

RL - Week - 7

L1 Intro to Funcⁿ Approximaⁿ

- q funcⁿ is stored as a lookup table

s	$q(s, a_1)$	$q(s, a_2)$
s_1	—	—
s_2	—	—
s_3	—	—
- if they are large no. of states \rightarrow not efficient, data sparsity, contin. state / actⁿ spaces, we can't generalize it.
- \therefore , we can use a parameterized representaⁿ \rightarrow value funcⁿ, policies, models, etc. $\hookrightarrow f(x) = \alpha x + \beta^2 x + \gamma$
- a NN \rightarrow fully parameterized parameters object

L2 Linear Funcⁿ Approximaⁿ

Least sq. $Q(s_t, a_t) = \sum (s_t, a_t; w_t)$

$$w_{t+1} = w_t - \frac{1}{2} \alpha \nabla_{w_t} [q_{\pi}(s_t, a_t) - Q(s_t, a_t)]$$

but the problem, I don't the target itself.

maybe use the TD target.

$$w_{t+1} = w_t - \frac{1}{2} \alpha \nabla_{w_t} \left[r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - \underline{Q(s_t, a_t)} \right]^2$$

TD error, δ_t

$$f(s_{t+1}, a_t; w_t) f(s_t, a_t; w_t)$$

- GD \rightarrow 1st order iterative optimization algo. to find the local minima of a diff. funcⁿ. It is in the direction of steepest descent.

$$f(x_0 + h) \approx f(x_0) + h \cdot f'(x_0)$$

$$L(\theta + \alpha \Delta \theta) \approx L(\theta) + \alpha \Delta \theta^T \nabla_{\theta} L(\theta)$$

- Gradient ascent — in the direction of steepest ascent.

- so we use semi-gradient method
(we ignore the gradient of TD target).

$$Q(s_t, a_t) = \phi^T(s_t, a_t) \times w_t$$

$$\delta_t = r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - \underline{Q(s_t, a_t)}$$

$$\underline{w_{t+1} = w_t + \alpha \delta_t \phi(s_t, a_t)}$$

taking gradient w.r.t only periodic values

L3 State & Action Representations

$$\phi(s) = \langle x_1, y_3, 1000 \rangle \quad N = 1000$$

↳ kind of a
6-dimensional

$$S = 0100 \quad W = 0010$$

$$Q(s, w) = w_1 x_1 + w_2 y_2 + w_3 z + w_4$$

$$\phi(s, N) = \langle x_3, y_3, x_3^2, y_3^2, x_3 y_3, 0 \rangle$$

↳ this is still a linear approximation,
even though square terms are,
we do not have to solve for it.

$$\phi_s(s, a) = \text{is (block-on-table)} = \begin{cases} 1 & T \\ 0 & F \end{cases}$$

↳ these are features of state-action
value pairs.

L4 More on Representations

Lookup table

	a_1	a_2	
hot	$q(\text{hot}, \text{ad})$	$\neg \text{ad}$	using linear func ⁿ approximator
cold	$q(\text{cold}, \text{res})$	$\neg \text{res}$	

$$\phi(\text{hot}, \text{ad}) = \langle 1, 0, 0, 0 \rangle$$

$$\phi(\text{hot}, \neg \text{ad}) = \langle 0, 1, 0, 0 \rangle$$

$$\phi(\text{cold}, \neg \text{res}) = \langle 0, 0, 0, 1 \rangle$$

$$q(s, a) = w_1 \phi_1(s, a) + \\ w_2 \phi_2(s, a) + \\ w_3 \phi_3(s, a) + \\ w_4 \phi_4(s, a)$$

$$q(\text{hot}, \text{ad.}) = w_1 \cdot 1 + 0 + 0 + 0 = w_1$$

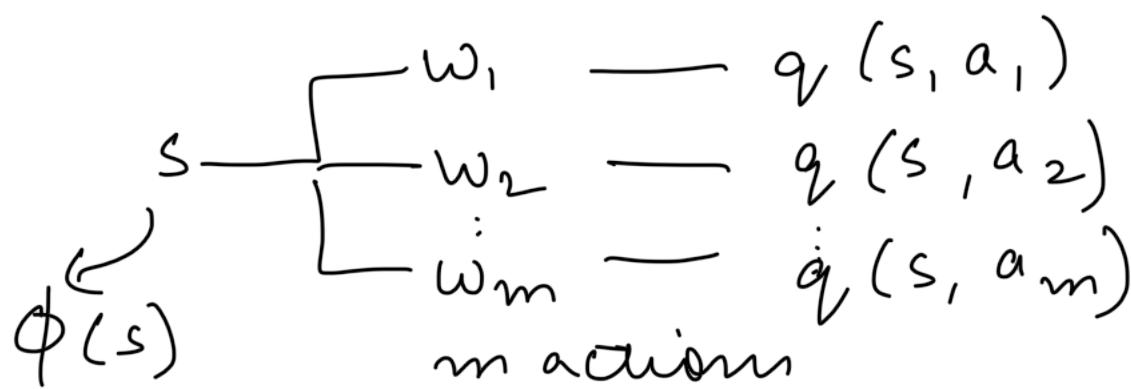
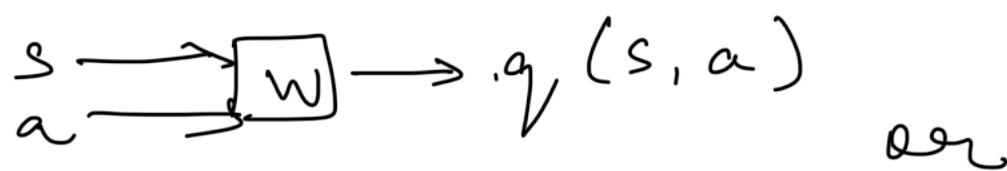
$$q(\text{hot}, \neg \text{ad}) = w_2$$

$$q(\text{cold}, \text{mes}) = w_3$$

$$q(\text{cold}, \neg \text{mes}) = w_4$$

another way,

$$\phi^T(s, a) w$$



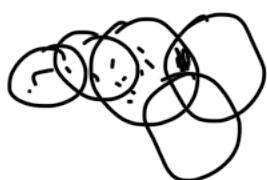
↪ these ϕ_s will be slowly learned in the process.

- we have to exploit the possibility the q values of near states would n't change much.

- Naive \rightarrow
 1. divide the grid 10×10 .
 2. abrupt Δ in Q value at the boundary.



- Coarse Coding \rightarrow avoid abrupt Δ s.
consider a circular world. Smooths the q values during the transition.

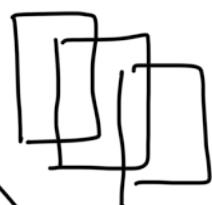


problem \rightarrow no uniform no. of "ON" bits representing states.

- Tile Coding \rightarrow like coarse but in grid world, more systematic.

No. of ON bits = no. of tiles used.

\hookrightarrow and each one of tile represents a feature.



\rightarrow leads to CMAC (cerebral model articulation controller)

powerful classifier, linear, no complex forward pass, in backward pass just update what is required.

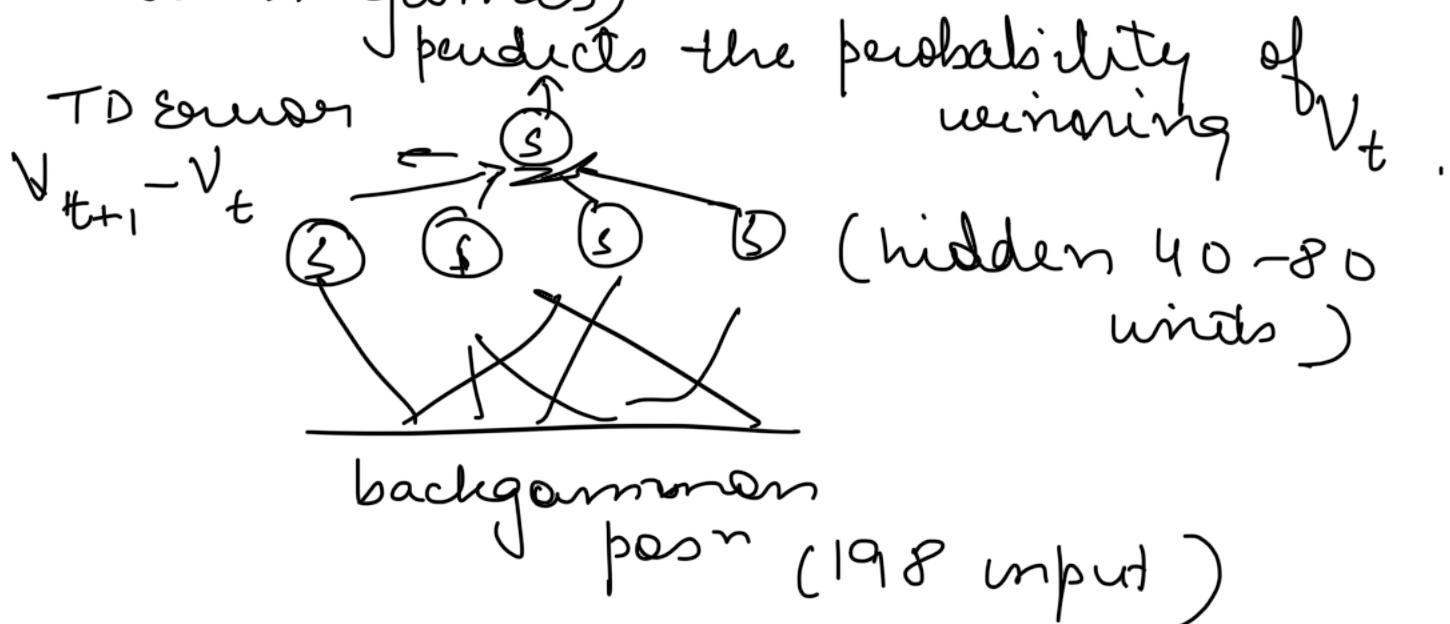
- In SGD \rightarrow the true gradient of loss funcⁿ is approximated by the gradient at a single eg.

- In usually, we do mini-batch GD, use small egs. This allows for efficient

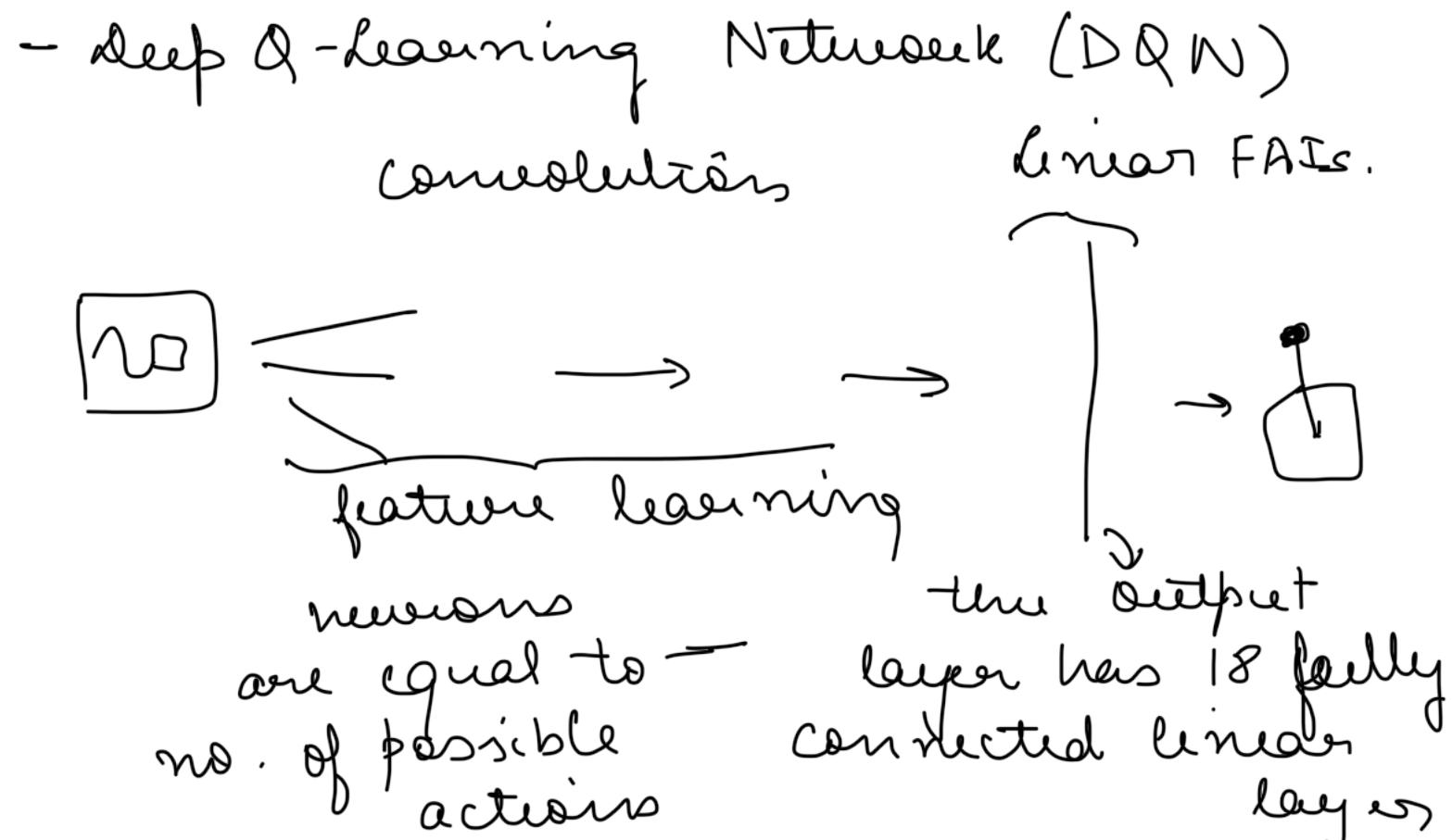
computaⁿ & smooth convergence (approx. of gradient, ..., still SGD).

L5 DQN (Deep Q-Networks)

- linear funcⁿs are restricted. They can only model linear funcⁿ/extrapolation.
- non linear approx. can model complex funcs, very powerful, compute & data^{(dis)ad_u} hungry. Features are learnt on the fly and not like tile coding.
- It generalizes well on unseen cases.
- Eg. TD-Gammon. Classical AI players can't play that good. Self play (1.4 millions games)



- what about Atari games? learnt from scratch, therefore, ϕ , were also learnt from scratch. ϕ was initialized randomly



- Finally we train using the TD error that also calculates the ϕ_s .

- Q-Network Learning

$$w_{t+1} = w_t - \frac{1}{2} \alpha \nabla_{w_t} [u_{t+1} + r \max_a q(s_{t+1}, a) - q(s_t, a_t)]^2$$

- target is non-stationary.
- correla" btw. samples.
- divergence issue \rightarrow current network is used to decide its own target.

- How to address this? \rightarrow replay memory, freeze target network.

- Replay memory \rightarrow fixes correla".
 ↳ build a data-set from the

agent's experience.

$$s_1, a_1, r_2, s_2$$

:

$$s_t, a_t, r_{t+1}, s_{t+1}$$

- frozen target network - non-stationary
↳ sample experience from dataset
slow down the process. w - frozen
for some time period.

$$\left[\{a_{t+1} + \gamma \max_a q(s_{t+1}, a; w^-)\} - q(s_t, a_t; w) \right]$$

- except for few games, DQNs do well than humans.

- Architecture of DQN

1. set of conv. layers which act as feature extractors
2. features are passed through FC layers.
3. output layer has $|A|$ no. of nodes.

uses new huber loss not sq. least sq. loss.

$$L_g(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta \\ \delta(|a| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

L6 Advanced Value-Based Methods

- Double DQN (DDQN)

1. add replay memory & target network.
2. make Q network a NN
3. add 2 more networks. Then we can use this target network to estimate the value while the online network is used for selecting the action..

$$y_t^{\text{DDQN}} = R_{t+1} + \gamma \underset{a}{\operatorname{argmax}} \underset{Q_t^-}{Q(S_{t+1}, a'; \theta_t)}$$

\swarrow target network params

- dueling DQN \rightarrow sometimes it's not necessary to estimate the value of each action choices in many states.

↳ Architecture -

1. Value network attends on both horizons of track as well as score.
2. the attention network concentrates only when there are inbound cars.

- so we estimate the Q value using both $A(s, a)$ & $V(s)$.
- Advantage \rightarrow calculates the benefit of taking an action.
- What is the maximum value of advantage funcⁿ?

$$Q(s, a) = V(s) + A(s, a)$$

↓ ↓ ↓
 to update value funcⁿ advantage
 all other funcⁿ funcⁿ
 things

(for most of the funcⁿ
 A will be 0, only
 V has to learn.)

- Aggregaⁿ layer \rightarrow a sep. layer is eng.
to combine $V(s)$ & $A(s, a)$
- Issue \rightarrow it is impossible to recover $V(s)$ & $A(s, a)$ from $Q(s, a)$
- Overcome \rightarrow agg. advantage is calculated!

$$Q(s, a; \theta, a, b) = V(s; \theta, b) + (A(s, a; \theta, a) - \frac{1}{A} \sum_{a'} A(s, a'; \theta, a))$$

Thus, we can now apply Q -learning techniques such as DDQN & PER with the dueling architecture.

- SARSA

$$Q(s_t, a_t) \quad \cdots \quad r Q(s_{t+1}, a_{t+1})$$

a_{t+1} introduces variance which slows convergence.

↓
Expected SARSA

we are approximating by just 1 sample.

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_{t+1} + \gamma \sum_a \pi(a | s_{t+1}) Q(s_{t+1}, a) - Q(s_t, a_t)]$$

it converges faster.

- Expected SARSA is ON Policy. The expectation is related to π . We are looking at the distribution of each action.