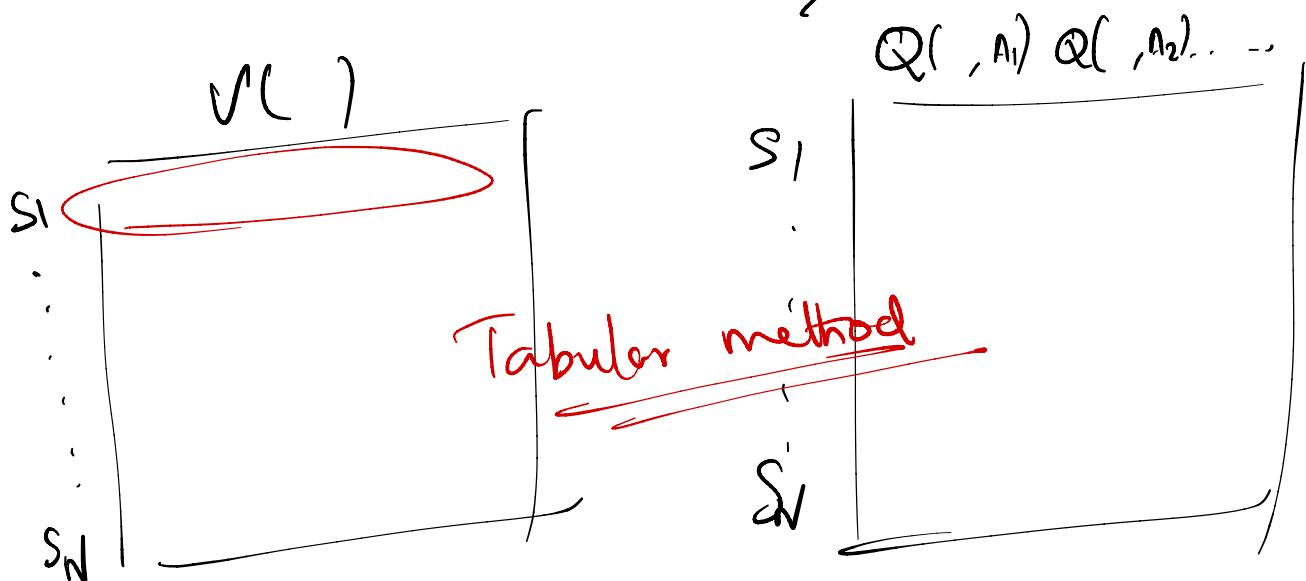
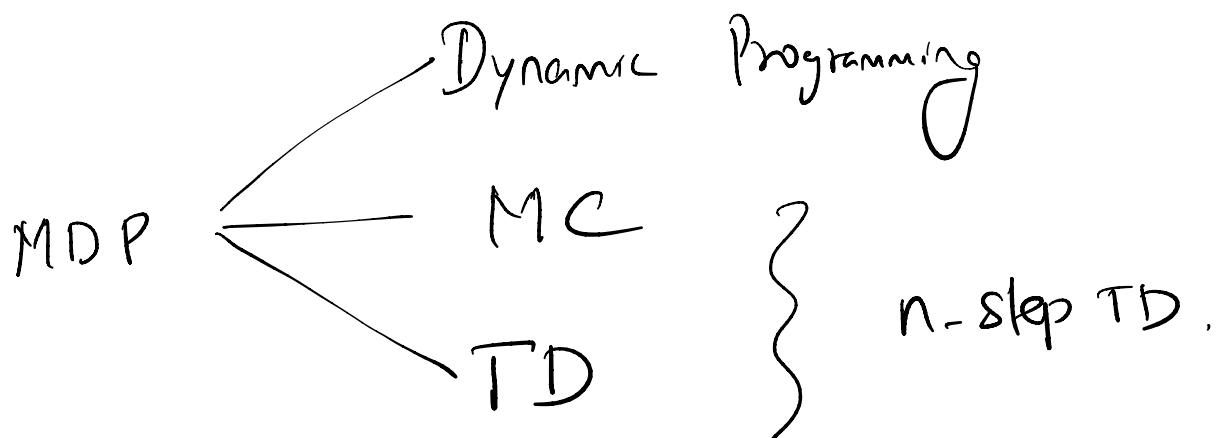


Lecture - 07

Value Function Approximation



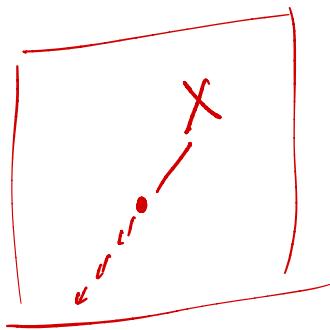
$$s_t \rightarrow x(s_t)$$

$$\hat{v}(s; \omega) = \omega^T x$$

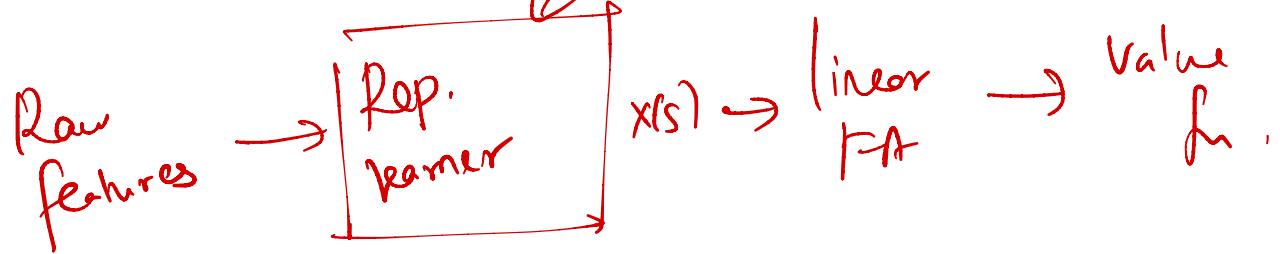
Non-linear function approximators:-

$$\hat{v}(s; \omega) = \omega^T x$$

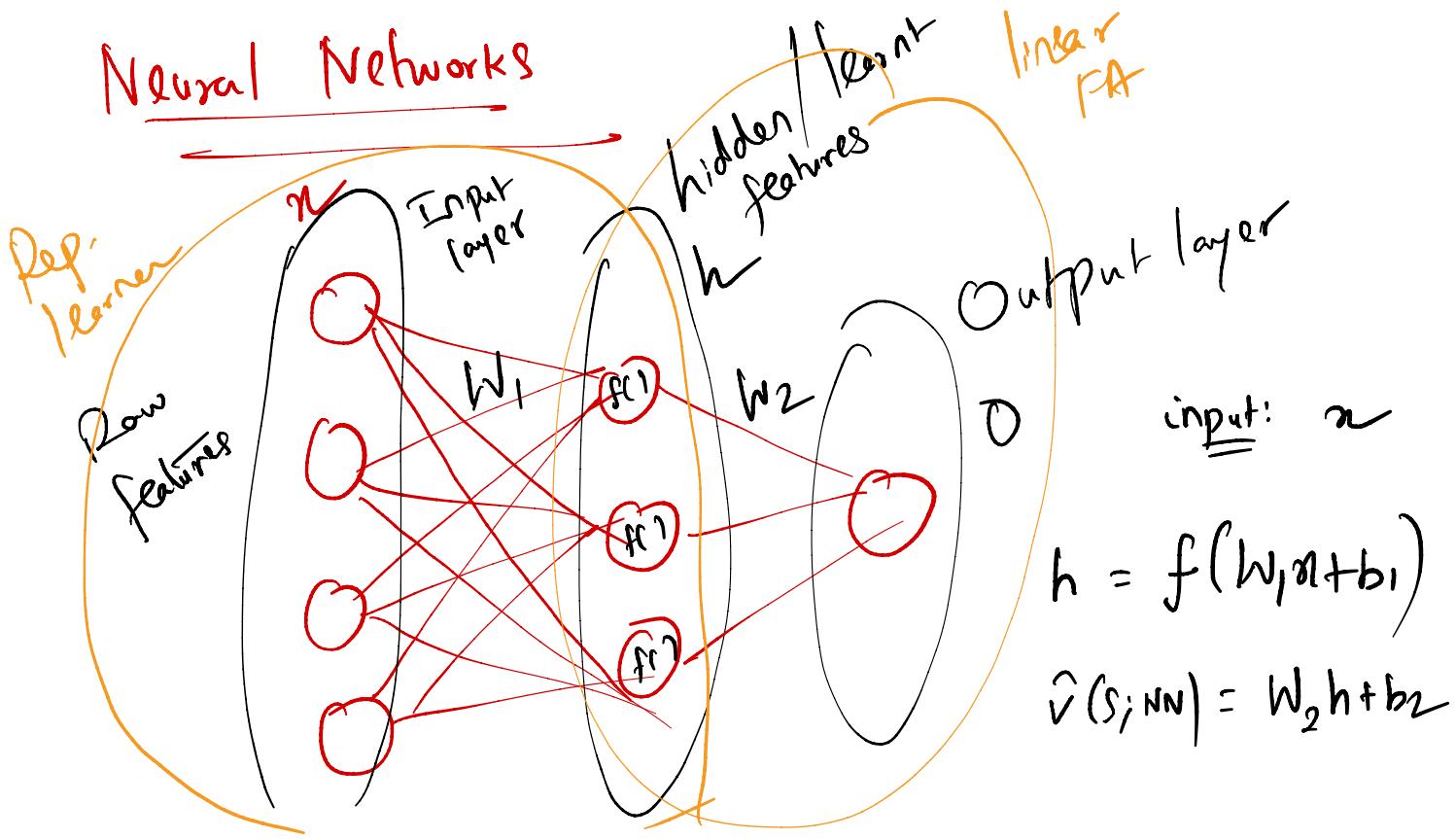
- ① Carefully handcrafted features.
- ② Interaction between the features is limited.



Representation Learning :- non-linear model



Neural Networks



$$h = f(w_1x + b_1)$$

$$\hat{v}(s; \theta) = w_2 h + b_2$$

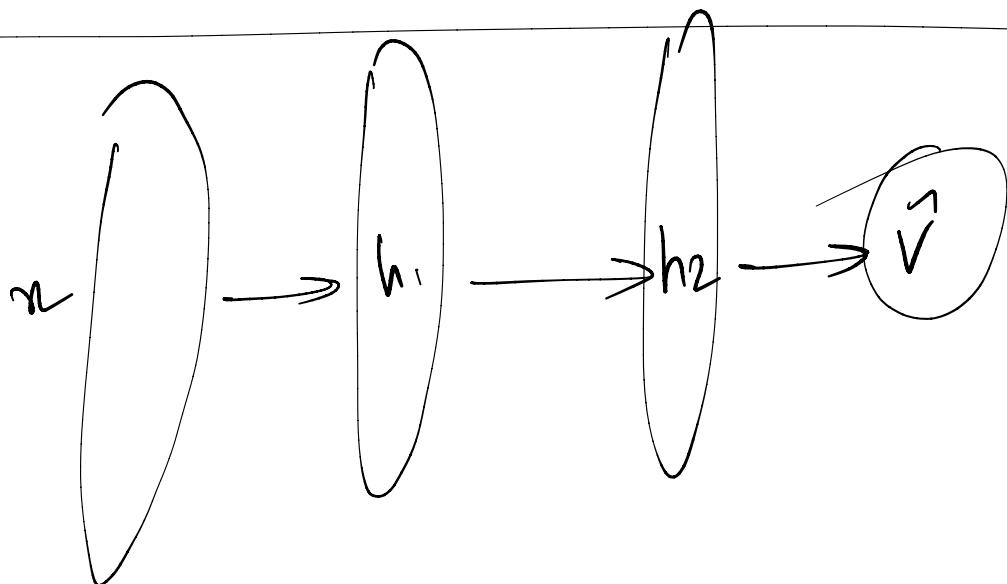
$$\theta = (w_1, b_1, w_2, b_2)$$

Non-linear activation fn:

1) Sigmoid $\sigma(a) = \frac{1}{1+e^{-a}}$

2) tanh $\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$

3) ReLU (a) $= \max(\theta, a)$
Rectified linear unit



Fully
Connected
Network:
Multi Layer
Perceptron
(MLP)

$$h_1 = \sigma(w_1n + b_1)$$

$$h_2 = \sigma(w_2h_1 + b_2)$$

$$o = w_3 h_2 + b_3$$

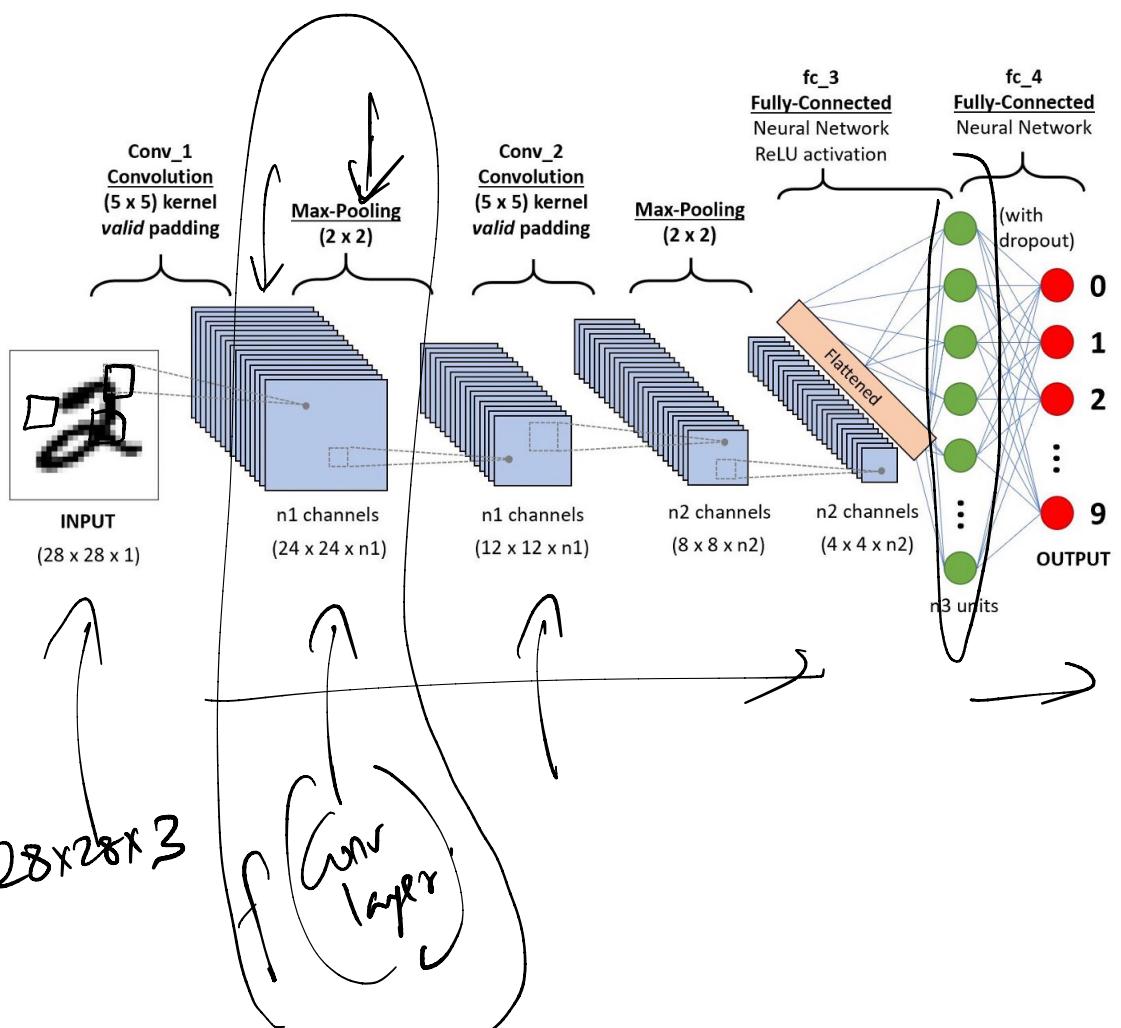
feed-forward.

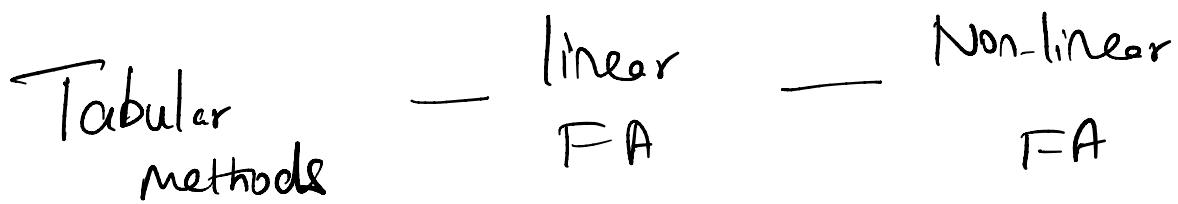
$$\text{loss} = \frac{1}{2}(o - v^{(s)})^2$$

Backpropagation Algorithm

$$\frac{\partial \text{loss}}{\partial w_2} = \frac{\partial \text{loss}}{\partial o} \frac{\partial o}{\partial h_2} \frac{\partial h_2}{\partial w_2}$$

Convolutional Neural Networks





Off-Policy methods with FA:-

Off-Policy Learning :-

goal: to learn a value for for
 a target policy π , given
 data due to a different
behavior policy b .

Prediction: π, b are static and given.

Control: π, b are changing.

$\pi \rightarrow$ greedy pol w.r.t. q .
 $b \rightarrow \epsilon$ -greedy.

2 challenges in off-Policy learning:

- ① target of the update [tabular, FA]
- ② distribution of updates [only in FA]

Semi-gradient methods:-

Per-step important Sampling ratio:

$$P_t = P_{t:t} = \frac{\pi(A_t | s_t)}{b(A_t | s_t)}$$

Semi-gradient off-Policy TD:-

— one-step, state-val. algo

$$\omega_{t+1} = \omega_t + \alpha P_t \delta_t \nabla \hat{v}(s_t; \omega_t)$$

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

Semi-gradient expected Sarsa :-

— one-step, action-value algo

$$\omega_{t+1} = \omega_t + \alpha \delta_t \nabla \hat{q}(s_t, a_t; \omega_t)$$

where $\delta_t = R_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) \hat{q}(s_{t+1}, a; \omega_t)$

— $\hat{q}(s_t, a_t; \omega_t)$

n-step semi-grad. expected Sarsa :-

$$\omega_{t+n} = \omega_{t+n-1} + \alpha P_{t+1} P_{t+2} \dots P_{t+n-1} \left[G_{t:t+n} - \hat{q}(s_t, a_t; \omega_{t+n}) \right]$$

$$\nabla \hat{q}(s_t, a_t; \omega_{t+n})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} +$$

$$\gamma^n \hat{q}(s_{t+n}, a_{t+n}; \omega_{t+n})$$

Examples for off-policy divergence :-

Example 1:



$$\begin{aligned} x(S_1) &= 1 \\ x(S_2) &= 2 \end{aligned}$$

$$w = 10.$$

$$\hat{v}(S_1; w) = 10 \quad \hat{v}(S_2; w) = 20$$

$$\gamma \approx 1 \quad \text{TD error} \sim 10 \quad \alpha = 0.1$$

$$w = 11$$

$$\hat{v}(S_1; w) = 11 \quad \hat{v}(S_2; w) = 22$$

$$\text{TD error} \approx 11 \quad w \approx 12.1$$

TD error:

$$\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}; w_t) - \hat{v}(S_t; w_t)$$

$$= 0 + \gamma 2w_t - w_t$$

$$\delta_t = (2\gamma - 1) w_t$$

off-Policy Semi-grad. TD update:

$$\omega_{t+1} = \omega_t + \alpha P_t \delta_t \nabla J(s_t; \omega_t)$$

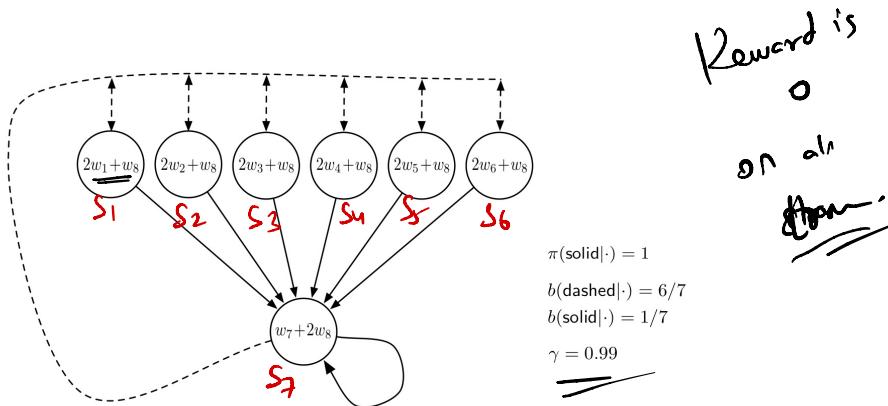
$$= \omega_t + \alpha \gamma (\alpha \delta_t) \omega_t \cdot /$$

$$\omega_{t+1} = \underbrace{(1 + \alpha \gamma)}_{\text{if } (1 + \alpha \gamma) > 1} \omega_t$$

$$\gamma > 0.5$$

Ex 2:

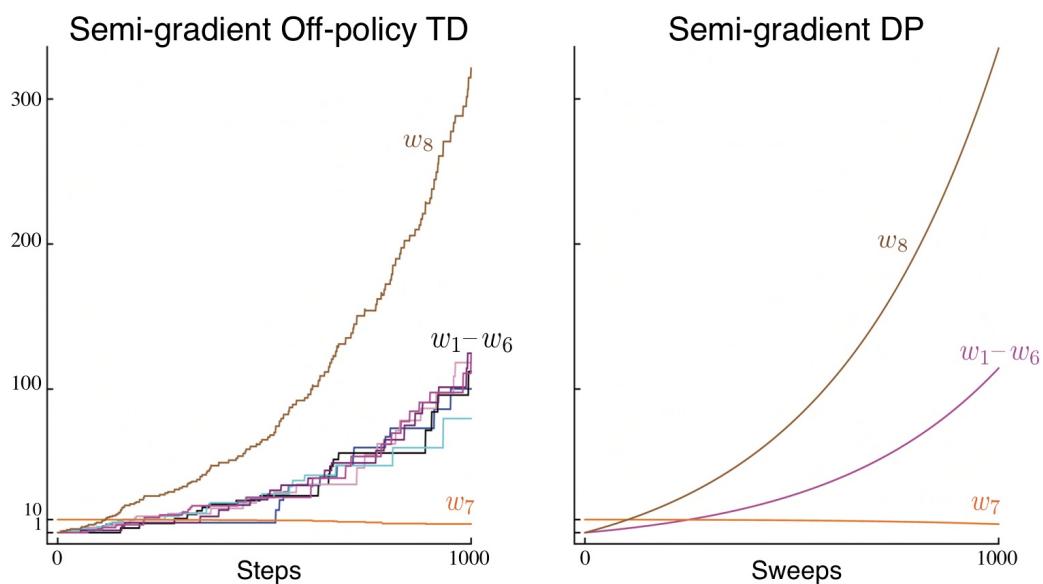
Baird's Counter example.



$$w \in \mathbb{R}^8 \quad s_1 = 2w_1 + w_8$$

$$x(1) = (2, 0, 0, 0, 0, 0, 1)^T$$

$$\sqrt{\pi(s)} = 0 \quad \text{for all } s \Rightarrow w = 0$$



Semi-gradient DP

$$w_{k+1} = w_k + \frac{\alpha}{|S|} \sum_s \left(E_h [R_{t+1} + \gamma V(s_{t+1}; w_k) | s_t = s] - V(s; w_k) \right) \nabla V(s; w_k)$$

F A , Bootstrappy , off-Policy training .

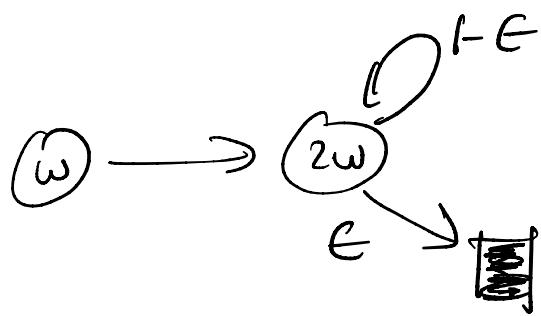
→ Stochastic ? X

→ asynch. ? X

→ one-step return ? X

Ex 3:

Tsitsiklis & Van Roy :



$$w_{k+1} = \underset{w \in \mathbb{R}}{\operatorname{argmin}} \sum_{s \in S} \left[V(s; w) - \mathbb{E}_{\pi} \left[R_{t+1} + \gamma \hat{V}(s_{t+1}; w_k) \right] \right]$$

$$= \underset{w \in \mathbb{R}}{\operatorname{argmin}} (w - \gamma w_k)^2 +$$

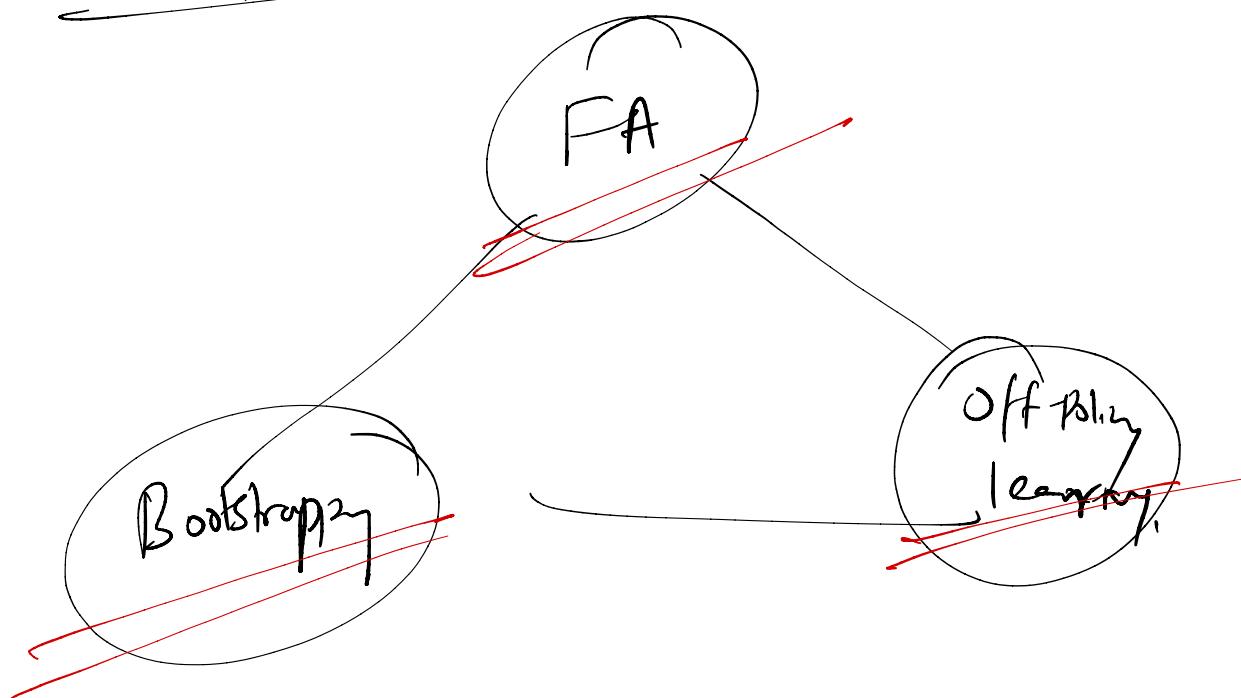
$$(2w - (1-\gamma) \gamma w_k)^2$$

$$= \frac{6-4\gamma}{5} \gamma w_k.$$

$\{w_k\} \rightarrow$ diverge when

$$\gamma > \frac{5}{6-4\gamma}$$
 and w_0 .

the deadly triad:



On/off-Policy	Algo	Tabular	linear FA	Non-linear FA
---------------	------	---------	-----------	---------------

on-policy	MC	✓	✓	✓
	TD	✓	✓	✗

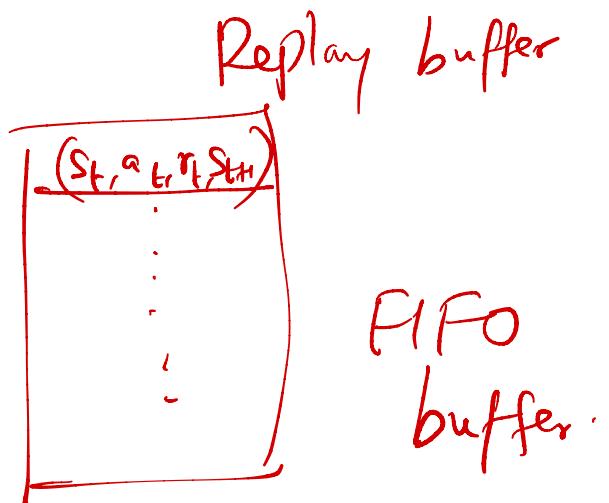
off-Poly	MC	✓	✓	✓
	TD	✓	✗	✗

Deep Q-Networks (DQN):

Q learning + NN FA.

① Experience Replay:

DQN →



- data efficiency
- breaks correlation & reduces the variance in update
- averages the behav. distrib

② Target Network:

$$\delta_t = R_{t+1} + \gamma \max_a Q(s_{t+1}, a; \theta_t^-) - Q(s_t, a_t; \theta_t)$$



θ_t^-



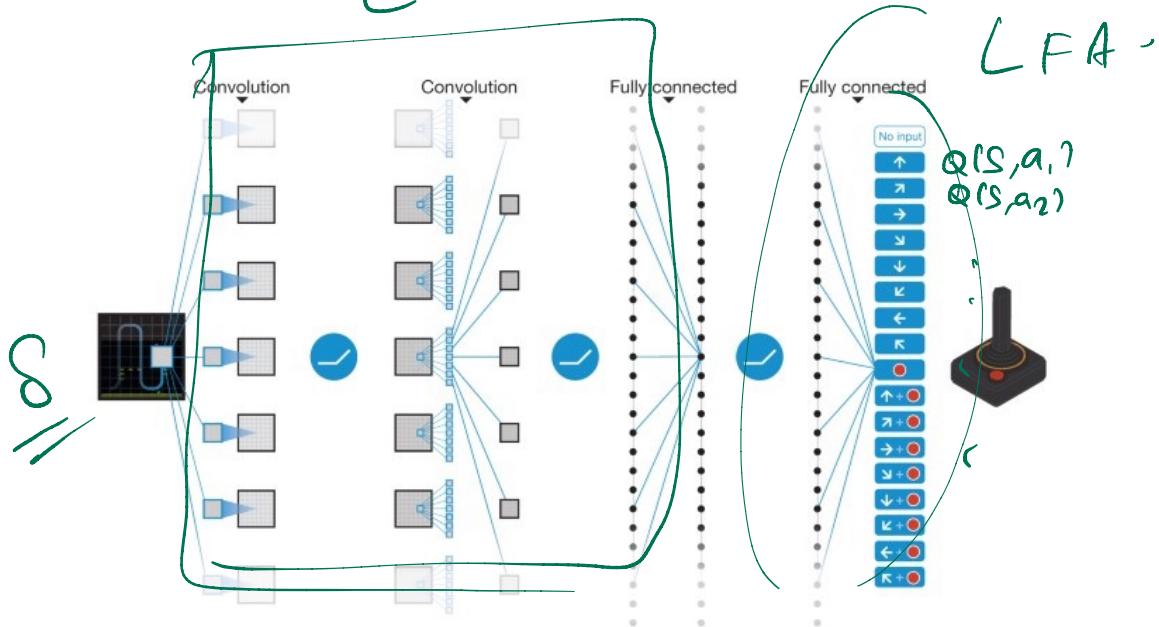
③ Reward clipping.

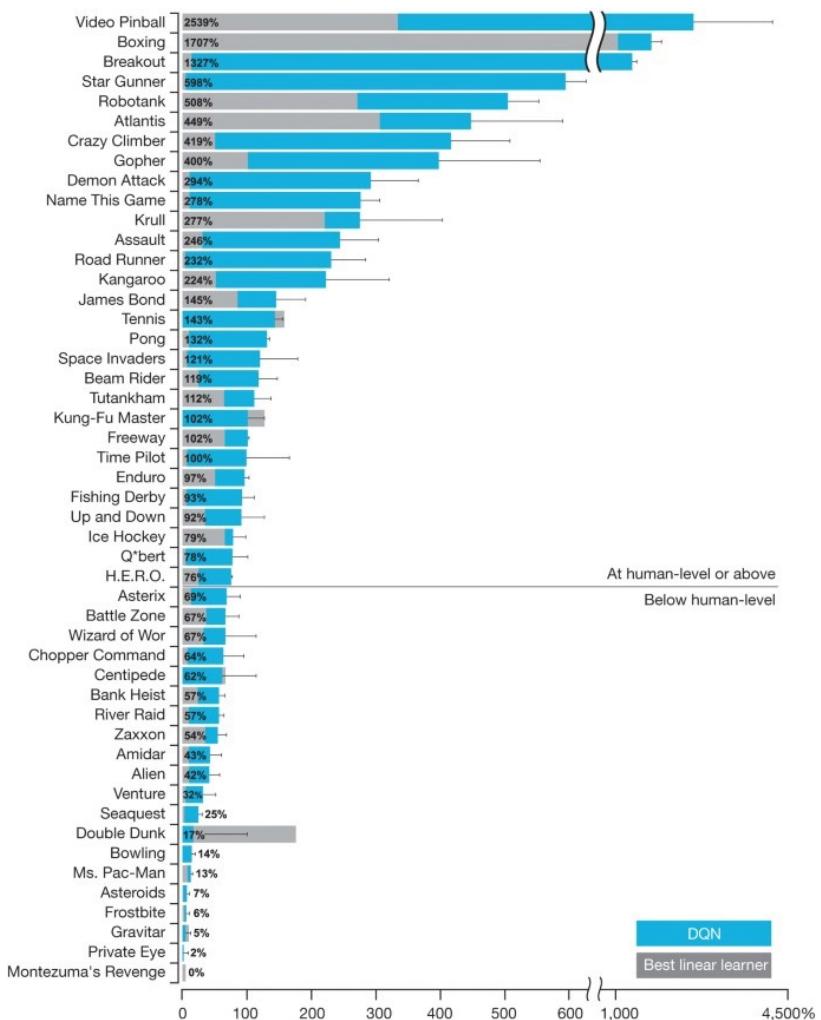
$$(-1, 1)$$

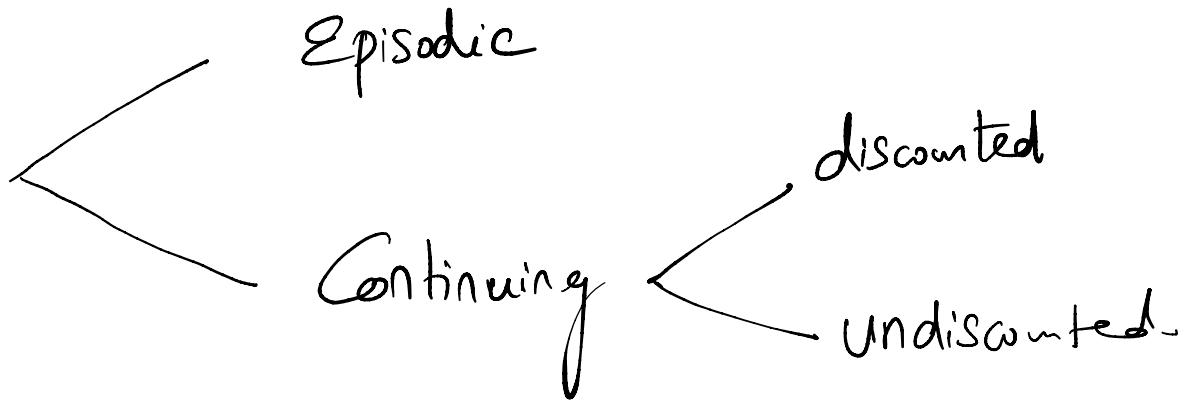
Atari :



CNN







Average reward setting:

$$\begin{aligned}
 \gamma(\pi) &= \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t | S_0, A_0:t-1 \sim \pi] \\
 &= \lim_{t \rightarrow \infty} \mathbb{E}[R_t | S_0, A_0:t-1 \sim \pi] \\
 &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_{s',r} P(s'|r|s,a) \ r
 \end{aligned}$$

$\mu_\pi(s)$ — Steady state distnb.

$$\mu_\pi(s) = \lim_{t \rightarrow \infty} \Pr\{S_t=s | A_0:t-1 \sim \pi\}$$

(Ergodic assumption)

Steady State distn

$$\sum_s M_{\pi}(s) \sum_a \pi(a|s) P(s'|s,a) = M_{\pi}(s')$$

Return

$$G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

differential return.

Bellman eqn

$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{r,s'} P(s'|r|s,a) \left[r - r(\pi) + V_{\pi}(s') \right]$$

$$q_{\pi}(s,a) = \sum_{r,s'} P(s'|r|s,a) \left[r - r(\pi) + \sum_{a'} \pi(a'|s') q_{\pi}(s',a') \right]$$

$$V_{\pi}(s) = \max_a \sum_{r,s'} P(s'|r|s,a) \left[r - \max_{\pi} r(\pi) + V_{\pi}(s') \right]$$

$$q_{\pi}(s,a) = \sum_{r,s'} P(s'|r|s,a) \left[r - \max_{\pi} r(\pi) + \max_{a'} q_{\pi}(s',a') \right]$$

$$\delta_t = R_{t+1} - \bar{R}_{t+1} + \hat{v}(s_{t+1}; w_t) - \hat{v}(s_t; w_t)$$

Differential semi-gradient Sarsa for estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step sizes $\alpha, \beta > 0$, small $\varepsilon > 0$

Initialize value-function weights $\mathbf{w} \in \mathbb{R}^d$ arbitrarily (e.g., $\mathbf{w} = \mathbf{0}$)

Initialize average reward estimate $\bar{R} \in \mathbb{R}$ arbitrarily (e.g., $\bar{R} = 0$)

Initialize state S , and action A

Loop for each step:

Take action A , observe R, S'

Choose A' as a function of $\hat{q}(S', \cdot, \mathbf{w})$ (e.g., ε -greedy)

$$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$$

$$\bar{R} \leftarrow \bar{R} + \beta \delta$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$$

$$S \leftarrow S'$$

$$A \leftarrow A'$$

Depreciating discounted setting :-

~~+ $\gamma \pi(s')$~~

$$\cancel{\gamma \pi(s')}$$

$$J(\pi) = \sum_s \pi_\pi(s) v_\pi^\pi(s)$$

$$= \sum_s \pi_\pi(s) \sum_a \pi(a|s) \sum_{s'} p(s'|r|s,a) [r + \gamma v_\pi^\pi(s')]$$

$$= r(\pi) + \sum_s \pi_\pi(s) \sum_a \pi(a|s) \sum_{s',r} p(s'|r|s,a) \gamma v_\pi^\pi(s')$$

$$= r(\pi) + \gamma \sum_{s'} v_\pi^\pi(s') \sum_s \pi_\pi(s) \sum_a \pi(a|s) p(s'|s,a)$$

$$= r(\pi) + \gamma \sum_{s'} v_\pi^\gamma(s') \underbrace{\sum_a M_\pi(s) \sum_a \pi(a|s) p(s'|s,a)}$$

$$= r(\pi) + \gamma \sum_{s'} v_\pi^\gamma(s') M_\pi(s')$$

$$= r(\pi) + \gamma J(\pi)$$

$$= r(\pi) + \gamma r(\pi) + \gamma^2 J(\pi)$$

$$= r(\pi) + \gamma r(\pi) + \gamma^2 r(\pi) + \dots$$

$$= \frac{1}{1-\gamma} r(\pi)$$