

6. Value function approximation

6-1 large problem

① go 10^{170} states

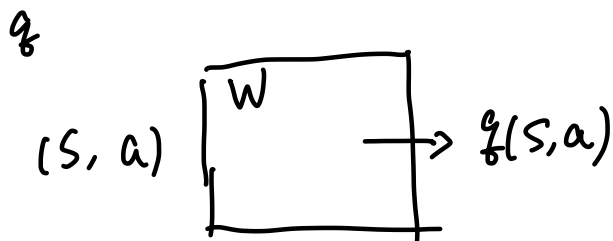
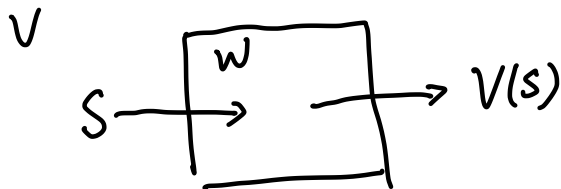
② Helicopter continuous states

	s_1	s_2	...
$V(s)$			

Impossible list all states in a table

How do scale-up?

$$V(s) \longleftrightarrow V(s, w)$$



update w
instead of V, q

6-2 Approximation function

① Linear function

② NN

6-3 Incremental Method

6-3.1 Gradient Descent (GD)

Define loss

$$J(w) = E_{\pi} [(V_{\pi}(s) - \hat{V}(s, w))^2]$$

Supervised learning $V_{\pi}(s)$: real $\hat{V}(s, w)$: approximation

GD $W^* = \arg \min_w J(w)$

$$\Delta W = -\frac{\alpha}{2} \nabla_w J(w) = \alpha E_{\pi} [(V_{\pi}(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)]$$

SGD Batch 1

$$\Delta W_i = \alpha (V_{\pi}(s_i) - \hat{V}(s_i, w)) \nabla_w \hat{V}(s_i, w)$$

Batch |B|

$$\Delta W = \frac{1}{|B|} \sum \Delta W_i$$

6.3.2 Linear combination of feature

$$\hat{V}(s, w) = X(s)^T w = \sum_{j=1}^n X_j(s) w_j$$

$$\Delta W = \alpha (V_{\pi}(s) - \hat{V}(s, w)) X(s)$$

$$\text{if } X(s) = \begin{pmatrix} 1 (s=s_1) \\ 1 (s=s_2) \\ \vdots \\ 1 (s=s_n) \end{pmatrix}$$

previous Table Method

6.3.3 Study $V_a(s)$

① MC $\Delta W = \alpha(G_t - \hat{V}(S_t, w)) \nabla_w \hat{V}(S_t, w)$

② TD(0) $\Delta W = \alpha(R_{t+1} + \gamma \hat{V}(S_{t+1}, w) - \hat{V}(S_t, w)) \nabla_w \hat{V}(S_t, w)$

③ TD(λ) $\Delta W = \alpha(G_t^\lambda - \hat{V}(S_t, w)) \nabla_w \hat{V}(S_t, w)$

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

$$V(s) = E(G_t | S_t = s)$$

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

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} G_t^{(n)} \lambda^{n-1}$$

Dataset

For MC: $\{ \langle S_i, G_i \rangle \}_{i=1}^T$

TD: $\{ \langle S_i, R_{i+1} + \gamma V(S_{i+1}) \rangle \}_{i=1}^{T-1}$

TD(λ): $\{ \langle S_i, G_i^\lambda \rangle \}_{i=1}^{T-1}$

Very important problem

← This is not an exact GD?

Why?

$$J(w) = E_{\pi} [L((V_{\pi}(s) - \hat{V}(s, w))^2)]$$

$$TD(0) \quad J(w) = E_{\pi} [(R_{t+1} + \gamma \hat{V}(S_{t+1}, w) - \hat{V}(S_t, w))^2]$$

Exact GD, Gradient on $\hat{V}(S_{t+1}, w), \hat{V}(S_t, w)$

6.4 Control problem

$$\hat{q}_{\pi}^{\gamma}(S, A, w) \rightarrow q_{\pi}^{\gamma}(S, A)$$

$$J(w) = E_{\pi} [L(\underbrace{q_{\pi}^{\gamma}(S, A)} - \hat{q}_{\pi}^{\gamma}(S, A, w))^2]$$

$$\text{For MC: } \Delta w = \alpha (G_{t+1} - \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)) \nabla_w \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)$$

$$TD(0): \Delta w = \alpha (R_{t+1} + \gamma \hat{q}_{\pi}^{\gamma}(S_{t+1}, A_{t+1}, w) - \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)) \nabla_w \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)$$

$$TD(\lambda): \Delta w = \alpha (q_t^{\pi} - \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)) \nabla_w \hat{q}_{\pi}^{\gamma}(S_t, A_t, w)$$

Convergence of Control Algorithms

Algorithm	Table Lookup	Linear	Non-Linear
Monte-Carlo Control	✓	(✓)	✗
Sarsa	✓	(✓)	✗
Q-learning	✓	✗	✗
Gradient Q-learning	✓	✓	✗

(✓) = chatters around near-optimal value function

6.5 Replay or not

$$\mathcal{D} = \{ \langle s_1, v_1^{\pi} \rangle, \langle s_2, v_2^{\pi} \rangle, \dots, \langle s_T, v_T^{\pi} \rangle \}$$

Not Replay

$\langle s_i, v_i^{\pi} \rangle$ one by one following Time order

Replay

every time randomly select
 $\langle s_i, v_i^{\pi} \rangle$ from \mathcal{D}

Apply SGD

6.6 Many states, but relative few actions

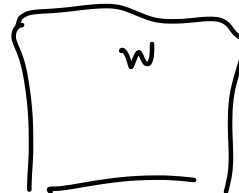
Experience Replay in Deep Q-Networks (DQN)

DQN uses **experience replay** and **fixed Q-targets**

- Take action a_t according to ϵ -greedy policy
- Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory \mathcal{D}
- Sample random mini-batch of transitions (s, a, r, s') from \mathcal{D}
- Compute Q-learning targets w.r.t. old, fixed parameters $\underline{w^-}$
- Optimise MSE between Q-network and Q-learning targets

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[\left(r + \gamma \max_{a'} Q(s', a'; w_i^-) - Q(s, a; w_i) \right)^2 \right]$$

- Using variant of stochastic gradient descent



$$\begin{array}{l} \text{第 } 1 = \\ \leftarrow W \quad W = \bar{w} \\ \quad \quad W + \Delta w \end{array}$$

$$\Delta w = \nabla_w L(\bar{w}, w)$$

$$\begin{array}{l} | \\ | \\ | \\ K \\ \leftarrow W \quad W + \Delta w \end{array}$$

$$\Delta w = \nabla_w L(\bar{w}, w)$$

$$\begin{array}{l} (k+1) \\ \quad \quad \bar{w} = w \\ W = W + \Delta w \end{array}$$

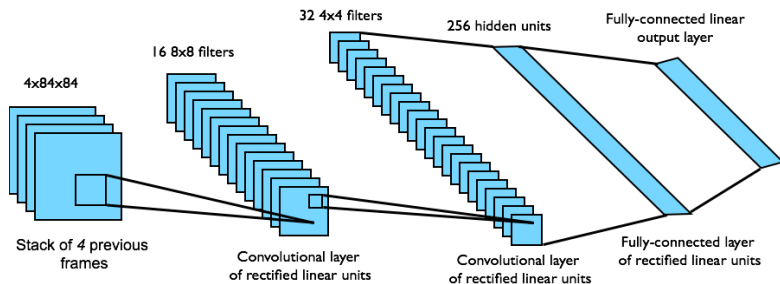
$$\Delta w = \nabla_w L(\bar{w}, w)$$

$$\begin{array}{l} | \\ | \\ | \\ W = \dots \end{array}$$

.....

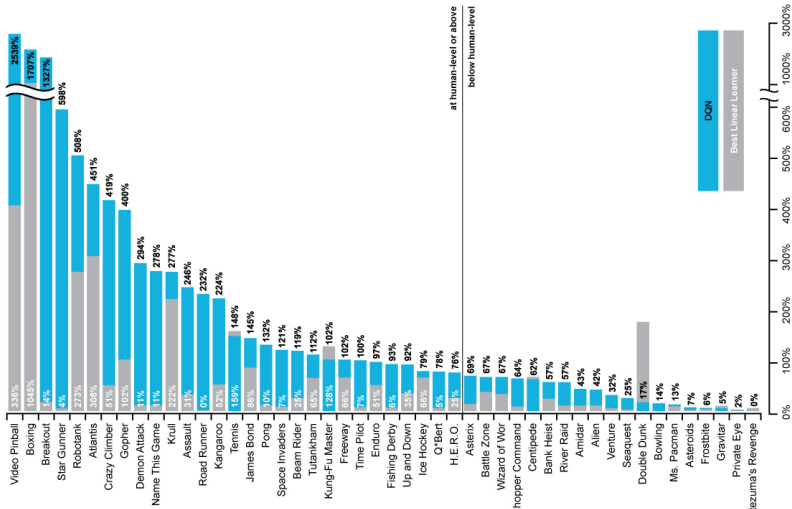
DQN in Atari

- End-to-end learning of values $Q(s, a)$ from pixels s
- Input state s is stack of raw pixels from last 4 frames
- Output is $Q(s, a)$ for 18 joystick/button positions
- Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

DQN Results in Atari



How much does DQN help?

	Replay Fixed-Q	Replay Q-learning	No replay Fixed-Q	No replay Q-learning
Breakout	316.81	240.73	10.16	3.17
Enduro	1006.3	831.25	141.89	29.1
River Raid	7446.62	4102.81	2867.66	1453.02
Seaquest	2894.4	822.55	1003	275.81
Space Invaders	1088.94	826.33	373.22	301.99