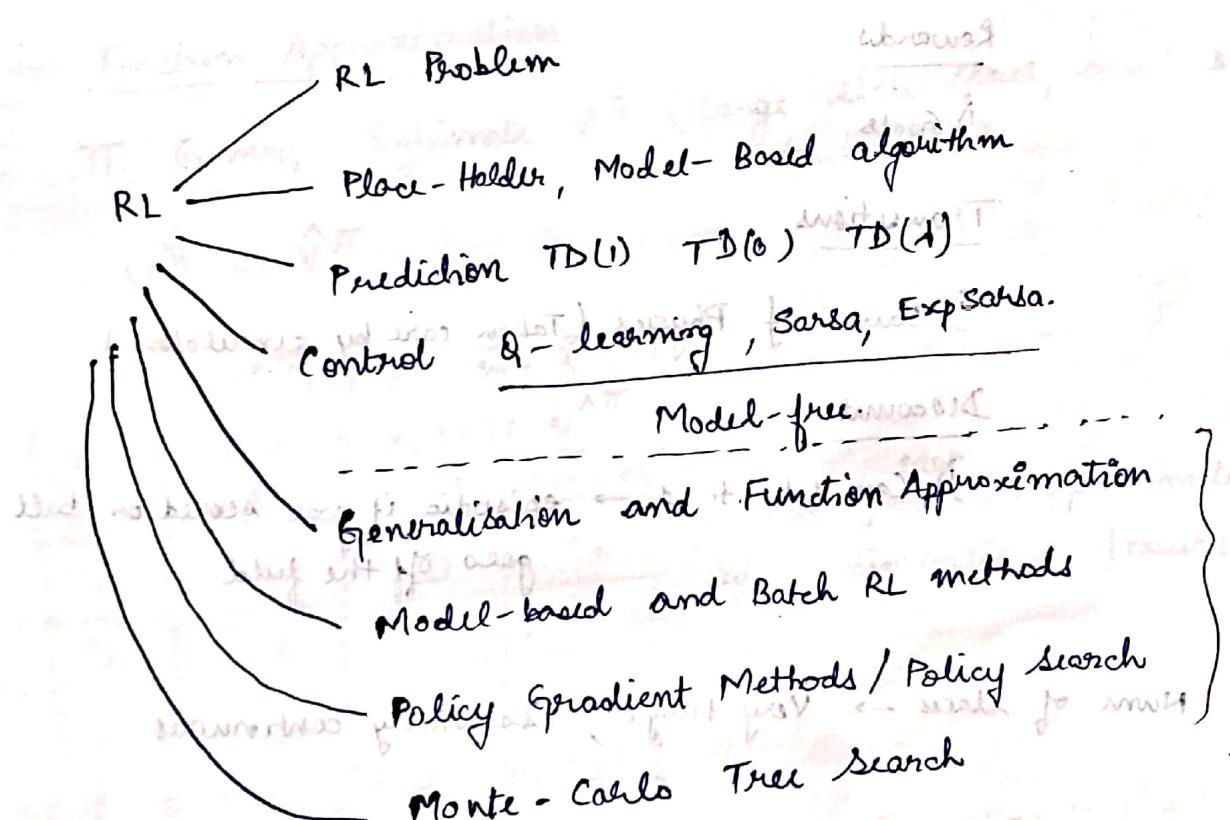
$$e(s) \leftarrow \begin{cases} e(s) + 1 & // \text{Accumulating trace} \rightarrow \text{Every-visit MC} \\ 1 & // \text{Replacing trace} \rightarrow \text{First-visit MC} \end{cases}$$

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

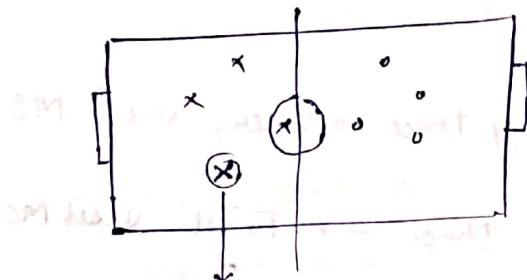
$$v(s) \leftarrow v(s) + \alpha_2 s e(s)$$

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

## 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.

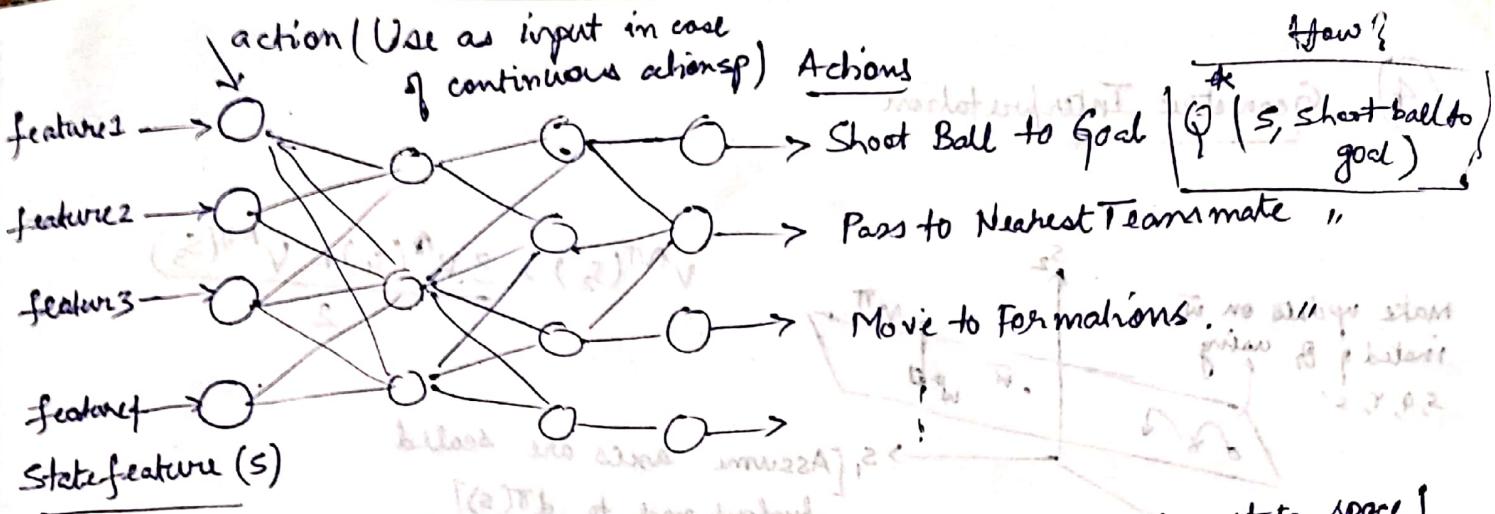
⇒ Extract Features :

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 ✗

Convergence?

### ③ Linear Function Approximation

$\pi$  Given, estimate  $v^\pi$  (Large state space, can't store in table).

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

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

$v^\pi$	$s$	$x_1(s)$	$x_2(s)$	$v^\pi$
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"  $(w_1, w_2)$ ?

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

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

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

#### (4) Geometric Interpretation

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

make update on  $\bar{w}$  using  
initiated  $\theta$  using  
 $s, a, r, s'$

Assume axes are scaled  
proportional to  $d\pi(s)$ .  
I want state  $s_2$  to get updated (1 -> 2 plus 3 times of 3).  
so that it has to proceed plus

$\times$  plateaued

Set experienced

reinforcement reinforcement

note that  $\theta$  does not affect  $V^{\pi}$  (not in  $\theta$ )

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

$$(2) \hat{x} \cdot \bar{w} \approx (2)V^{\pi}$$

intermediate

(state) dimensions as pi

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

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

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

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

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

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

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

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

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

$$(x_{new} - \hat{x}) + \gamma \left[ r + \beta \hat{V}^{\pi} \left( (w - \bar{w})^T \hat{x} \right) \right] \text{ where } \beta = ((w - \bar{w})^T \hat{x})$$

$$= (x_{new} - \hat{x}) + \gamma \left[ r + \beta \hat{V}^{\pi} \left( (w - \bar{w})^T \hat{x} \right) \right]$$

$$= (x_{new} - \hat{x}) + \gamma \left[ r + \beta \hat{V}^{\pi} \left( (w - \bar{w})^T \hat{x} \right) \right]$$

Oct - 14th

probabilistic MLE

## ① Stochastic Gradient Descent

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

Follow  $\pi$

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

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

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.

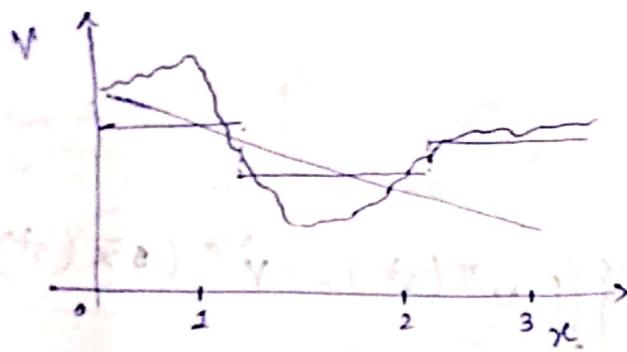
Option - 3 : Use  $\lambda$  - returns  
Linear TD(1).

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

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

TD(1)  $\rightarrow$  Convergence not proved, but no divergence also not proved

## ② Tile Coding



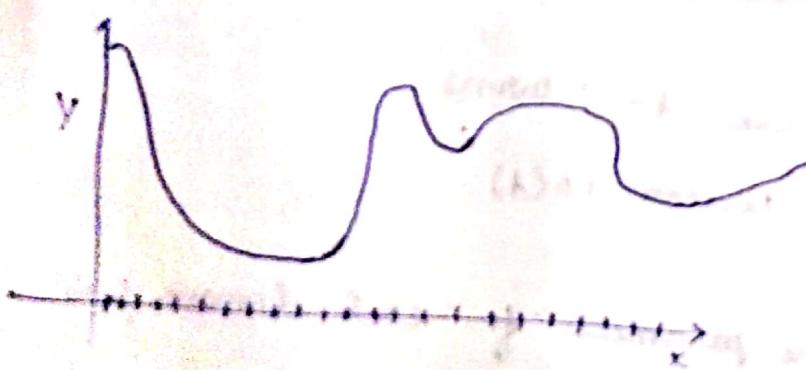
$$\hat{V} = mx + c$$

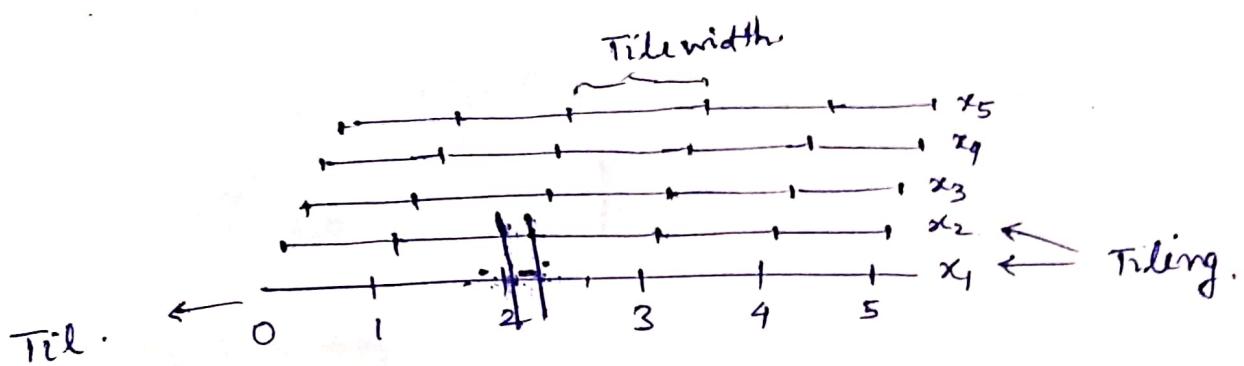
$$x_1 = \begin{cases} 1 & 0 \leq x < 1 \\ 0 & \text{o/w} \end{cases}$$

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

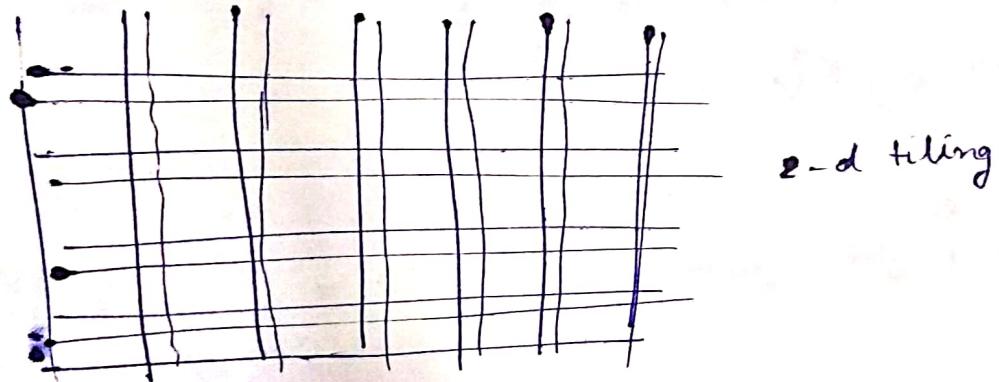
$$x_3 = \begin{cases} 1 & 2 \leq x < 3 \\ 0 & \text{o/w} \end{cases}$$

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





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



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

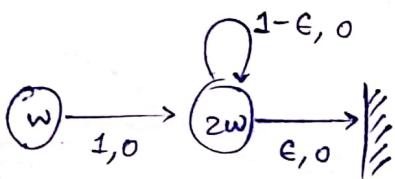
$\pi^t$  // Control

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

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

# 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

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

Root-Cause

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

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

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

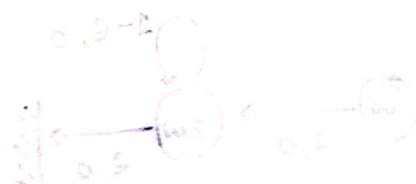
$$\Rightarrow \boxed{w_{k+1} = \frac{2}{5} \gamma w_k (3-2\epsilon)}$$

May Diverge

## (2) Approximate Policy Improvement

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

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



$s_2$

$$\text{under } \frac{1}{2} \text{ stabilize } \frac{6}{5} \text{ word } \rightarrow \text{state value is } 0.7$$

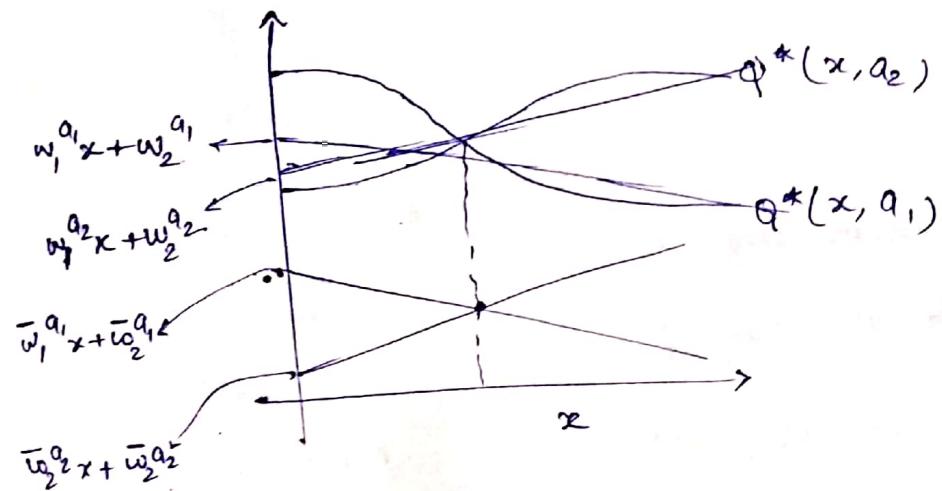
$$\text{Greedy } (Q^\pi) = (a_1, a_2)$$

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

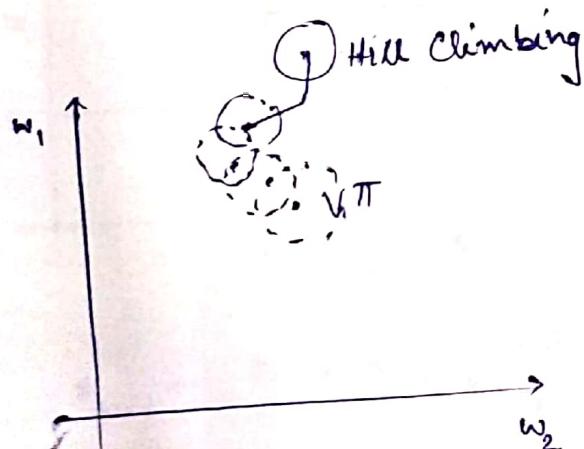
## (3) Summary of results

Method	Tabular		Linear FA	Non-Linear FA
	Time	Space		
TD(0)	C, exact	C, bounded error	NK	NK
TD(λ)	C, exact	C, bounded error	NK	NK
TD(1)	C, exact	C, best	C, Local Optimal	NK
SARSA(0)	C, opt	Chattering	NK	NK
SARSA(λ)	NK	NK	NK	NK
SARSA(1)	NK	NK	NK	NK
Q-Learning(0)	C, opt	NK	NK	NK

## ④ Policy search methods



$$Q = w_1 x + w_2$$



: Evolutionary  
Algorithm

: Particle Swarm  
Optimisation

: Grid-Search

## ① Introduction

$\pi_\theta$

Policy

If  $\theta < 6$ :

$a_1$  gets reward of 5

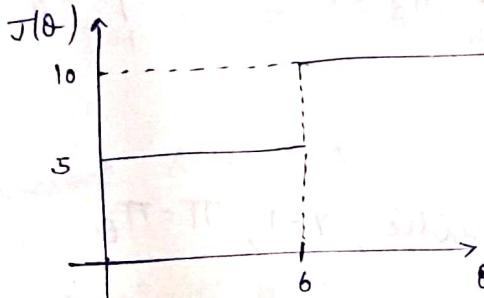
Take  $a_1$ ,

$a_2$  gets reward of 10

Else

Take  $a_2$ .

$J(\theta)$ : reward achieved by  $\pi_\theta$

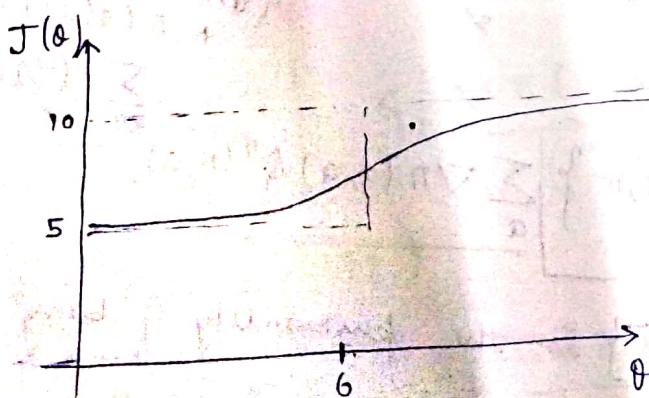


$$\theta_{t+1} \leftarrow \theta_t + \alpha_t \nabla_{\theta} J(\theta) \Big|_{\theta=\theta_t}$$

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

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

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



Start state  $s_0$

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

common choice of  $\pi_\theta$

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

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

## ② Policy Gradient Theorem

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

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

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

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

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

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

Proportionality constant is expected length of the episode.

Policy gradient theorem

$$\theta_{t+1} \leftarrow \theta_t + \alpha_2 (\nabla J(\theta))$$

### (3) Policy Gradient Algorithm

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

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

$$= E_\pi \left[ Q^\pi(s^t, a^t) \cdot \frac{\nabla \pi(s^t, a^t)}{\pi(s^t, a^t)} \right]$$

$$= E_\pi \left[ G_t \frac{\nabla \pi(s^t, a^t)}{\pi(s^t, a^t)} \right]$$

REINFORCE (Williams, '92)

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

$$\textcircled{1} \quad \nabla J(\theta) = \nabla V^\pi(s^0) = \sum_{x \in S} \sum_{k=0}^{\infty} P_h \{ s_k \rightarrow x, a_k, \pi \} \sum_a \nabla \pi(x, a) \hat{Q}^\pi(s_k, a)$$

Follow  $\pi$ , starting at  $s^0$

$$s_0 = s^0 a^0 r^0 s^1 a^1 r^1 - \dots - s^{T-1}$$

$$\text{Return Estimate} = f(s^0) + f(s^1) + \dots + f(s^{T-1})$$

$$IE[\text{Estimate}] = \sum_{k=0}^{\infty} IE[f(s^k)]$$

Loop forever:

Generate episode  $e = s^0 a^0 r^0 \dots r^{T-1} s^T$ , following  $\pi_\theta$

~~for~~ Loop for  $t = 0, 1, \dots, T-1$

$$g \leftarrow \sum_{k=t}^{T-1} y^{k-t} r^k$$

$$\theta_{\text{new}} \leftarrow \theta_{\text{new}} + \gamma^t g \nabla \ln \pi_\theta(s^t, a^t)$$

Not  $\theta_{\text{new}}$   
(as in SB)

①

Baseline

$Q^\pi$

$$Q^\pi = \sum_a \pi(s, a) Q^\pi(s, a)$$

$a_1$	$a_2$	$a_3$
$s_1$	105	79

$$s_2 \quad 10 \quad 6 \quad 13 \quad 12$$

$$s_3 \quad -50 \quad -60 \quad -50 \quad -55$$

$$\nabla \pi(s^t, a^t) \cdot Q^\pi(s^t, a^t)$$

Might be high.  $\rightarrow$  cause high variance

② Actor

Critic

Actor

Let  $b: S \rightarrow \mathbb{R}$  be any function of state

$$\begin{aligned}\nabla J(\theta) &\propto \sum_s \mu^{\pi, s_0}(s) \sum_a \nabla \pi(s, a) \partial^{\pi}(s, a) \\ &= \sum_s \mu^{\pi, s_0}(s) \sum_a \nabla \pi(s, a) \left\{ \partial^{\pi}(s, a) - b(s) \right\} \\ &\quad \text{(fixing } b \text{)} \\ &\sum_s \mu^{\pi, s_0}(s) \sum_a \nabla \pi(s, a) b(s) \\ &= \sum_s \mu^{\pi, s_0}(s) b(s) \nabla \left( \sum_a \pi(s, a) \right) = 0.\end{aligned}$$

If you don't know something, what do you do in reinforcement learning? Stick in an estimator.

**REINFORCE** with baseline

$$\theta_{\text{new}} \leftarrow \theta_{\text{new}} + \alpha \sum_{t=0}^{\infty} \gamma^t (G_t - b(s_t)) \nabla \ln \pi_{\theta}(s_t, a_t)$$

$\beta$  the advantage  $G_t - b(s_t)$

$\gamma$  discount factor

$\nabla \ln \pi_{\theta}(s_t, a_t)$ , which is itself estimated (empirically).

## ② Actor-Critic Methods

Critic Update  $w$  such that  $\hat{V}_w$  is a good estimate of the value function of the current policy.

$$\text{Actor: } \theta_{\text{new}} \leftarrow \theta_{\text{new}} + \alpha \sum_{t=0}^{\infty} \gamma^t (r_t + \gamma \hat{V}_w(s_{t+1}) - \hat{V}_w(s_t)) \times \nabla \ln \pi_{\theta}(s_t, a_t).$$

$s \xrightarrow{a} r s' a'$   $\pi$   $s a r$  length of trajectory = sample complexity  
 $Q(s, a) \leftarrow Q(s, a) + \alpha \{ r + \gamma \max_a Q(s', a) - Q(s, a) \}$   
 Computation complexity  $\uparrow$   $(\text{length of trajectory})^2$   $\propto$   $\sum_{s,a} (\pi(s,a))^2$

### ③ Batch RL

$$Q \leftarrow 0$$

Repeat for ever:

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

Follow  $\pi$  for  $N$  episodes, gather data  $D$   $\{ (s_i, a_i, r_i, s_{i+1}) \}_{i=1}^N$

$$D = (s_i, a_i, r_i, s_{i+1})^L$$

$$Q \leftarrow \text{BatchUpdate}(D, Q) \quad // \text{Optional to use } Q.$$

E.g.

Batch Update Experience Replay ( $D, Q$ )  
Followed by  $M$  iterations

Repeat  $M$  times

Pick  $i \in L$  Uniformly at random

$$Q(s_i, a_i) \leftarrow Q(s_i, a_i) + \alpha \{ r + \gamma \max_a Q(s'_i, a) - Q(s_i, a) \}$$

Return  $Q$ .

$$(t_0, f_0) \xrightarrow{\pi} (t_1, f_1) \xrightarrow{\pi} \dots \xrightarrow{\pi} (t_M, f_M)$$

## Batch Update Fitted Q Iteration (D)

$Q_0 \leftarrow 0, i \leftarrow 0$

Assume predictor  $\leftarrow$  Supervised Learning  $((x_j, y_j)^\top)$

Repeat H times:

for  $j = 1$  to L

$\bar{x}_j = \text{featureVector}(s_j, a_j)$

$y_j = r_j + \gamma \max_{a'} Q_i(s_{j+1}, a')$

$Q_{i+1} \leftarrow \text{Supervised Learning } (\bar{x}_j; y_j)$

$i \leftarrow i + 1$

Return  $Q_i$ .

Nov -1

① Models  $S, A, T, R, \gamma$

$T, R \rightarrow \text{Model}$

Distributional Models :  $T(s, a, s')$  is known/given  $\rightarrow D$

Sampling Model :  $s' \sim T(s, a)$

sampling operator

Planning

More powerful,

Not always available.

systems where

set of states too large, or  
complex dynamics to model  
the probability exactly.

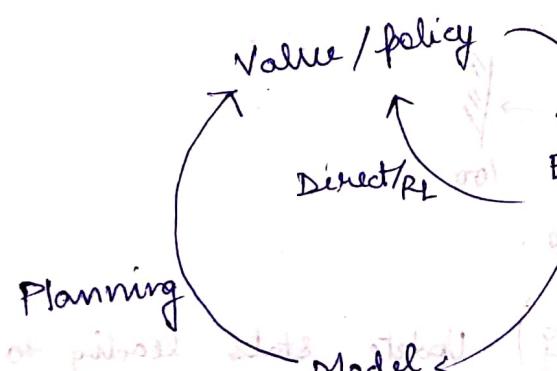
Board games, which involve  
a roll of dice. Distribution  
of states very difficult as  
compared to just sampling  
a new state.

② Concept / Idea

Experience  $\rightarrow$  Value Function  $\rightarrow$  Policy

Model  $\rightarrow$  simulated experience  $\rightarrow$  value function  $\rightarrow$  policy

Model  $\rightarrow$  simulated experience



- Model  $\Rightarrow$
- 1) Less experienced
- 2) More computation
- 3) Useful in a changing world

estimate of future state

beginning and when value

Thinking : Fast and Slow

Dyna - Q

Initialise  $Q(s, a)$ ,  $Model(s, a)$

Loop forever

Direct RL

$s \leftarrow$  current state

$a \leftarrow \epsilon\text{-greedy}(a, s)$

Take action  $a$ . Get next state  $s'$ , reward  $r$ .

$Q(s, a) \leftarrow Q(s, a) + \alpha \{ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \}$

Update Model ( $Model, s, a, r, s'$ )

Loop  $n$  times:

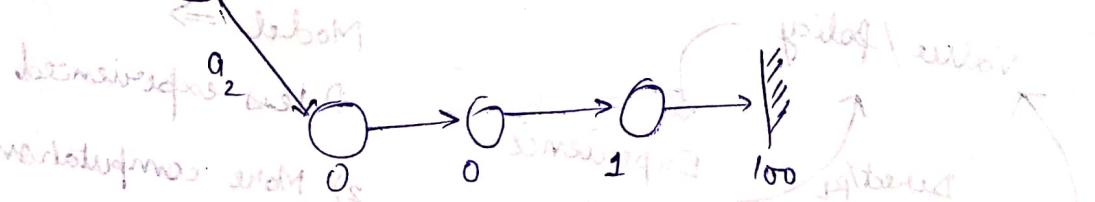
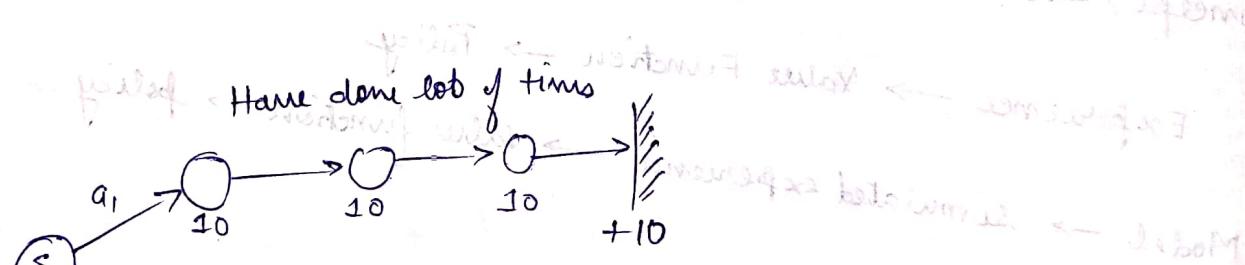
Simulation

$s \leftarrow$  random previously observed state

$a \leftarrow$  random previously taken action from  $s$

$r, s' \leftarrow Model(s, a)$

$Q(s, a) \leftarrow Q(s, a) + \alpha \{ r + \gamma \max_{a'} Q(s', a') - Q(s, a) \}$

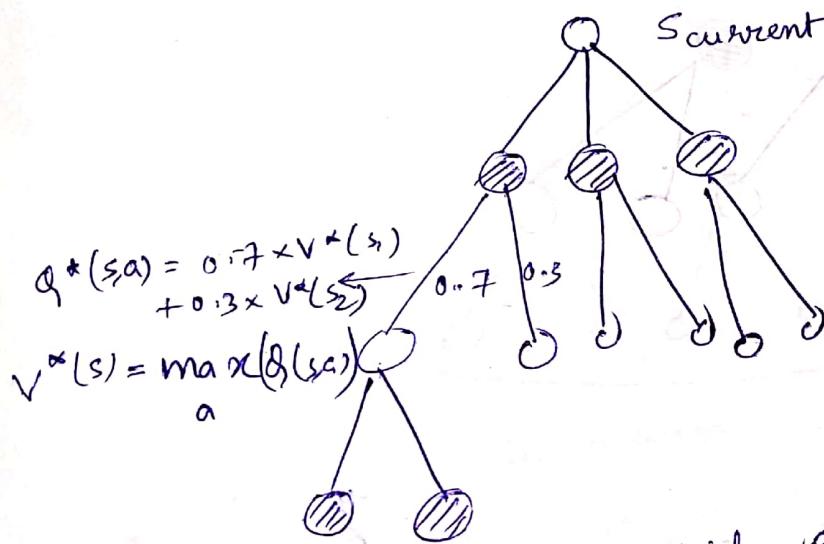


Then this happens.

Prioritised Sweeping: Update stats leading to states whose values have changed the most.

## ④ Action Selection at run time

### Forward simulation / Tree Search

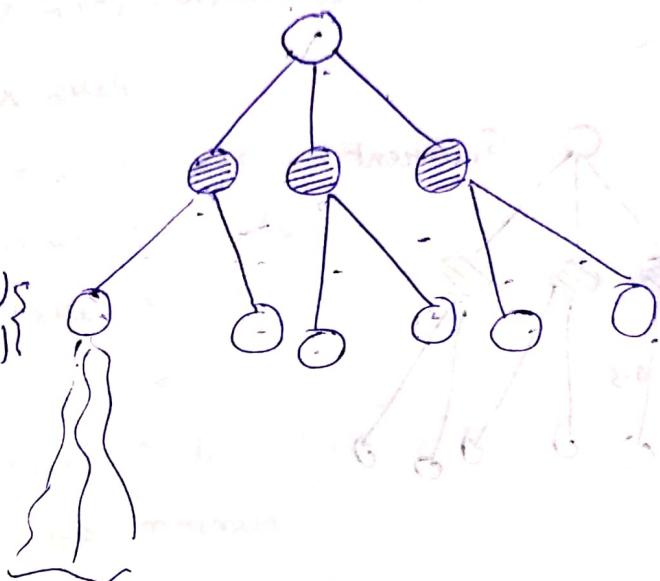


Suppose we want to pick an action. we need a tree  
 $\epsilon$ -optimal action. we need a tree  
 with depth  $f\left(\frac{1}{\epsilon}, \frac{1}{\gamma}\right)$

b size of tree

## ① Roll-out Algorithms

Ideal to use  $V^*(s)$   
Practical to use  $\sim V^\pi(s)$



Approximate with

$V^\pi(s)$ , where  $\pi$  is  
a good policy. ~~random policy~~

In practice, use rollout  
after expansion to a

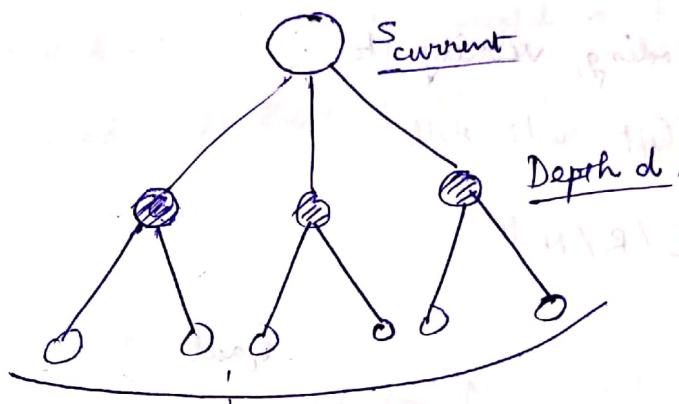
reasonable depth. Hope that

a randomly good policy  $\pi$  is  
also able to distinguish b/w actions.

Problem: Expanding  
lower Q values  
same no. of times as  
higher Q values.

## ② UCT (MCTS)

Happening at run time



- Aim is to find  $\underset{a}{\operatorname{argmax}} Q^*(s_{\text{current}}, a)$ .
- Store  $(s, a)$  pairs reachable in  $d$  steps in table/memory.
- For each pair maintain:
  - $Q(s, a)$ , the empirical average of cumulation reward based on max, expectation, rollouts.
  - $UCB(s, a) = Q(s, a) + C_p \sqrt{\frac{\ln(t)}{n(s, a)}}$ 
    - No. of times you have called your simulator.
    - No. of visits to  $(s, a)$ .
  - when in state  $s$ , sample action  $a = \underset{a}{\operatorname{argmax}} UCB(s, a)$ .
  - when at leaf, follow roll-out policy.
  - Update values along simulate path.

Non-stationary handled by large values of  $C_p$ .

Focus attention on "high-value" states.

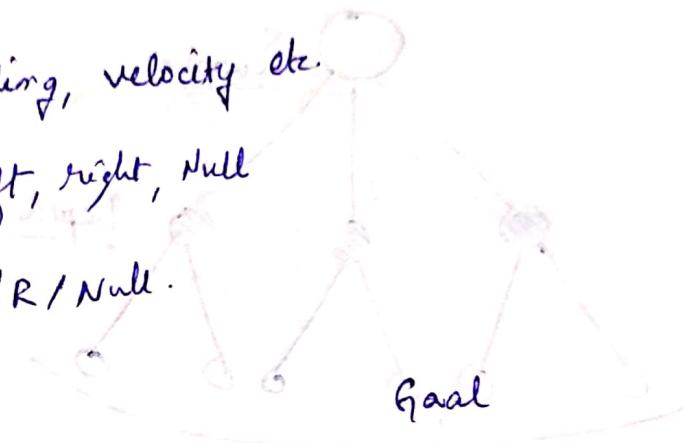
### ③ Reward Shaping

Bicycle

State: Angle, heading, velocity etc.

Action: Turn, left, right, null

Reward L/R/Null



start

→ For every time step.

$$\frac{1}{R}$$

at every time step.

Good reward (e.g. 100) Bad reward (e.g. -100)

Goal

$v_g$

Start

$$r_t = a + b \gamma^t (s^t)$$

if reward (e.g. 100) if bad reward (e.g. -100)

(0.2) at chipping off

(0.2) at hitting a dot in water

(0.2) at water = 0 marks

operating two-blade water jet to water

May think

prev. water level

$$M = SATR^\gamma \quad \} \text{ Under what conditions}$$

$$M' = SATR'^\gamma \quad } \text{ will an } \text{oppo} \text{ optimal policy of } M \text{ also be optimal for } M'?$$

Definition:  $F: S \times A \times S \rightarrow R$  is called a potential-based functioned if there exists  $\phi: S \rightarrow R$  such that  $\forall s \in S, a \in A, F(s, a, s') = \gamma \phi(s') - \phi(s)$

Theorem (Ng, M —, Russel, 95).

If  $\pi^*$  is an optimal policy for  $M = (S, A, T, R, \gamma)$ , then for envy potential based shaping  $f^N F$ ,  $\pi^*$  is also an optimal policy for  $M_F = (S, A, T, R + F, \gamma)$

Necessity

If  $F$  is not a potential-based shaping  $f^N$ , then there exists  $N = SATR^\gamma$ ,  $\pi^*$  such  $\pi^*$  is an optimal policy for  $M$ ,

$\pi^*$  is not an optimal policy  $f^N$ .

$$M = (S, A, T, R + F, \gamma)$$