

Reinforcement Learning

Fundamental:

feedback may be delayed
time matters - ?

Markov Decision Process

bunch of ~~is~~ states and actions

State transitions

Probabilistic \rightarrow 2nd order

1st order MDP, only current state determining

= current and prev

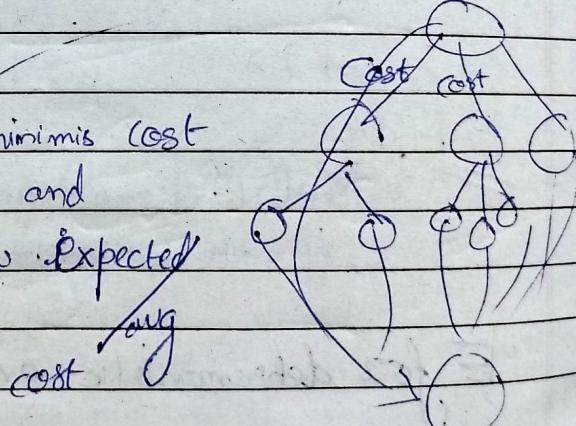


Stochastic shortest path. lena ha.
Uncertainty

Have to make

a policy lookup table ~~when~~ when this state do this

Now goal is to minimize cost
for a policy and
get to know expected



How? Bellman Eq,

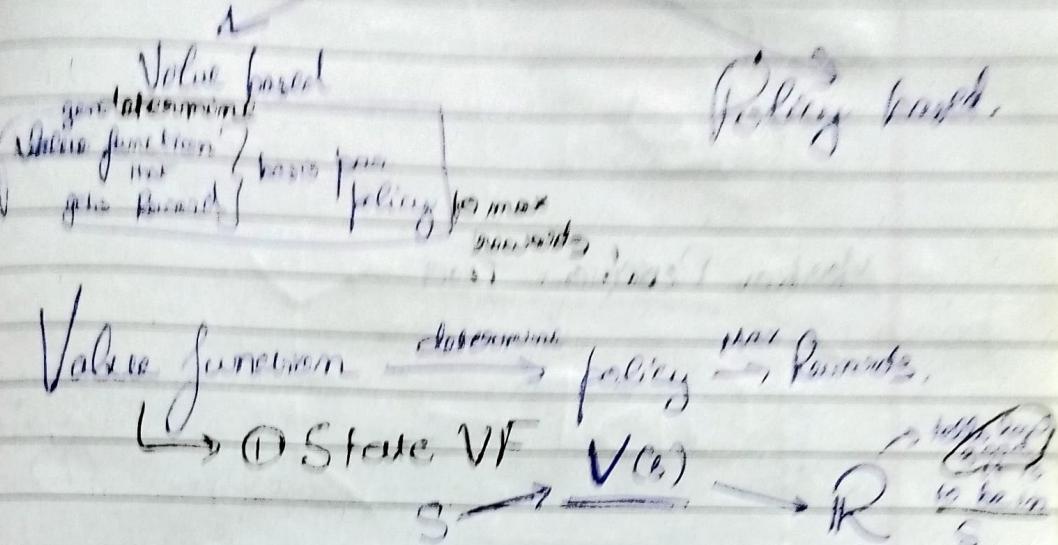
determine how good a choice is.

Why not use

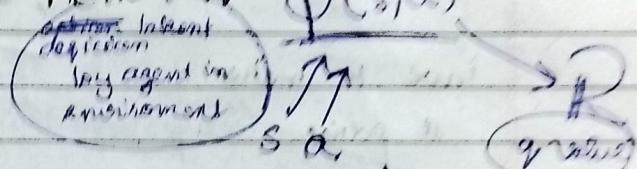
Dijkstra? when deterministic graph
No probability

(Page 14)
Date: / /

Reinforcement Learning



(2) State-Action VF $Q(s, a)$



BE: argmax

argmax to do
in s and to
take action a

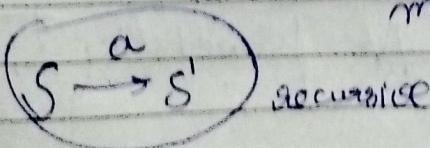
what long-term reward to expect
if i take best action in a state } Values of state

RE for deterministic env No probability

$$V(s) = \max_a (R(s, a) + \gamma V(s'))$$

Value of state

pick action that maximizes

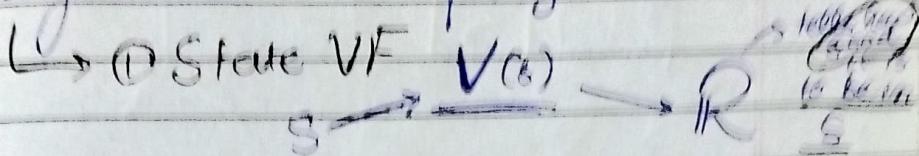


R.L. Algos

Value based
Value function
not basis for
policy to make
decisions

Policy based

Value function \rightarrow policy \rightarrow Rewards.



(2) State-Action VF $Q(s,a)$

action taken
by agent in
environment

s, a

R
(rewards)

Row good to tie
in s and be
taken action, a

BE. (Engelse)

what long-term reward is expected

if take best action for a state

E for deterministic env. No probability

$$V(s) = \max_a (R(a,a) + \gamma V(s'))$$

Value of state

pick
action that
maximizes

where you are
at

different
function

$s \xrightarrow{a} s'$
accuracy

~~sooner than the bank will pay~~
 Discount factor? $\gamma \rightarrow$ long term rewards
 $\gamma = \frac{1}{1 + r}$

Q: How to move which action at step S_t ?

~~R~~

Agent gets state s_t at instant t from environment
 on basis of π and policy
 applies a_t to environment.

Environment responds by ^{transition} s_{t+1} and also gives reward r_{t+1}
 action gets R_{t+1}

The Agent

• Control policy that maps S_t to A_t

$\pi(A_t | S_t)$

~~P~~

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

Given $\gamma = \infty$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots - \infty$$

\downarrow discount

~~Want to max~~ G_t

Now as R, A, π, S all is stochastic

We do not have defind G_t

Can find expectation value of G_t .

ret sequence of control actions that
maximise

$E(G_t)$

Dynamics of MDP

is the probability ^(conditional) of reaching state s' and getting r at t if already in s at $t-1$ and a done

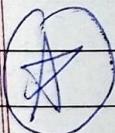
$$p(s'|s, a)$$

~~Prob S~~

↳ State transition prob.

$$p(s'|s, a)$$

reaching s' given in s at $t-1$ and action a



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

Policy $\pi(a|s)$ Prob of selecting a when at s .

$\pi(\cdot)$ is a function

gives us probability

$a|s$

or

a

goal is to get a good $\pi(\cdot)$

Temporal Difference Learning

using different fit algo. (i.e.)

Temporal Learning SDF

SDF

(α gradient
update)

Let's look into it.

$s_1 \rightarrow s_2$	-1	$+1$
$s_1 \rightarrow s_3$	-10	$+10$
$s_2 \rightarrow s_1$	-1	$+1$
$s_2 \rightarrow s_3$	-10	$+10$

Let's take sum of all states now.

0

sum of all states

s_1	$V(s_1)$
s_2	0
s_3	0
0	0

the agent take action a to
(s_2) then

$$V(s_1) = R(s_2) + V(s_2)$$

discuss.

Reward (asymmetric)

$$= -1 + 0$$

$$= -1 \text{ but } V(s_1) \text{ expected } = 0$$

Then

Temporal Difference Error = -1

using reward
first step

$$V(s_1)_{\text{actual}} -$$

$$V(s_1)_{\text{expected}}$$

Obs - Exp

Now more better, learned.

$$\textcircled{A} \quad V(s_1)_{\text{act}} = V(s_1) + \alpha \cdot (\text{TOE})$$

Repeat

Learning Rate

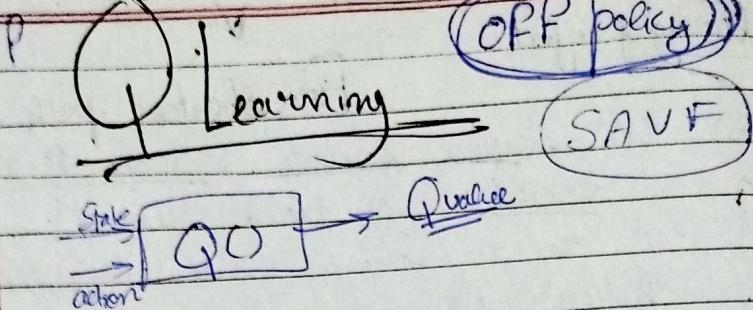
after 1 time step.

After learned, agent will take path that will result in
weighted $\sum Q$ value · get the target policy

Behaviour Policy is what is used to learn.

Page No.:

Date: / /



let q	-1	-1	-1	↓
	-1	-1	-1	
	-1	-10	+10	

over Q values in a table 2 find such Q as
option policy mix.

Q	at left	Right	down	up
s_1	-0.5	0.7	1.3	
s_2	0.3	0.75	2	-0.5
s_3				

let s_1 then randomly s_2

$a = \text{Right}$

max Q then m

$$Q(s_1, \text{Right})_{\text{obtained}} = R(s_2) + \gamma \max_a Q(s_2, a)$$

$$= -1 + 0.1 \times 2 = -0.8$$

i.e. down

check w/ for which
 a in state s_2

we have max
 $Q(s_2, a)$

$$\text{TDE} = Q(s_1, \text{Right})_{\text{obt}} - Q(s_1, \text{Right})_{\text{Exp}} \\ = -0.8 - 0.7 = -1.5$$

$$\text{Updated Value } Q(s_1, \text{Right})_{\text{new}} = Q(s_1, \text{R}) + \alpha \text{TDE}$$

Q. Learning project

Code)

Page No.:
Date: / /

We use negative rewards even in to normal traversal blocks because if +ve, then to maximise value, current agent would keep on walking on blocks.

So lab -ve for path, and huge -ve for forbidden area

Deep Q-Learning

NN + Q-Learning

CNN

ANN for input data with special structure

means img of a cat

→ Img, Vid, Tex-

Cat because all the pixels

are arranged in that specific way

CNN captures

special relationships

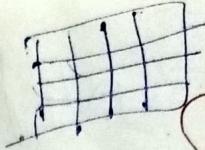
Working

① Generating set of feature maps

data for
CNN

② Simplify FM by pooling

③ Flatten the pFM



④ Connect pFM to NN



FM

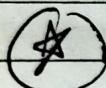
1st step \Rightarrow Apply filters / feature detector

determine if a particular part ^{kernel} of image contains a specific feature

Result is a
FM

Contains a specific feature

Collection of FM = Convolutional layer.



Eg Apply 3x3 filter to every 3x3 section of img.
Multiply each element of filter to corresponding cell in image. Add get degree of overlap
(divide by total cells in filter)

→ Why No ANN for img.?

Too many inp. weights (Computation)

Treats local pix same as far apart pix

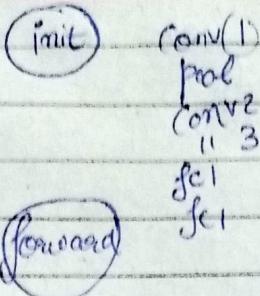
Sensitive to location of obj

→ Max Pooling = reduce size of each feature map
Take the max value from the section ^{count}
→ Window jumps by its size ^(2x2 for eg) rather 1px

→ Flattening = 2D pFM to 1D vector.
all pFMs in one vector

Feature layer

Input layer to NN

Class ConvNet (nn.Module)Exploration vs Exploitation tradeoff

ϵ -greedy strategy / implementation note

Tip to make it better: Start from 1 then decay to 0.01

DQL new best friend

Friendship with value iteration ended.

Image inp from eg game

feed as a state to NN get q values of corresponding to actions in that state.

Experience Replay & Replay Memory (Training)

Experience of agent at (t)

$$p_{\mathcal{E}_t} = (s_t, a_t, r_{t+1}, s_{t+1})$$

Replay memory gets containing these e_t from ∞ episodes

Usually stores last N e_t .

Why Cost?

Page No.:

Date: / /

We randomly sample from the ~~existing money set~~
to train ~~network~~

Called ~~Eq~~ Rely

Why? not normally using ~~sequential approach~~

To break ~~correlation~~ between ~~consec samples~~
Globally, samples highly correlated \Rightarrow inefficient learning

db weights initialization how?

Page No.:

Date: / /

- init display mem capacity N
- init weights
- for episode
 - start state
 - action
 - execute
 - reward, & next state
 - store experience

Training Starts

- Take any grand experience and feed forward
- Calc Loss (Net diff q and Bellman q)
optimal!

Again Loss

$$\Rightarrow q_{**}(s,a) = q(s,a)$$

$$\Rightarrow E[R_{t+1} + \gamma \max_a q_{*}(s',a)]$$

1st pass
feed s' to NN and get max q for a' .

2nd pass required for optimal value.

$$- E\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}\right]$$

Then grad descent

A new contender appears!!

(newest network)

target

Why? Issues w/ older approach

What?

THE 2nd PASS!

As for generating Optimal q w/ billion we run the NN again for Step 1, but it uses the same noisy weights that in the first place we were trying to improve

When update happens, target q value (2nd pass val)

No update in

q , value $\xrightarrow{\text{approach}}$ $\begin{matrix} \checkmark \\ \times \end{matrix}$ but $\begin{matrix} \checkmark \\ \times \end{matrix}$ $\xrightarrow{\text{also moves away}}$

Solution? ~~3rd pass (not on same NN)~~

Obtain $\text{target } q$ from Separate target NN

2nd pass now in Target NN!

(it's a clone) and its weight only

update after a custom timestep

New super parameters

N
ow, coming to 2048 !!

keep biggest in corner ! (lower) More reward if higher
bottom full

2048 is stochastic! so more interesting

Approach ① Monte Carlo Tree Search.

(A) heuristic driven search algo.

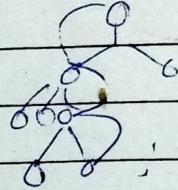
Search by
itself

Tree search + q-learning

→ explore all ($\epsilon = 1$)

→ get best (suboptimal)

→ keep looking, updates q-values



not much as not efficient

similar to Minimax

Look

→ MCTS only searches couple layer deep

Steps ① Tree traversal.

② Node expansion

Root node

③ Rollout

④ Backprop

Leaf nodes
2nd position