

Chapter 1 :- Introduction to Reinforcement learning

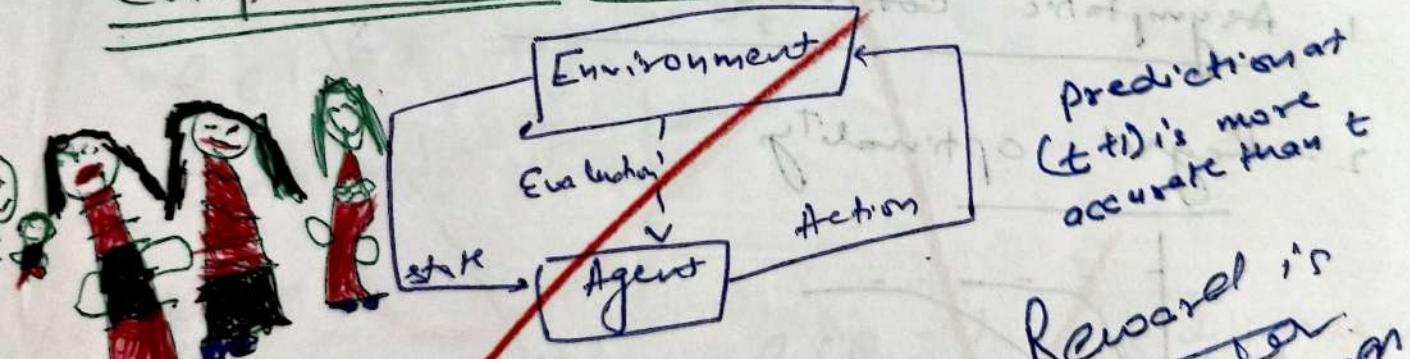
learning

Learn through system by interacting with system

what is reinforcement learning

- Learning about stimuli and action based on rewards and punishments alone.
- No detail supervision available
- Trial and error learning
- Delayed reward
- Sequence of actions required to obtain reward
- Associate learning required
 - Need to associate action to states
- Learn about policies not just actions
- Typically in a stochastic world.

Chapter 2 :- RL framework and application



~~prediction at $(t+1)$ is more accurate than t~~

~~Reward is scalar~~

- Learn from close interaction
- stochastic environment
- Noisy delayed scalar evaluation.
- Maximizing a measure of long term performance

Temporal

CAE, (2, 3)

Chapter 3

Immediat Reinforcement

Learning problem

long the way

outcome is immediate

e.g. dog training

t

R_t

a_t

Exploration

Exploitation

dilemma

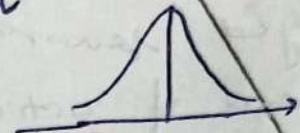
Reward
Payoff
Cost

something

$$A = \{1, \dots, n\}$$

for each of the action $q_{\pi}(a)$

$q_{\pi}(a)$ denotes payoff
of action a



Chapter 4:

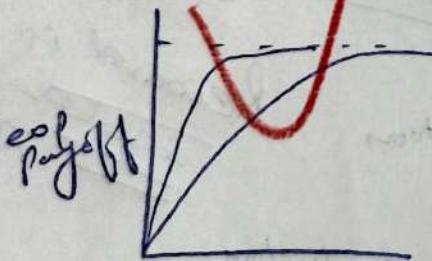
Bandit Problems

multi arm Bandit Problem

Solution?

1. Asymptotic correctness

2. Regret optimality



(2) PAC optimality

absolute
guarante

$(\mathcal{B}, \delta) \sim \text{PAC}$

Probability
Probability
Approximately

Probability

$$P(q_{\pi}(a) > (1-\epsilon) q^*(a^*)) \geq (1-\delta)$$

with
some
probability

$$P(q_{\pi}(a) > q^*(a^*) - \epsilon) \geq (1-\delta)$$

Chapter 5 Value function Based Methods (Q_t)

- where we do not know q^* and we try to predict q^* using indicator function

$$Q_t(a) = \frac{1}{t} \sum_{i=1}^t I(a_i = a) \cdot R_{t+i} \quad \text{This give average}$$

$$\sum I(a_i = a)$$

Greedy $\arg \max_a Q_t(a) \rightarrow$ take action that give maximum exploitation

Exploration Policy

ϵ -greedy

$P_t(1-\epsilon)$ argmax $Q_t(a)$
 \therefore uniformly from

get exploration

46-3

~~softmax~~

$$\frac{Q_t(a)}{\sum_{b=1}^n Q_t(b)}$$

$$= P_t(Q_t^{-1})$$

if difference is very very small
then softmax used

$$\frac{e^{Q_t(a)/\beta}}{\sum_{b=1}^n e^{Q_t(b)/\beta}}$$

β is temperature parameter

$$P(a_t | a^*) = 1 - \epsilon + \epsilon / A$$

$$\hat{Q}_t(a) = \frac{\hat{Q}_{t-1}(a)_{\text{no}} + R_t \underset{\text{if } a_t = \text{act}}{\cancel{Q_t(a)}}}{n_{at+1} \underset{\text{if } a_t = \text{act}}{(a_t = a_{\text{act}})}}$$

$$\begin{aligned} b^{a_t = a_{\text{act}}} \\ &= \frac{\hat{Q}_t(a_{\text{act}})_{\text{no}} + R_t}{n_{at+1}} \end{aligned}$$

$$= \hat{Q}_{t-1}(a) \left(1 - \frac{1}{n_{at+1}} \right) + \frac{R_t}{n_{at+1}}$$

$$= \hat{Q}_{t-1}(a) + \frac{1}{n_{at+1}} (R_t - \hat{Q}_{t-1}(a))$$

Non stationary problems

$$\hat{Q}_t(a) = \hat{Q}_{t-1}(a) + \alpha [R_t - \hat{Q}_{t-1}(a)]$$

Chapter 6 UCB1 exploration

MA Bandit problem

Upper confidence Bound (UCB1)

There are K arms

0/1 Just play each arm once

loop play arm j that maximizes

$$\text{Expected Payoff} = Q_i + \sqrt{\frac{2 \ln n}{\text{no. of trials}}}$$

no. of trials

Expected Payoff

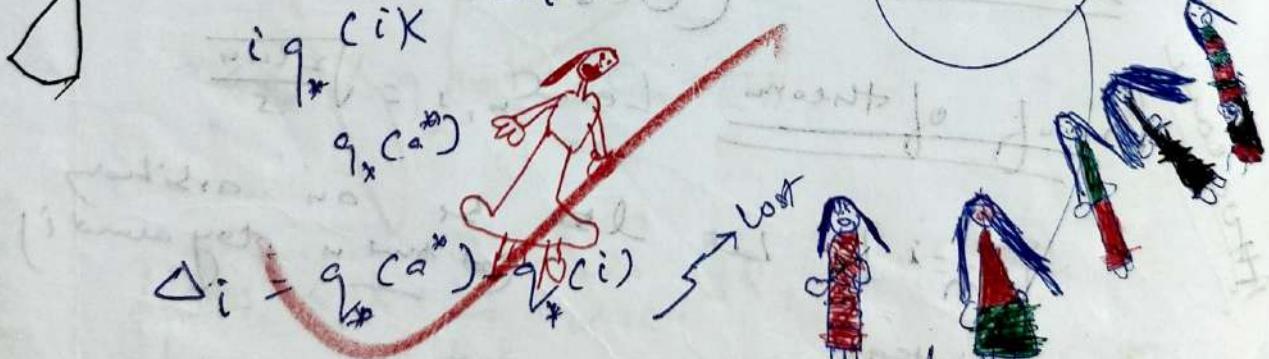
Chapter 7: $\xrightarrow{\text{concentration}}$ Bound
 $\xrightarrow{\text{fairly trivial}}$ arm index

Theorem: for all $K > 1$, if UCB1 is

run on K arms having arbitrary reward distribution, with support in $[0, 1]$

then its expected regret after any number of plays n is at most magic constant

$$8 \sum_{i=1}^K \left(\frac{\ln n}{T_i} \right) + C \left(\frac{n}{K} \right)^2 \left(\sum_{j=1}^K \Delta_j \right)$$



~~no. of time arm i played in first n trials~~
 $T_i(n)$ no. of time arm i played in first n trials

$$\text{Regret}_n = \sum_i E[T_i(n)] \Delta_i$$

$$E[X_{i,n}] = q^*(i)$$

we show that

$$E[\tau_j(n)] \leq \frac{8}{\Delta_j^2} \ln n + c$$

Cheinoff - Hoeffding bound

Let x_1, \dots, x_n be r.v with common range

$$[0, 1] \text{ s.t. } E[x_{t+1}, x_{t-1}] = \mu$$

let $s_n = \frac{x_1 + \dots + x_n}{n}$, then $\forall \epsilon \geq 0$

$$P(s_n \geq \mu + \epsilon) \leq e^{-2\epsilon^2 n}$$

$$P(s_n \leq \mu - \epsilon) \leq e^{-2\epsilon^2 n}$$

Chapter 8 Theorem 1 proof
(UCB1 theorem)

Proof of theorem

$$\text{Let } c_{n,i} = \sqrt{\frac{\ln n}{s}}$$

$\{I_{n=i}\}$ let I_m be an arbitrary
(at time m play arm i)

an integer

$$\tau_i(n) = 1 + \sum_{m=K+1}^n \{I_{m=i} \mid \tau_i(m) \geq l\}$$

$$\leq l + \sum_{m=K+1}^n \{I_{m=i} \mid \tau_i(m) \geq l\}$$

$$\leq l + \sum_{m=k+1}^n \left\{ Q(a^*) (c_m, T_{a^*}(m-1)) \leq Q(a^*) \right. \\ \left. + (c_{m-1}, T_i(m-1)) \geq R \right\}$$

$$\leq l + \sum_{m=k+1}^n \left\{ \min_{0 \leq s \leq m} (Q_s(a^*) + c_{m-1}, T_s(a^*)) \right. \\ \left. \leq \max_{1 \leq s_i \leq m} (Q_{s_i}(i) + c_{m-1}, T_i(s_i)) \right\}$$

over count

$$\leq l + \sum_{m=1}^n \sum_{s=1}^{m-1} \sum_{s_i=1}^{m-1} \left\{ Q(a^*) + c_{m-1}, T_{a^*}(s) \leq \right. \\ \left. Q_{s_i}(i) + c_{m-1}, T_i(s_i) \right\}$$

$$Q(a^*) \leq q_a(a^*) - [c_{m-1}, T_{a^*}(s)] > \{ \dots \}$$

↳ lower value s_i ↳ confidence value

$$Q_{s_i}(i) \geq q_a(i) + c_{m-1}, T_i(s_i)$$

$$q_a(a^*) < q_a(i) + 2c_{m-1}, T_i(s_i)$$

$$Q(a^*) \leq l + \sum_{m=k+1}^n \left\{ \min_{0 \leq s \leq m} (Q_s(a^*) + c_{m-1}) \right. \\ \left. \leq \max_{1 \leq s_i \leq m} (Q_{s_i}(i) + c_{m-1}, T_i(s_i)) \right\}$$

$$q_a(i)$$

! :: also same answer

$$\textcircled{1} \quad P(Q_s, c_a^*) \leq q_s(c_a^*) - C_{m,T} \epsilon^{(s)} \leq n^{-4} \quad (\text{m})$$

$$\textcircled{2} \quad P(Q_{s_i}(i) \geq q_s(i) + C_{m,T} \epsilon^{(s)}) \leq n^{-4}$$

\textcircled{3} For $\delta = \sqrt{8 \epsilon n m} / \Delta_i^2 \geq 1$, \textcircled{3} is false!

$$q_s(c_a^*) - q_s(i) - 2C_{m,T} \epsilon^{(s)}$$

$$= q_s(c_a^*) - q_s(i) - \Delta_i = 0$$

$$E[\tau_i(c_n)] \leq \left\lceil \frac{8 \ln n}{\Delta_i^2} \right\rceil + \sum_{m=1}^{k-1} \sum_{s=1}^{\ell-1} \sum_{\tau_s(c_n) / \Delta_i^2} 2^{m-4}$$

$$\leq \frac{8 \ln n}{\Delta_i^2} + 1 + \frac{\pi^2}{3}$$

Chapter 9 PAC Exploration and naive Algo Proof (PAC Bounds)

Chernoff-Hoeffding Bound: Let x_1, \dots, x_n be d.r.v with common range $[0, 1]$, s.t

$$E[x_t | x_1, \dots, x_{t-1}] = \mu. \quad \text{Let } s_n = \frac{x_1 + \dots + x_n}{n}$$

then for $\epsilon > 0$,

$$P[s_n \geq \mu + \epsilon] \leq e^{-2\epsilon^2 n}$$

$$P[s_n \leq \mu - \epsilon] \leq e^{-2\epsilon^2 n}$$

→ →

Asymptotic
Regret

Mazurov Inequality

$$P_r(X \geq a) \leq \frac{E(X)}{a}$$

Union Bound

$$P(\cup_i A_i) \leq \sum_i P(A_i)$$

media action eliminated

PAC Bounds for MAD

Input $\epsilon > 0, \delta > 0$

output "An arm"

for each $a \in A$ do

Sample ' a ' l -times,

$$l = \frac{2}{\epsilon^2} \cdot \ln\left(\frac{2k}{\delta}\right)$$

Let $Q(a)$ be the average reward

of a

end

Output $a = \operatorname{argmax}_{a \in A} Q(a)$

cycle action elimination
of reinforcement
learning

Theorem 1: Naive (ϵ, δ) bound $(\epsilon, \delta) - \text{PAC}$

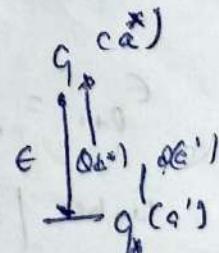
algo, with arm sample complexity

$$\mathcal{O}(C \sqrt{\epsilon^2} \cdot \log(1/\delta))$$

Proof Let a' be an arm s.t. $q_{\pi}^{*}(a') < q_{\pi}^{*}(a^*) - \epsilon$

$$P(Q(a') > Q(a^*)) \leq P(Q(a') > q_{\pi}^{*}(a') + \epsilon/2)$$

$$\text{or } Q(a^*) < q_{\pi}^{*}(a^*) - \epsilon/2$$



$$\leq P(\varphi(c^*)) > q_{\frac{1}{2}}(c^*) + \epsilon_{1/2} + P(\varphi_{c^*}) < q_{\frac{1}{2}}(c^*) - \epsilon_{1/2}$$

$$\leq 2e^{-2} \frac{\epsilon^2}{4} l$$

$$\text{Subst. } P(\varphi(c^*)) > q_{\frac{1}{2}}(c^*)) \leq \frac{\delta}{k}$$

Summing over all a^* ,

$$\text{Prob. of failure} \leq (k-1) \frac{\delta}{k} \leq \delta$$

Lecture 10 Median elimination Algo

Input $\epsilon > 0, \delta > 0$; process in rounds.

Output An arm

$$\text{Set } S_1 = A, \epsilon_1 = \epsilon/4, \delta_1 = \delta/2, l = 1$$

repeat

Sample every arm $a \in S_l$ for

$$\frac{1}{(\epsilon/4)^2} \cdot \log\left(\frac{3}{\delta_1}\right) \text{ let } \varphi_a(c) \text{ denote } \left(\frac{\epsilon L}{2}\right)^2 \cdot \log\left(\frac{3}{\delta_1}\right)$$

in value

$$\leftarrow \frac{1}{\epsilon L^2} \cdot \log\left(\frac{3}{\delta_1}\right)$$

find among the median of $\varphi_a(c)$, denoted $\log m_l$

$$S_{l+1} = S_l \setminus \{a : \varphi_a(c) < \log m_l\}$$

$$\epsilon_{l+1} = \frac{3}{4} \epsilon_1, \quad \delta_{l+1} = \frac{\delta}{4}, \quad l = l+1$$

$$\liminf |S_l| = 1$$

Theorem The MEA (ϵ, δ) is a (C, δ) PAC

Algo: Sample complexity $O\left(\frac{K}{\epsilon^2} \cdot \log\left(\frac{1}{\delta}\right)\right)$

Lemma For every phase l in MET, we have

$$P\left[\max_{j \in S_l} q_j^*(j) \leq \max_{i \in S_{l+1}} q_i^*(i) + \epsilon_l\right] \geq 1 - \delta_l$$

Proof :- wlog, consider $l=1$ $E_1 = [\phi(a_1^*) < q_1(a_1^*) - \epsilon_{1/2}]$

$$\epsilon_1 = [\phi(a_1^*) < q_1(a_1^*) - \epsilon_{1/2}]$$

$$P[E_1] \leq \delta_{1/3} \quad \text{True best arm in 1^{th} level}$$

$$P[Q_1(j) \geq Q_1(a_1^*) \mid Q_1(a_1^*) > q_1(a_1^*) - \epsilon_{1/2}] \leq \delta_{1/3}$$

1st arm i , not ϕ optimal

Let # bad₁ be the no. of bad arms from which

(1) holds

$$E[\#\text{bad}_1 \mid Q_1(a_1^*) \geq q_1(a_1^*) - \epsilon_{1/2}] < 1\delta_1$$

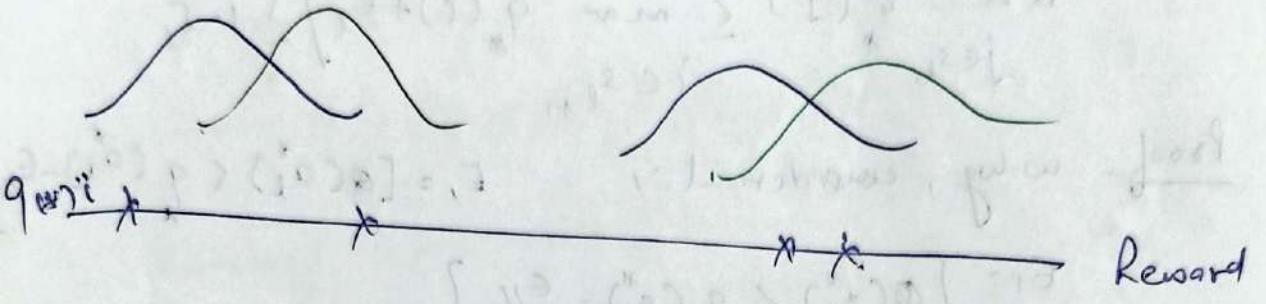
$$P[\#\text{bad}_1 \geq \frac{1\delta_1}{2} \mid Q_1(a_1^*) \geq q_1(a_1^*) - \epsilon_{1/2}]$$

$$\leq \frac{1\delta_1(1\delta_1)}{1\delta_1/2}$$

$$= \frac{2\delta_1}{3}$$

Chapter 11 Thompson Sampling

OR
Posterior Sampling



Chapter 12 Policy search

$$\bar{\pi}_t(a) \stackrel{D}{=} \Pr(a_t = a)$$

Binary bandit (only two outcomes 0 or one)

$$R_t \in \{0, 1\}$$

$$\text{if } R_t = 1$$

$$\bar{\pi}_{t+1}(a_t) = \bar{\pi}_t(a_t) + \alpha [1 - \bar{\pi}_t(a_t)]$$



$$\pi_{t+1}(a') = \pi_t(a')(\alpha)$$

$a' \neq a_t$

If $R_t = 0$?

$$\pi_{t+1}(a_t) = \pi_t(a_t)(1-\beta)$$

$$\pi_{t+1}(a') = \pi_t(a') + \beta[1 - \pi_t(a')]$$

if $\alpha = \beta$, $L_{R=0}$

$\alpha \gg \beta$, $L_{R \neq 0}$

$\beta = 0$, $L_{R=1} \rightarrow$ inaction

(go playing
against paper)

Policy gradient (Parametrization ^{lower} policy gradient)

- Depends on some set of parameters, θ
- Preferably weight of $A^{n \times n}$
- Modified parameters $\theta' = \theta + \alpha \nabla \pi(a; \theta)$

Chapter 13

Reinforce

$\pi(\cdot; \theta)$

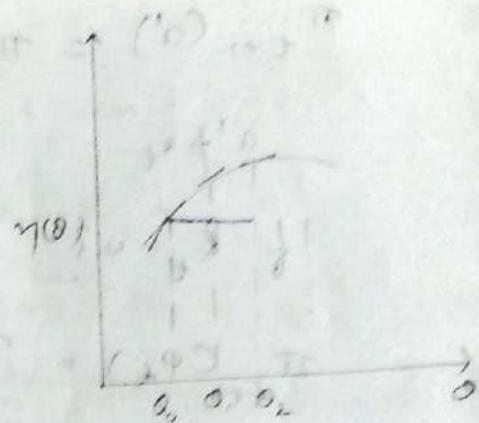
$$\gamma(\theta) = E[R_t]$$



$$\sum_a q(a), \pi(a; \theta)$$

$$\theta \leftarrow \theta + \alpha \nabla \gamma(\theta)$$

$$\begin{aligned}\nabla \cdot \eta^{(0)} &= \sum_a q_{\pi}(a) \nabla_0 \pi(a, \theta), \\ &= \sum_a \left(q(a) \frac{\nabla_0 \pi(a, \theta)}{\pi(a, \theta)} \right) \pi(a, \theta) \\ E_{\pi(\cdot; \theta)} \left[\frac{q(a) \nabla_0 \pi(a, \theta)}{\pi(a, \theta)} \right] \\ &\approx \frac{1}{n} \sum_{t=1}^n R_t \frac{\nabla_0 \pi(a_t, \theta_t)}{\pi(a_t, \theta_t)} \quad [\text{gradient}]\end{aligned}$$



Incremental Version: $\theta_{n+1} = \theta_n + \Delta \theta_n$

$$\Delta \theta_n = \alpha_n R_n \frac{\nabla \pi(a_n, \theta_n)}{\pi(a_n, \theta_n)}$$

Rewards
of good
otherwise bad
reinforcement
in direction opposite
to gradient

$$\Delta \theta_n = \alpha_n R_n \frac{\partial \ln \pi(a_n, \theta_n)}{\partial \theta}$$

$$\Delta \theta_n = \alpha_n (R_n - b_n) \frac{\partial \ln \pi(a_n, \theta_n)}{\partial \theta}$$

Characteristic
eligibility-
which have
the highest
gradient

Reinforce

Two actions are reinforced

$$\pi(a, \theta) = \begin{cases} \theta_a & a=1 \\ 1-\theta & a=0 \end{cases}$$

$$\frac{\partial \ln \pi}{\partial \theta} = \frac{a-\theta}{\theta(1-\theta)}, \quad \alpha = \rho(0(1-\theta), b=0)$$

$$\Delta \theta_n = \rho(a=1) \cdot R_n$$

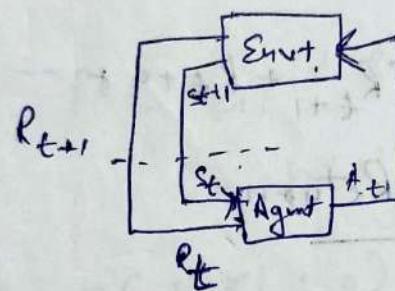
$$\pi(a, s, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} \cdot e^{-\frac{(a - \mu)^2}{2\sigma^2}}$$

Chapter 14 Contextual Bandits

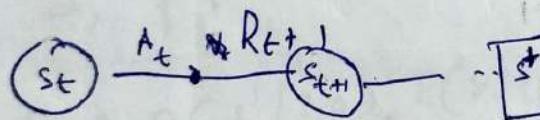
- action can have parameters
- Parameterized Bandit

$$Q(s) = \alpha_1 \phi_1 + \alpha_2 \phi_2 + \alpha_3 \phi_3 + \dots$$

Chapter 15: Full RL introduction



r_t Reward signal
from environment



$t=1$ decision 1
 $t=2$ decision 2
Chapter 3 in the book

Reward

- scalar quantities :- Just a number
- outside the direct control of agent
- bounded
- frequent

Policy $\pi_t(a|s) = \Pr(A_t = a | S_t = s)$

$$\pi_1 \neq \pi_2$$

$$\Pr(A_1 = a_1 | S_1 = s_1)$$

$$= \Pr(A_1 = a_1 | S_1 = s_1)$$

Chapter 16: Returns, Value functions and MDPs

Return value function and MDP

Return:

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

total return

$$G_t = R_{t+1} + \delta R_{t+2} + \delta^2 R_{t+3} + \dots$$

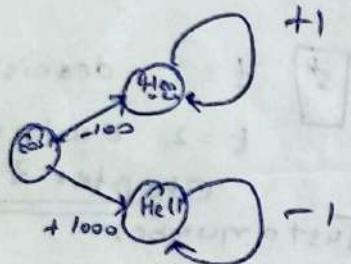
Discounted return $0 \leq \delta < 1$

near sighted
agents \rightarrow far sighted.

$$G_t = \lim_{N \rightarrow \infty} \frac{1}{N} [R_{t+1} + R_{t+2} + \dots + R_{t+N}]$$

Average reward Return:

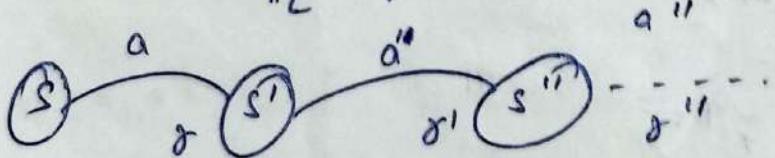
Level of return
to the correlation



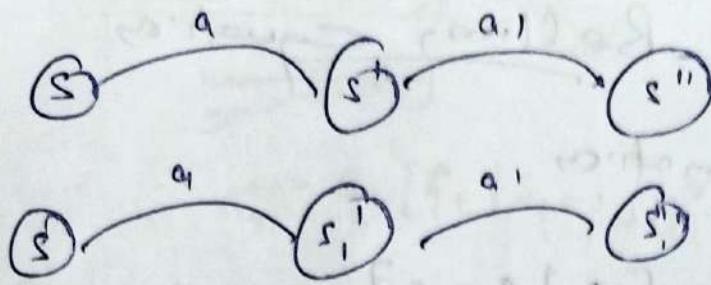
Value function:-

$v^\pi(s)$ is value function associated with policy π

$$v^\pi(s) = E_\pi \{ G_t | s_t = s \}$$



$$q^\pi(s, a) = E_\pi \{ G_t | s_t = s, a_t = a \}$$



How could you characterize the problem.

$$P_r(s_{t+1}, r_{t+1} | s_t, a_t, \dots)$$

$$= P_r(s_{t+1}, r_{t+1} | s_t, a_t)$$

$$\cancel{P_r(s_{t+1})} \quad P_r(s_{t+1} = s', r_{t+1} = r | s_t = s, a_t = a) \\ = p(s', r | s, a)$$

$$S, A, p(s', r | s, a), r$$

$$S, A, p(s' | s, a), E[r | s, a, s']$$

Markov Decision Process

Chapter 17: MDP Modeling

$$s, a, p(s', r | s, a), E[r | s, a, s']$$

$$r_{can} = (1 + b)p \cdot (2)$$

$$r_{out \text{ of battery}} = -\beta$$

$$r_{recharge} = 0$$

- | | |
|---------------|-------------------------|
| <u>State</u> | - charge: low, high |
| <u>Action</u> | - No. of cans |
| | - Achieve seek can & |
| | - Wait for cans β |
| | - Recharge |

Y > R

Chapter 18 Bellman equations

Value functions

$$v^\pi(s) = E_\pi [q_t | s_t = s]$$

$$q^\pi(s) = E_\pi [q_t | s_t = s, a_t = q]$$

$$\begin{aligned} v^\pi(s) &= E_\pi [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s] \\ &= E_\pi [R_{t+1} + \gamma v^\pi(s') | s_t = s] \end{aligned}$$

$$v^\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) \{ E[\gamma | s, a, s'] + v^\pi(s') \}$$

Bellman eqn. for v^π, q^π

$$q^\pi(s, a) = \sum_{s'} p(s'|s, a) \{ E[\gamma | s, a, s'] + \sum_a \pi(a|s') q^\pi(s', a) \}$$

$$v^\pi(s) = \sum_a \pi(a|s) \cdot q^\pi(s|a)$$

Chapter 19 Bellman Optimality eqn

Optimal Value Function:

w.r.t π best of all π

$$\max_{\pi} [V_{\pi}]$$

$$\max_{\pi} E_{\pi} \{ R_t | S_t = s \} + \gamma$$



$$\pi^* = \arg \max_{\pi} V^{\pi}(s) \quad \forall s$$

optimal policy

$$V^{\pi^*} \triangleq N^* = \max_{\pi} V^{\pi}(s) \quad \forall s$$

↑
optimal value function for

$$q^*(s, a) = \max_{\pi} \sum_{s'} p(s'|s, a) \cdot [E[r|s, a, s'] + \gamma V^{\pi}(s')]$$

$$= \sum_{s'} p(s'|s, a) \cdot [E[r|s, a, s'] + \gamma \max_{\pi} V^{\pi}(s')]$$

$$V^*(s) = \max_a q^*(s, a)$$

$$q^*(s) = \max_a \sum_{s'} p(s'|s, a) [E[V|s, a, s'] + \gamma V^*(s')]$$

Bellman optimality eqn for V^*

$$\pi^*(s) = \arg \max_a \sum_{s'} p(s'|s, a) [E[r|s, a, s'] + \gamma V^*(s')]$$

$$\pi^*(s) = \arg \max_a q^*(s, a)$$

Chapter 20 Cauchy Sequence and Green's Equation

MDPs $s, \pi, p(s'|s, a), E[r|s, a, s'] \&$

Bellman eqn:

$$\therefore \hat{v}^\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) [E[r|s, a, s'] + \gamma v^\pi(s')]$$

$$v^\pi(s) = \max \sum_{s'} p(s'|s, a) [E[r|s, a, s'] + \gamma v^\pi(s')] \quad \text{finite MDPs}$$

$\pi(a|s), \pi(s) \rightarrow \text{Def. policy}$

v^π is a vector with $|S|$ components

$$\|v\| = \sup_{s \in S} |v(s)|$$

$$\begin{aligned} \|x\| &= 0 \text{ iff } x = 0 \\ \|Px\| &= \alpha \|x\| \\ \|x+y\| &\leq \|x\| + \|y\| \end{aligned}$$

complete normed vector space

Cauchy sequence

x_1, x_2, x_3, \dots
for every $\epsilon > 0$ $\exists n \in \mathbb{N}$

$\exists N \in \mathbb{Z}^+ \text{ st } \forall m, n > N$

$$\|x_m - x_n\| < \epsilon$$

If every cauchy sequence is convergent
then the vector space is complete

$$\hat{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) E[r|s, a, s']$$

$$\left[= \sum_{s'} |u(s'|s, \pi(s)) \cdot E[r|s, \pi(s), s'] \right]$$

$$p_{\pi}(j|s) = \sum_a \pi(a|s) \cdot p(j|s,a)$$

$$\{ = p(j|s, \pi, \omega)$$

π is a $|s|$ dimensional vector p_{π} is a $|s| \times |s|$ matrix dim stochastic matrix

$$0 \leq v < 1$$

$$\pi + \gamma P_{\pi} \cdot v = v^{\pi}$$

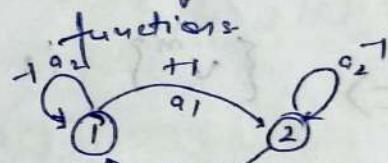
$$\Rightarrow v^{\pi} = (I - \gamma P_{\pi})^{-1} \pi$$

Google it!
Relationship of eigenvalues to invertibility of matrix

Chapter 21 completely non vector space is Banach space

Banach fixed pt. theorem

V = space of all value functions



$$L_{\pi} : V \rightarrow V$$

$$L_{\pi} v = \pi + \gamma P_{\pi} v$$

$$L_{\pi} v^{\pi} = v^{\pi} \Rightarrow v^{\pi} \text{ is a fixed point of } L_{\pi}$$

Banach fixed pt. theorem

(T, U, T_U, T_V)

contraction mapping

Suppose T is a Banach space and $T: U \rightarrow U$ is a contraction mapping then

$T: U \rightarrow U$ is a contraction mapping then $T: U \rightarrow U$ is a contraction mapping then

\exists an unique v^* in U . sc. $Tv^* = v^*$ and for arbitrary v^0 in U , the sequence $\{v^n\}$ defined by $v^{n+1} = T v^n$, converges to v^*

T is a contraction, if \therefore apply in fixed point space

Proof

$$\|Tu - Tv\| \leq \lambda \|u - v\|, 0 \leq \lambda < 1, \forall u, v \in U$$

$$\|v^{n+m} - v^n\| \leq \|v^{m+n} - v^{n+m-1}\| + \|v^{n+m-1} - v^n\|$$

$$\|v^{n+m} - v^n\| \leq \sum_{k=0}^{m-1} \|v^{n+k+1} - v^{n+k}\|$$

$$= \sum_{k=0}^{m-1} \|T^{n+k} v' - T^{n+k} v^0\|$$

$$\leq \sum_{k=0}^{m-1} \lambda^{n+k} \|v' - v^0\|$$

$$= \frac{\lambda^n (1 - \lambda^m)}{1 - \lambda} \|v' - v^0\|$$

$\Rightarrow \{v^n\}$ is cauchy

Chapter 22 Convergence Proof

$$0 \leq \|Tv^* - v^*\| \leq \|Tv^* - v^n\| + \|v^n - v^*\| \quad \text{triangle inequality}$$

$$= \|T v^* - T v^{n-1}\| + \|v^n - v^*\|$$

$$\leq \lambda \|v^* - v^{n-1}\| + \|v^n - v^*\|$$

$$\|v^* - v^n\| \rightarrow 0 \quad \text{as } n \rightarrow \infty$$

$$0 \leq \|Tv^* - v^*\| \leq 0$$

$$\Rightarrow \overline{T v^*} = v^*$$

Let u^*, v^* be 2 fixed points

$$\|Tu^* - Tv^*\| \leq \lambda \|u^* - v^*\|$$

$$\Rightarrow u^* = v^*$$

Let $u \neq v$ be in V

$$L_\pi u(s) = \delta_\pi(s) + \sum_{j \in s} \gamma^j p_\pi(j|s) \cdot u(j)$$

$$L_\pi v(s) = \delta_\pi(s) + \sum_{j \in s} \gamma^j p_\pi(j|s) v(j)$$

$$\text{Let } L_\pi u(s) = L_\pi(v(s))$$

$$0 \leq L_\pi v(s) - L_\pi u(s) \leq \delta_\pi(s) + \gamma \sum_j p_\pi(j|s) v(j) - \delta_\pi(s) - \gamma \sum_j p_\pi(j|s) u(j)$$

$$\leq \gamma \|v - u\| + \sum_j p_\pi(j|s) [v(j) - u(j)]$$

$$= \gamma \|v - u\|$$

Other way (con)

$$|L_\pi v(s) - L_\pi u(s)| \leq \gamma \|v - u\|$$

$\Rightarrow L_\pi$ is a contraction

Chapter 23! L_p convergence

L_π is a contraction.

V space of value function

$v \in V$

v/s bounded function

$$v^* = \max_{\pi} \{ \delta_\pi + \gamma p^\pi v^* \} \triangleq L v^*$$

$$v_{cs}^* = \max_a \{ E[\delta] s, a] + \gamma \sum_j p(c_j | s, a) v^*(c_j) \}$$

$$L v \triangleq \max_{\pi} \{ \delta_\pi + \gamma p^\pi v \}$$

Claim L is a contraction

$$\text{Let } v^* = \max_a \{ E[\delta] s, a] + \gamma \sum_j p(c_j | s, a) v(c_j) \}$$

$$0 \leq L v(s) - L u(s) \leq E[\delta] s, a^*] + \gamma \sum_j p(c_j | s, a^*) - E[\delta] s, a^*] + \gamma \sum_j p(c_j | s, a^*)$$

$$= \gamma \sum_j p(c_j | s, a^*) [v(c_j) - u(c_j)]$$

$$\leq \gamma \|v - u\| \sum_j p(c_j | s, a^*)$$

$$\gamma \|v - u\| \leq L v(s) - L u(s) \leq \gamma \|v - u\|$$

$$|L v(s) - L u(s)| \leq \gamma \|v - u\| \quad \forall s$$

$\Rightarrow L$ is contraction



Value iteration

- value iteration
1. Select $v^0 \in V$, pick an $\epsilon > 0$ Set $n=0$
 2. for each $s \in S$, compute $v^{n+1}(s)$ by

$$v^{n+1}(s) = \max_a \left\{ E[r|s,a] + \gamma \sum_{s'} p(s'|s,a) v^n(s') \right\}$$

$$\{ v^{n+1} = Lv^n \}$$
 3. If $\|v^{n+1} - v^n\| < \epsilon (1-\gamma)/2\gamma$
 go to step 4, otherwise increment n by 1, go to step 2
 4. for each $s \in S$

$$\pi(s) = \arg \max_a \left\{ E[r|s,a] + \gamma \sum_{s'} p(s'|s,a) v^{n+1}(s') \right\}$$

stop

Chapter 24 Value iteration proof

Theorem: Let $v^0 \in V$ $\epsilon > 0$, $\{v^n\}$ be derived from

$$v^{n+1} = Lv^n$$

Then $\|v^n\|$ converges in norm to v^*

b. \exists a finite n at which (3) is met for
 $\|v^n - v^*\| \leq \epsilon_1$

c. π defined by (4) is ϵ -optimal

d. $\|v^{n+1} - v^*\| \leq \epsilon_2$ when (3) holds

○ suppose (3) is met for some n, α, π
 satisfies (C1) then

$$\|v^n - v^*\| \leq \|v^n - v^{n+1}\| + \|v^{n+1} - v^*\|$$

$$L_\pi v^{n+1} = E[\alpha(s, \pi(s))] + \delta \sum_{s'} p(s'|s, \pi(s)) v^{n+1}(s)$$

$$\rightarrow \|v^n - v^{n+1}\| = \|L_\pi v^n - v^{n+1}\|$$

$$\leq \|L_\pi v^n - L v^{n+1}\| + \|L v^{n+1} - v^{n+1}\|$$

$$= \|L_\pi v^n \cdot L_\pi v^{n+1}\| + \|L v^{n+1} - L v^n\|$$

$$\leq \gamma \|v^n - v^{n+1}\| + \gamma \|v^{n+1} - v^n\|$$

$$\|v^n - v^{n+1}\| \leq \frac{\gamma}{1-\gamma} \|v^{n+1} - v^n\|$$

$$\|v^n - v^*\| \leq \sum_{k=0}^{\infty} \|v^{n+k+1} - v^{n+k}\|$$

$$= \sum_{k=0}^{\infty} \|L_\pi v^n - L_\pi v^{n+k}\|$$

$$\leq \sum_{k=0}^{\infty} \gamma^{k+1} \|v^{n+1} - v^n\|$$

$$= \frac{\gamma}{1-\gamma} \|v^{n+1} - v^n\|$$

Linear convergence with rate γ

Rates of convergence for different algorithms

Chapter 25: Policy Iteration

Select π_0 , set $n=0$

Policy evaluation. $v^{\pi_n} = (I - \gamma P_{\pi_n})^{-1} r_{\pi_n}$

3. Policy Improvement Choose π_{n+1}

$$\pi_{n+1} \in \arg \max_{\pi} \{r_{\pi} + \gamma P_{\pi} v^{\pi_n}\}$$

4. stop if $\pi_{n+1} = \pi_n$, $\pi_n = \pi^*$

Let π_{n+1} satisfy (3)

$$r_{\pi_{n+1}} + \gamma P_{\pi_{n+1}} v^{\pi_n} \geq r_{\pi_n} + \gamma P_{\pi_n} v^{\pi_n} \cdot \pi_n$$

$$r_{\pi_{n+1}} \geq (I - \gamma P_{\pi_{n+1}})^{-1} v^{\pi_n}$$

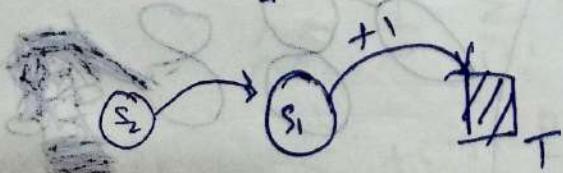
$$(I - \gamma P_{\pi_{n+1}})^{-1}$$

$$v^{\pi_{n+1}} > v^{\pi_n}$$

Chapter 26

Dynamic Programming

$$v^{n+1}(s) = \max_a \left\{ E[r | s, a] + \gamma \sum_{s'} p(s' | s, a) \cdot v^n(s') \right\}$$



Asynchronous DP

RTDP (Real time dynamic programming)

(P_{t+1}^{π}) GPI

- way the error

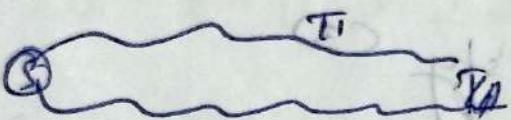
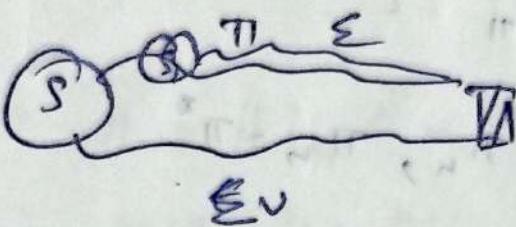
PE : Policy evaluation

PJ : Policy improvement

Chapter 27 Monte carlo

$$v^\pi(s) = E_{\pi}[G_t | s_t = s]$$

when we have no
access to MDP



Monte carlo methods for policy evaluation

sample models

Exploring Model

First visit MC

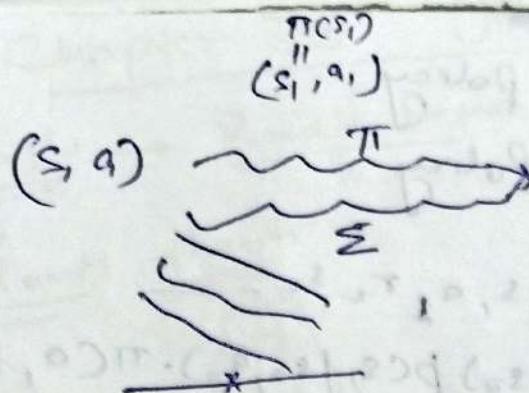
Every visit MC

Chapter 28

Centered Monte carlo

GPI (Generalized policy iteration)

9.



used q function for controlled

Epsilon soft policy

~~off policy~~

$$E[f(x)] = \sum_x p(x) f(x) \approx \frac{1}{n} \sum_{i=1}^n f(x_i)$$

$x \sim p$

$$= \sum_x \frac{p(x) \cdot q(x)}{q(x)} \cdot f(x)$$

$$= \sum_x \frac{p(x)}{q(x)} \cdot f(x) \cdot q(x)$$

$$= E_{x \sim q} \left[\frac{p(x)}{q(x)} \cdot f(x) \right] \approx \frac{1}{n} \sum_i f(x_i) \frac{p(x_i)}{q(x_i)}$$

Importance factor / weight

Chapter 29 off policy MC

$$E_{x \sim p} [f(x)] = E_{x \sim q} \left[f(x) \frac{p(x)}{q(x)} \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^n f(x_i) \cdot \frac{p(x_i)}{q(x_i)}$$

$$\text{weighted is } = \frac{\sum_{i=1}^n f(x_i) \frac{p(x_i)}{q(x_i)}}{\sum_{i=1}^n \frac{p(x_i)}{q(x_i)}}$$

π :- Estimation Policy

μ :- Behavior Policy

$a_i \leftarrow s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots$

$$p(s_{t+1}) = \pi(a_0 | s_0) p(s_1 | s_0, a_0) \cdot \pi(a_1 | s_1)$$

$$p(s_2 | s_1, a_1), \dots$$

$$q(s_{t+1}) = p(s_0) \cdot \mu(a_0 | s_0) p(s_1 | s_0, a_0) \mu(a_1 | s_1)$$

$$p(s_2 | s_1, a_1)$$

$$\frac{p(s_{t+1})}{q(s_{t+1})} = \frac{\pi(a_0 | s_0) \pi(a_1 | s_1)}{\mu(a_0 | s_0) \mu(a_1 | s_1)}$$

$$= \frac{\prod_{j=0}^T \pi(a_j | s_j)}{\prod_{j=0}^T \mu(a_j | s_j)} = \frac{\prod_{j=0}^T \frac{\pi(a_j | s_j)}{\mu(a_j | s_j)}}{\prod_{j=0}^T \mu(a_j | s_j)}$$

$$\sqrt{\pi(s)} = \sqrt{\sum_{i=1}^n \frac{\prod_{j=0}^{(i)} \frac{\pi(a_j | s_j)}{\mu(a_j | s_j)}}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}}$$

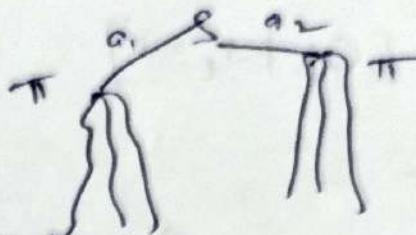
$$\sqrt{\sum_{i=1}^n \frac{\prod_{j=0}^{(i)} \frac{\pi(a_j | s_j)}{\mu(a_j | s_j)}}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}} = \sqrt{\sum_{i=1}^n \frac{1}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}}$$

$$\sqrt{\sum_{i=1}^n \frac{1}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}} = \sqrt{\sum_{i=1}^n \frac{1}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}} = \sqrt{\sum_{i=1}^n \frac{1}{\prod_{j=0}^{(i)} \mu(a_j | s_j)}}$$

Chapter 30 UCT

Google → Roulette based Monte carlo planning

Rollouts theorem



search (state, depth)

if terminal (state) the return 0 // ^{Terminal} reward

if leaf (state, depth) the return Eval(state)

Action - selected Action (state, depth)

(next state ~~reward~~) = Simulate (state, action)

$g = \text{reward} + \gamma \text{ search}(\text{next state}, \text{depth}+1)$

update values (state, action, g, depth)

return g / \sqrt{n}

$$CP \left[\frac{2 \ln t}{n} \right]$$

chaw lectures

$$CP \left[\frac{\ln t}{m} \right]$$

Chapter 31 TD(0)

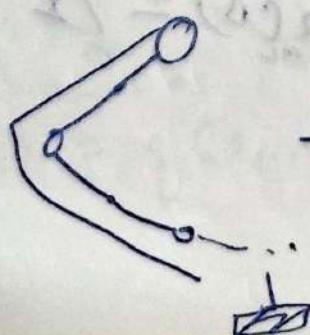
Temporal Difference method.



full backup

Bootstrap

MC



- sample backup
- no bootstrap

TD



- Sample Backup,
- Bootstrap.

$$\stackrel{MC}{=} E[G_t | s_t = s]$$

 G_t

$$\hat{G}_{\text{new}}(s) = \hat{f}_{\text{old}}(s) + \alpha \left[G_t^{\text{old}} + \underbrace{\varepsilon_t}_{\substack{\text{current} \\ \text{target}}} - \underbrace{\hat{f}_{t+1}(s_{t+1})}_{\text{current}} \right]$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots$$

$$= R_{t+1} + \gamma G_{t+1}$$

$$\approx R_{t+1} + \gamma \mathcal{V}(s_{t+1})$$

$$\hat{f}_{\text{new}}(s_t) = \hat{f}_{\text{old}}(s_t) + \alpha \left[f_{t+1}^{\text{old}} + \gamma \hat{V}_{\text{old}}(s_{t+1}) - \hat{f}_{t+1}(s_t) \right]$$

TD(0) algorithm

$$\hat{V}_{\text{old}}(s_t) + \alpha \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+2} + \gamma^3 \hat{V}_{\text{old}}(s_{t+3}) \right]$$

$$A, 1, B, 0, (\cancel{B, 0}, \cancel{B, 0})$$

$$A, 0, 1, B, 0$$

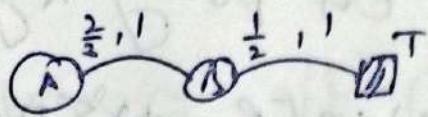
$$B, 0$$

$$B, 1$$

$$A, 1, B, 1$$

$$N_{mc}(A) = 1/2$$

$$N_{mc}(B) = 1/2$$



$$\frac{v(B)}{TD} = \frac{1}{2}$$

$$V_{TD}(A) = \frac{2}{3} + \gamma v(B) \quad \gamma = 1$$

$$= \frac{2}{3} + \frac{1}{2} = \frac{7}{6}$$

Chapter 32 TD(0) control

- Evaluation
- Sample backup & bootstrap

Control

$$q_{\text{new}}(s, a) = q_{\text{old}}(s, a) + \alpha \left[R_{t+1} + \gamma q_{\text{old}}(s_{t+1}, a_{t+1}) - q_{\text{old}}(s, a) \right]$$

 Target current
 { } { }
 Error

$$TD_{\text{error}} = \delta$$

Start with $q(s, a)$

Pick a state, so $t=0$

Pick action a_t in s_t , acc. to

$q(s_t, \cdot)$, ϵ -greedy

→ Apply a_t and sample s_{t+1}, R_{t+1}

- Pick a_{t+1} , acc. to $q(s_{t+1})$ (Greedy)

$$q(s_t, a_t) = q(s_t, a_t) + \alpha [R_{t+1} + \gamma q(s_{t+1}, a_{t+1}) - q(s_t, a_t)]$$

$$t \leftarrow t+1$$

SARSA - On Policy Algo.

$$* q(s_t, a_t) = q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_{a'} \pi(a_{t+1}, a') - q(s_{t+1}, a_t)]$$

Chapter 33

Q-Learning :-

$$q^*(s_t, a_t) =$$

$$E[R_{t+1} + \gamma q^*(s_{t+1})]$$

$$q^*(s_t, a_t) E[R_{t+1} + \gamma q^*(s_{t+1})]$$

$$= E[R_{t+1} + \gamma \max_{a'} q^*(s_{t+1}, a')]$$

$$q(s_t, a_t) = q(s_t, a_t) + \alpha [R_{t+1} + \max_{a'} q(s_{t+1}, a') - q(s_t, a_t)]$$

Start with $q(s, a)$

Pick a state s_0 , $t=0$

Pick action a_t in s_t according to $\pi(s_t, a_t)$; greedily

Apply a_t & sample $s_{t+1}, R_{t+1}, a_{t+1}$

$$q(s_t, a_t) = q(s_t, a_t) + \alpha [R_{t+1} + \max_{a'} q(s_{t+1}, a') - q(s_t, a_t)]$$

Q-Learning

ϵ greedy $\epsilon = 0.1$

π -learning off-policy

$$R_{t+1} + \gamma R_{t+2} + \gamma^2 \dots$$

Chapter 34

After state

Chapter 35 Eligibility traces
to reduce sample complexity

One way trying to speed up convergence
How you can debug the code

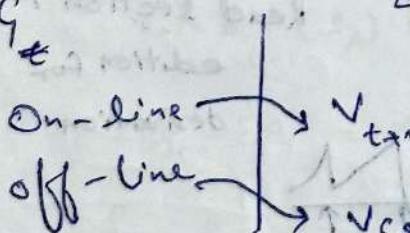
TDC(0) \rightarrow TDC(λ)

$$G_t^{mc} = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T = G_T$$

$$G_t^{TD} = R_{t+1} + \gamma V(s_{t+1}) = G_t^{(1)}$$
 Truncated return

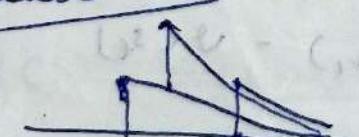
$$= R_{t+1} + \gamma R_{t+2} + \gamma^2 V(s_{t+2}) = G_t^{(2)}$$

$$G_t^{(n)} = \alpha [G_{t+n} - v_t(s)]$$



Backward view

no changes until an episode come to end



Forward view

$$G_T^\lambda = 0.5 G_t^{(1)} + 0.5 G_t^{(2)}$$

eligibility mechanism
immediate updates

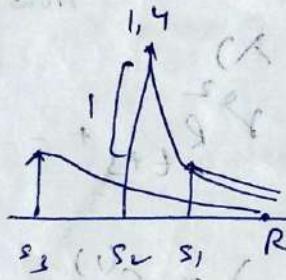
+ trace mechanism

$$G_t^{\lambda} = \gamma \lambda G_t^{(1)} + \text{bias}$$

$$G_t^{\lambda} = (1-\lambda) [G_t^{(1)} + \lambda G_t^{(2)} + \lambda^2 G_t^{(3)} + \dots]$$

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

Chapter 36 Backward view of eligibility traces



$$\epsilon_t(s) = \begin{cases} (\gamma \lambda) \epsilon_{t-1}(s) & \text{if } s_t \neq s \\ (\gamma \lambda) \epsilon_{t-1}(s) + 1 & \text{if } s_t = s \end{cases}$$

Accumulating trace

Read section 7.4 first edition for derivation

$\rightarrow 1$ if $s_t = s$

Replacing traces.



$$\delta_t = R_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\Delta v_t(s) = \alpha \cdot \delta_t \cdot \epsilon_t(s)$$

$$TD(u) = TD(u + \delta_t)$$

Dutch trace

Read set 7.4
1st edition

$$\epsilon_t(s) = \begin{cases} (\gamma^\lambda) \epsilon_{t-1}(s) & \text{if } s_t \neq s \\ \beta(\gamma^\lambda) \epsilon_{t-1}(s) + 1 & \text{if } s_t = s \end{cases}$$

Chapter 37 :- Eligibility trace control

TDC(λ) - Estimation

control

SARSA(λ)

$$\epsilon_t(s, a) = \gamma^\lambda \epsilon_{t-1}(s, a) + 1 \quad \text{if } s_t = s, A_t = a$$

$$= \gamma^\lambda \epsilon_{t-1}(s, a)$$

1 if $s_t = s, A_t = a$

0 if $s_t = s, A_t \neq a$

$$\gamma^\lambda \epsilon_{t-1}(s, a) \quad s_t \neq s$$

$$\Delta Q_t(s, a) = \alpha \cdot s_t \cdot \epsilon_t(s, a)$$

$$s_t = R_{t+1} + \gamma^\lambda Q_t(s_{t+1}, A_{t+1}) - Q_t(s_t, A_t)$$

$Q(s)$ Watkins & Lamba

$$s_t = R_{t+1} + \max_{a'} Q_t(s_{t+1}, a') - Q_t(s_t, A_t)$$

$$s \xrightarrow{a} s' \xrightarrow{\max a'} s'' \xrightarrow{\max a''}$$

By executing greed policy
executing

$$E_t(s,a) = \gamma \lambda E_{t-1}(s_{t-1}, a_{t-1}) + \text{if } s_t = s \quad A_t = a$$

$$s.t. = \arg \max_{a'} Q_t(s_t, a')$$

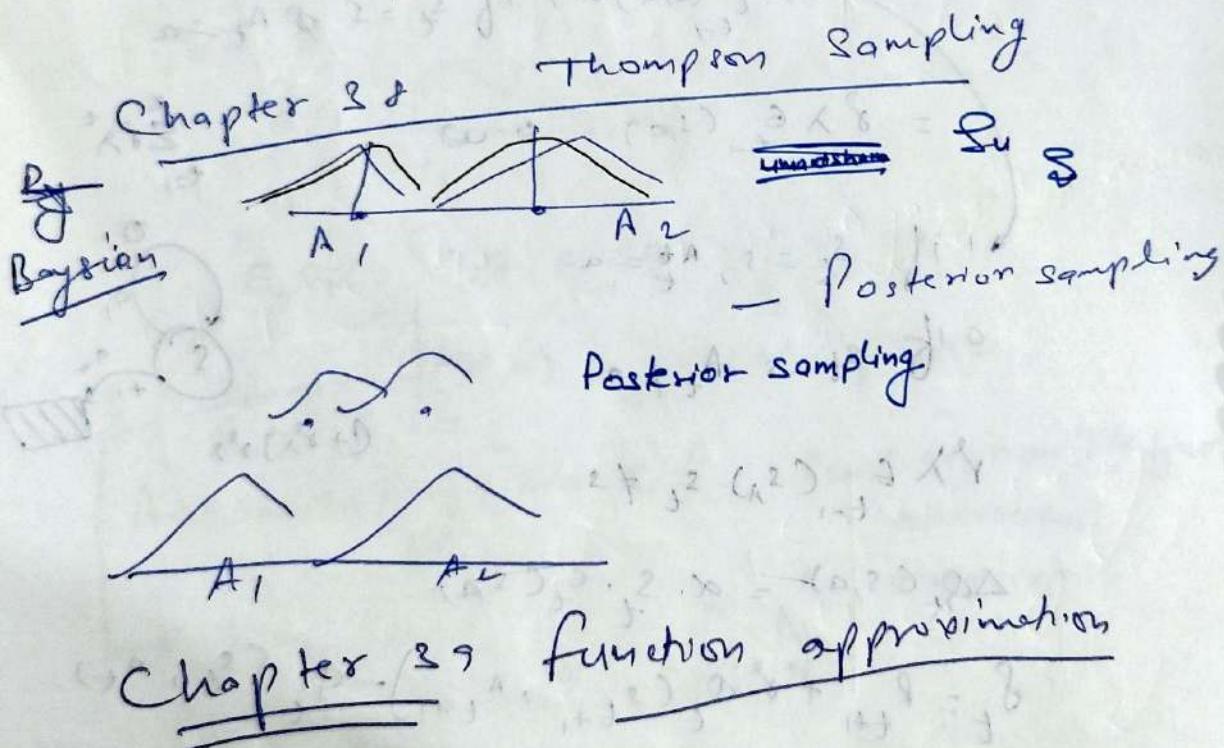
$$= \gamma \lambda E_{t-1}(s, a) \quad \text{if } A_t \neq a \quad A_t = \arg \max_{a'} Q_t(s_t, a')$$

$$= 0 \quad \text{if } A_t \neq \arg \max_{a'} Q_t(s_t, a')$$

Off Policy TD(λ) | $P_t(a) = \frac{\pi(s_t, a)}{\pi(s_t, a')}$

$$E_t(s) = \beta_t(\gamma \lambda E_{t-1}(s) + 1)$$

$$= P_t(\gamma \lambda E_{t-1}(s)) \quad \text{if } s_t = s$$



Lookup table

— memory
continuous state

— continuous state

— Generalize

— parameterized representation

— function Approximation

How output depend on t

$$v^\pi \quad \pi \hat{J}(s, \theta)$$

$$\sum_p^{\pi(s)} [\hat{v}(s_t, \theta) - v^\pi(s_t)]^2$$

$$\frac{1}{N} \sum_{t=1}^N [\hat{v}(s_t, \theta) - v^\pi(s_t)]^2$$

$$= \frac{1}{N} \sum_{t=1}^N [v^\pi(s_t) - \hat{v}(s_t, \theta)]^2$$

$$g_t$$

$$P_{t+1} + \gamma \hat{J}(s_{t+1}, \theta)$$

$$\text{Reg}: \langle n, f(n) : \rightarrow \hat{f}(n) \rangle$$

Rc

$$\begin{cases} s_t, A_t, R_{t+1}, s_{t+1}, A_{t+1}, R_{t+2}, \dots \\ \hat{J}(s_t, \theta) \end{cases}$$

$$\langle s_t, R_{t+1}, \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \rangle$$

$$\langle s_{t+1}, R_{t+2} + \gamma R_{t+3} + \dots \rangle$$

$$\langle s_t, R_{t+1} + \gamma \hat{J}(s_{t+1}, \theta) \rangle$$

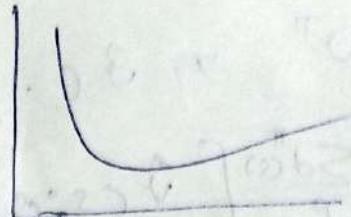
$$\langle s_{t+1}, R_{t+1} + \gamma \hat{J}(s_{t+1}, \theta) \rangle$$

$$\theta_{t+1} \leftarrow \theta_t + \alpha \nabla (R_{t+1} + \gamma \hat{J}(s_{t+1}, \theta_t) - J(s_{t+1}, \theta_t))^2$$

Chapter 40

Linear parameterisation

$$s \rightarrow (\phi_1(s), \phi_2(s), \phi_3(s), \phi_4(s))$$



$$\Phi(s)$$

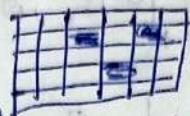
$$\bar{\Phi}$$

$$\hat{\mathbf{v}}(s, A) = \theta^T \mathbf{f}$$

for linear

$$\theta_{t+1} \leftarrow \theta_t$$

Linearly independent



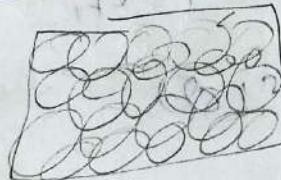
values of features are linearly independent

sq

Gangto

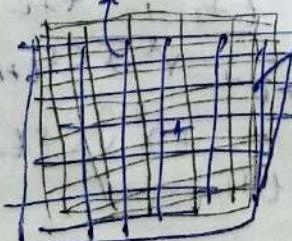
Chapter 41

state aggregation method



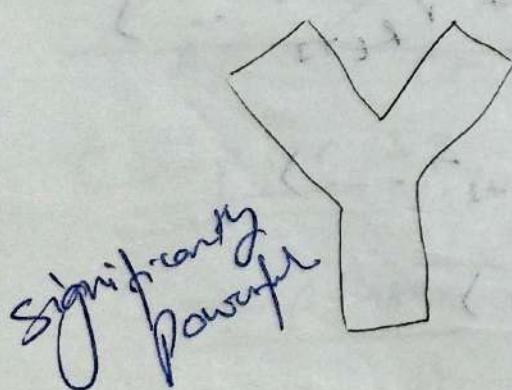
coarse coding

Tile coding



Tiling.

multiple
step grids



significantly powerful

CMAc - Cerebellar Model Action control

Chapter 42

Approximation and eligibility traces

trace

$$\delta_t = R_{t+1} + \gamma v(s_{t+1}) - v(s_t)$$

$$\vec{e}_t = \gamma \lambda \vec{e}_{t-1} + \nabla_{\theta_t} v(s_t)$$

$$\vec{\phi}_{t+1} = \vec{\phi}_t + \alpha \delta_t e_t$$

control w.r.t.

$$Q_t(s, a)$$

SARSA

$$\phi(s, a)$$

(Small no. of actions)

$$(\phi^T \phi_a) \delta_a = Q(s, a)$$

Large no. of actions.

- $\phi(s, a)$ - joint features

- after states

SARSA(λ)

$$\vec{e}_t = \gamma \lambda \vec{e}_t + \nabla_{\theta_t} \delta_t Q(s_t, a_t)$$

$$s - \phi(s) = \langle 0.01, 1.0, 1 \rangle$$

$$c_t^{(1)} = \gamma \lambda c_{t-1}^{(1)}$$

$$c_t^{(2)} = \gamma^2 \lambda c_{t-1}^{(2)}$$

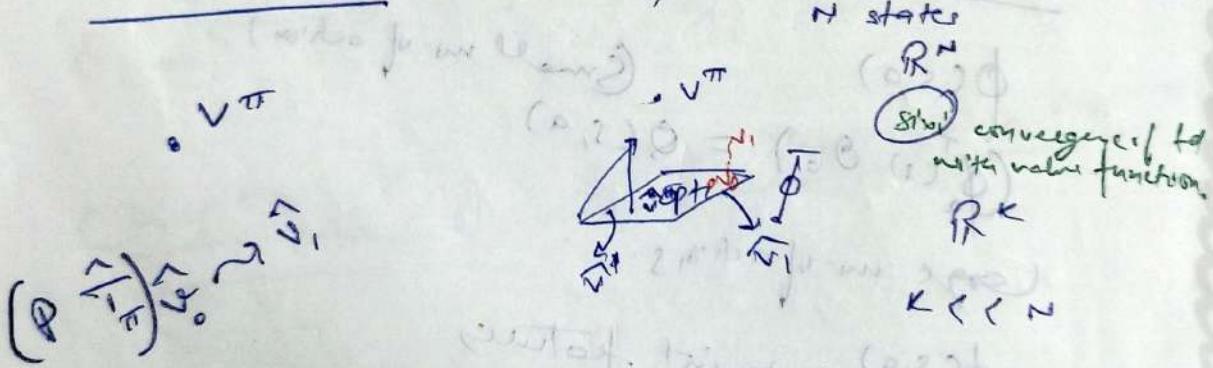
$$c_t^{(0)} = 1 \boxed{+ \gamma \lambda c_{t-1}^{(3)}}$$

$$c_t^{(4)} = 1 \boxed{+ \gamma \lambda c_{t-1}^{(4)}}$$

$$c_t^{(5)} = \gamma^3 \lambda c_{t-1}^{(5)}$$

$$c_t^{(6)} = 1 \boxed{+ \gamma^3 \lambda c_{t-1}^{(6)}}$$

Chapter 4.3 LSTD & LSTDQ



$$(P \hat{\pi}) \hat{V}_0 \rightarrow \hat{V}$$

$$\langle A_1, \dots, A_k, Y \rangle_{s=1}$$

$$\langle C\phi_r^{(s)}, \phi_u^{(s)} \cup \cdot \rangle$$

$$v(s) = \hat{P}_{ss} \hat{v} = \hat{v}$$

$$(P \hat{\pi}) \hat{V} = \hat{V}$$

$$\hat{V} = Cr + \gamma P \hat{\pi} \hat{V}$$

$$\hat{\pi}^{\pi} = \bar{\phi}(\bar{\phi}^T \bar{\phi})^{-1} \bar{\phi}^T (R^{\pi} + \gamma P^{\pi} \hat{\pi}^{\pi})$$

$$\hat{\Phi}^{\pi} = \bar{\phi}(\bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T (R^{\pi} + \gamma P^{\pi} \hat{\pi}^{\pi})$$

$$\hat{\phi}\theta^{\pi} = \bar{\phi}(\bar{\phi}^T \bar{\phi})^{-1} \bar{\phi}^T (R^{\pi} + \gamma P^{\pi} \hat{\phi}\theta^{\pi})$$

θ^{π} = weighted version

$$\underbrace{\bar{\phi}^T \bar{w} (\bar{\phi} - \gamma P^{\pi} \bar{\phi}) \theta^{\pi}}_{A} = \underbrace{\bar{\phi}^T \bar{w} R^{\pi}}_b$$

A

$$A: \theta^{\pi} = b$$

$$A = \bar{\phi}^T \bar{w} (\bar{\phi} - \gamma P^{\pi} \bar{\phi})$$

$$= \sum_s \phi(s), w(s) \cdot (\bar{\phi}(s) - \gamma \sum_{s'} P(s, \pi(s), s') \phi(s'))^T$$

$$= \sum_s w(s) \sum_{s'} P(s, \pi(s), s') \underbrace{[\phi(s) - \gamma(\phi(s) - \gamma \phi(s'))^T]}_{\text{del}}$$

$$\textcircled{*} D = \{(s_i, a_i, r_i, s'_i)\}_{i=1, 2, \dots, L}$$

$$\tilde{A} = \frac{1}{L} \sum_{i=1}^L [\phi(s_i) (\phi(s_i) - \gamma \phi(s_{i+1}))^T]$$

$$b = \bar{\phi}^T \bar{w} R^{\pi}$$

$$= \sum_s \phi(s) w(s) \cdot \sum_{s'} P(s, \pi(s), s') R(s, \pi(s), s')$$

$$= \sum_s w(s) \sum_{s'} P(s, \pi(s), s') [\phi(s), R(s, \pi(s), s')]^T$$

$\hat{b} = \sum$

$$\hat{b} = \frac{1}{L} \sum_{i=1}^L \phi(s_i) r_i$$

$$\hat{\pi} = \frac{1}{L} \sum_{i=1}^L \sum_{\hat{\pi}} [\phi(s_i, a_i) (\phi(s_i, a_i) - \gamma \phi(s'_i, \pi s'_i))]$$

$$\hat{b} = \frac{1}{L} \sum_{i=1}^L \phi(s_i, a_i) r_i$$

Chapter 44 LSTD and fitted Q

LSTD / LSTDQ

$$\pi_{LSTD}^{(t+1)} = \arg \max_a \hat{\phi}^T \pi^{(t)}(s, a)$$

$$\pi_{t+1} = \arg \max_a \hat{\phi}^T \hat{\theta}^{(t)}$$

$$q_j = \phi(s_1) \dots \phi(s_j) \left(\begin{array}{c} s_1, \dots, s_j \\ a_1, \dots, a_j - \pi_{t+1} \end{array} \right)$$

$$a_j q_1(s_1) \dots q_k(s_k)$$

$$\hat{\phi}(s_i, a_i) r_i + \gamma \max_a \hat{Q}_0(s'_i, a)$$

Chapter 45 DQN and fitted Q Iteration

fitted Q-iteration

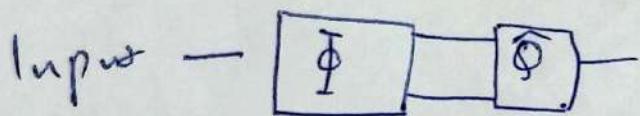
$\langle s_i, a_i, r_i, \pi_i \rangle$

$\langle \hat{Q}(s_i, a_i), \text{target}_i \rangle$

$$\text{target}_i = r_i + \gamma \max_a \hat{Q}(s'_i, a)$$

Neural fitted Q. FNFQ

DQN



(18)

"Online", off-policy algo

Replay

Memory $\langle s_i, a_i, \hat{s}_i, r_i \rangle$ structure

Transition replay