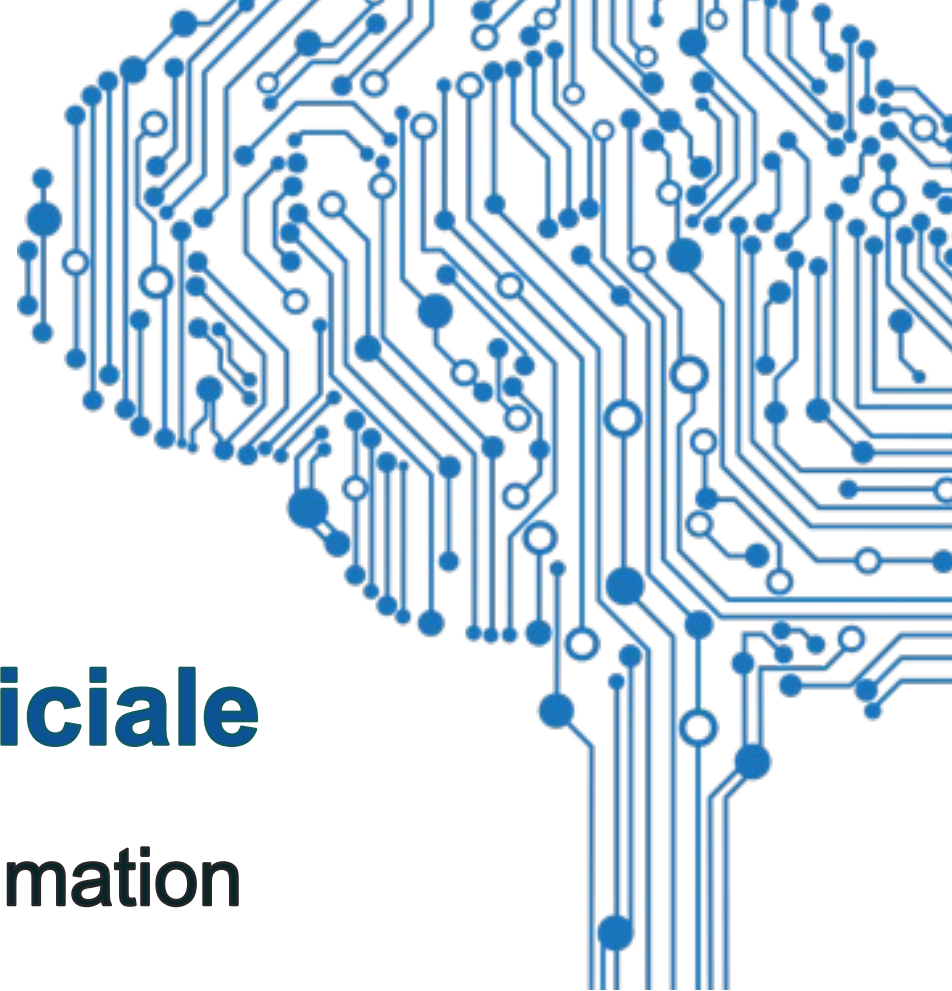




UNIVERSITÀ DEGLI STUDI DI SALERNO  
**DIPARTIMENTO DI INFORMATICA**



# **Intelligenza Artificiale**

## **Value Function Approximation**

# Metodi Batch

# Batch Reinforcement Learning

---

- ▶ Gradient descent semplice ed accattivante
- ▶ Ma non è efficiente dal punto di vista del campionamento
- ▶ I metodi batch cercano di trovare la value function più adatta
- ▶ Data l'esperienza dell'agente (dati di addestramento)

# Least Squares Prediction

---

- ▶ Data una value function approximation  $\hat{v}(s; \mathbf{w}) \approx v_{\pi}(s)$  e l'esperienza  $D$  costituita da coppie  $\langle stato, valore \rangle$

$$\mathcal{D} = \{\langle S_1, v_1^{\pi} \rangle, \langle S_2, v_2^{\pi} \rangle, \dots, \langle S_T, v_T^{\pi} \rangle\}$$

- ▶ Quali parametri  $\mathbf{w}$  offrono il miglior risultato di adattamento della value function  $\hat{v}(s; \mathbf{w})$ ?
- ▶ Gli algoritmi **Least Squares** individuano il vettore dei parametri  $\mathbf{w}$  che minimizza l'errore quadratico medio tra  $\hat{v}(s; \mathbf{w})$  e i valori target  $v_{\pi}(s)$

$$LS(\mathbf{w}) = \sum_{t=1}^T (v^{\pi} - \hat{v}(s_t; \mathbf{w}))^2 = \mathbb{E}_{\mathcal{D}} \left[ (v^{\pi} - \hat{v}(s; \mathbf{w}))^2 \right]$$

# Experience Replay

# SGD con Experience Replay

---

- ▶ Data una value function approximation  $\hat{v}(s; \mathbf{w}) \approx v_\pi(s)$  e l'esperienza  $\mathcal{D}$  costituita da coppie  $\langle \text{stato}, \text{valore} \rangle$

$$\mathcal{D} = \{\langle S_1, v_1^\pi \rangle, \langle S_2, v_2^\pi \rangle, \dots, \langle S_T, v_T^\pi \rangle\}$$

- ▶ Ripetere
  1. Campiona stato e valore dall'esperienza  $\mathcal{D}$

$$\langle s, v^\pi \rangle \sim \mathcal{D}$$

2. Esegui lo stochastic gradient descent update

$$\Delta \mathbf{w} = \alpha (v^\pi - \hat{v}(s; \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s; \mathbf{w})$$

- ▶ Converge alla soluzione least squares

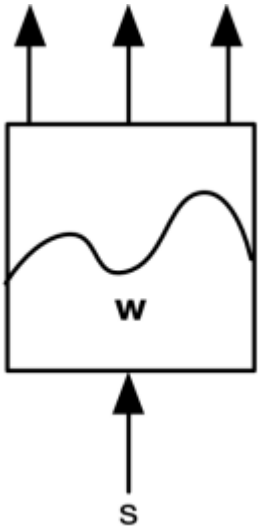
$$\mathbf{w} = \arg \min_{\mathbf{w}} LS(\mathbf{w})$$

# Deep Q-Networks (DQN)

---

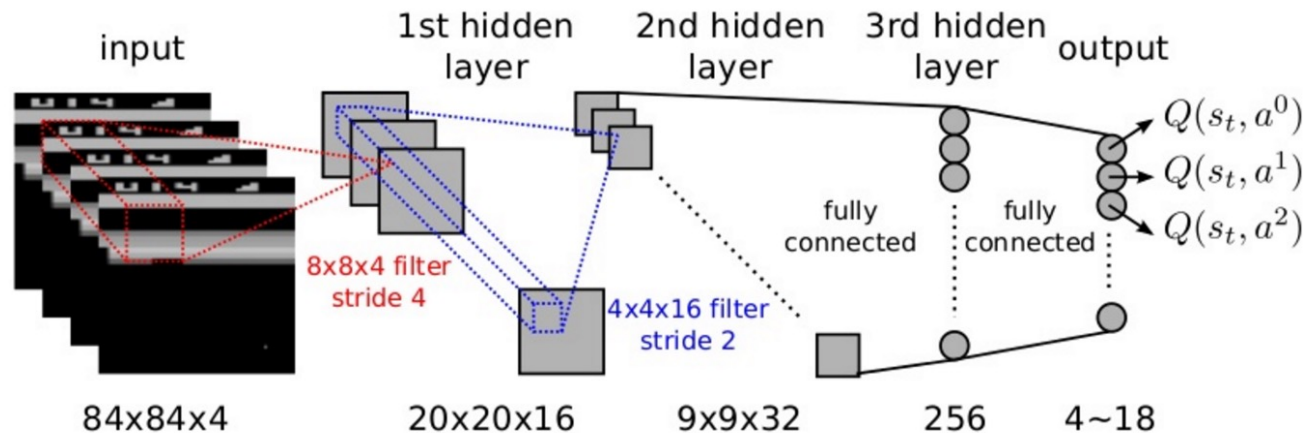
- ▶ DQN usa **experience replay** e **Q-target fissati**
- ▶ Esegue l'azione  $a_t$  in base alla policy  $\epsilon$ -greedy
- ▶ Memorizza le transizioni  $(s_t, a_t, r_{t+1}, s_{t+1})$  nella **replay memory  $D$**
- ▶ Campiona un **mini-batch casuale di transizioni  $(s, a, r, s')$**
- ▶ Calcola i Q-learning target rispetto ai vecchi parametri fissati  $w^-$
- ▶ Ottimizza la MSE tra la Q-network e i Q-learning target
$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[ \left( r + \gamma \max_{a'} Q(s', a'; w_i^-) - Q(s, a; w_i) \right)^2 \right]$$
- ▶ Utilizza una variante del ***stochastic gradient descent***

# Atari-DQN



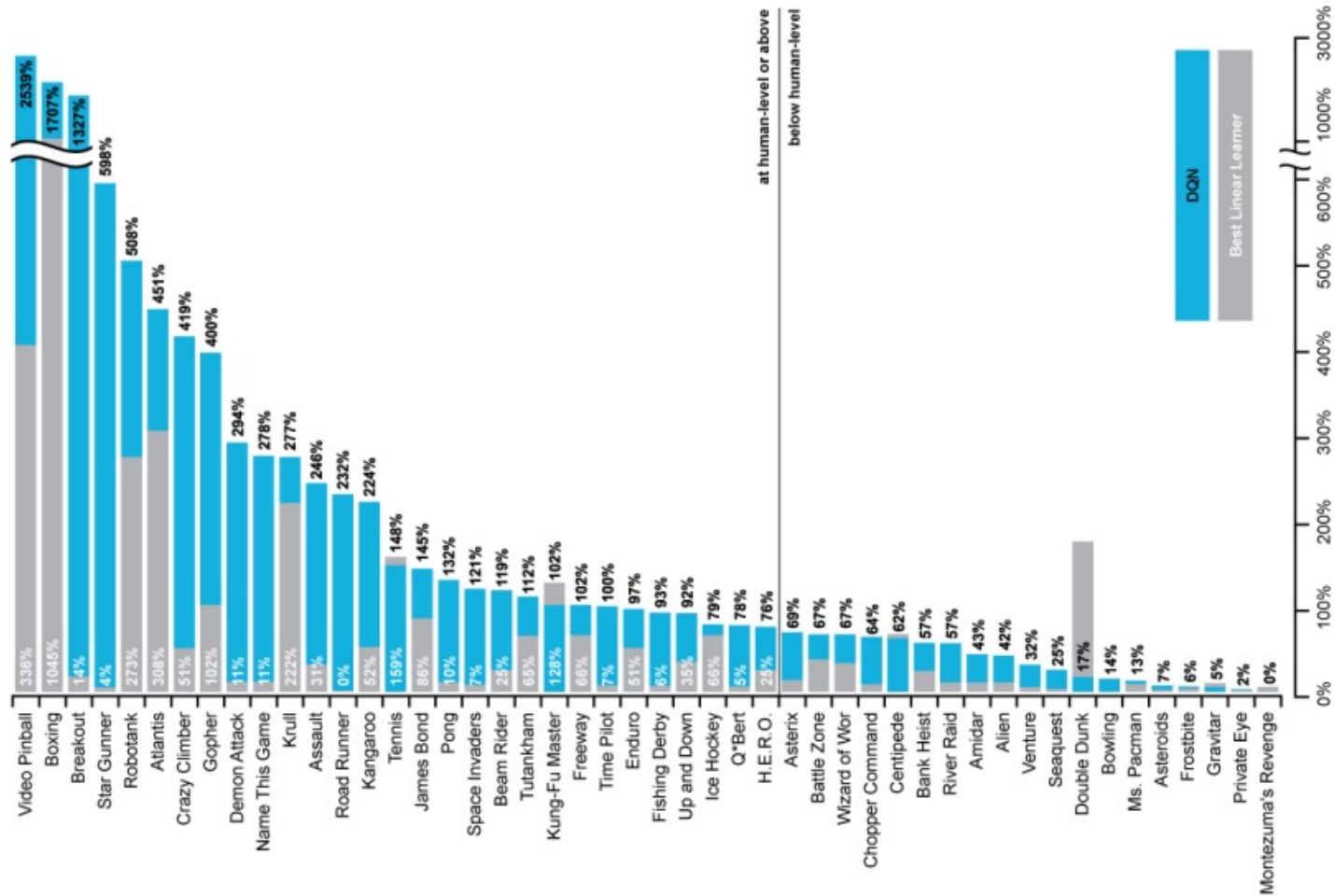
- ▶ Apprendimento end-to-end dei valori  $Q(s, a)$  dai pixel  $s$
- ▶ Lo stato di input  $s$  è uno stack di pixel grezzi degli ultimi 4 fotogrammi
- ▶ L'output è  $Q(s, a)$  per 18 posizioni di joystick/pulsante
- ▶ La ricompensa è la variazione del punteggio per quello step

L'Architettura di rete e gli iperparametri sono gli stessi per tutti i giochi





# Atari-DQN (Risultati)



# Least Squares Prediction and Control

# Linear Least Squares Prediction

---

- ▶ L'experience replay trova una soluzione least squares, ma può richiedere molte iterazioni
- ▶ Utilizzo della linear value function approximation  $\hat{v}(s; \mathbf{w}) = \mathbf{x}(s)^T \mathbf{w}$
- ▶ Possiamo risolvere la soluzione least squares direttamente

# Linear Least Squares Prediction - Batch

---

- ▶ Al valore minimo di  $LS(\mathbf{w})$ , l'expected update deve essere pari a zero

$$\begin{aligned}\mathbb{E}_{\mathcal{D}}[\Delta \mathbf{w}] &= 0 \\ \alpha \sum_{t=1}^T \mathbf{x}(S_t)(\mathbf{v}_t^{\pi} - \mathbf{x}(S_t)^T \mathbf{w}) &= 0 \\ \mathbf{w} &= \left( \sum_{t=1}^T \mathbf{x}(S_t) \mathbf{x}(S_t)^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) \mathbf{v}_t^{\pi}\end{aligned}$$

- ▶ Costi diretti ( $N^3$ ) e incrementali ( $N^2$ )

# Linear Least Squares Prediction - Algoritmi

---

- ▶ Non conosciamo il true value  $v_t^\pi$
- ▶ I dati di addestramento utilizzano campioni rumorosi o distorti di  $v_t^\pi$ 
  - ▶ Least Squares Monte-Carlo (**LSMC**) usa return  $v_t^\pi \approx G_t$
  - ▶ Least Squares TD (**LSTD**) usa TD target  $v_t^\pi \approx R_{t+1} + \gamma \hat{v}(S_{t+1}; \mathbf{w})$
  - ▶ Least Squares TD( $\lambda$ ) (**LSTD( $\lambda$ )**) usa  $\lambda$ -return  $v_t^\pi \approx G_t^\lambda$
- ▶ In ogni caso bisogna risolvere direttamente il punto fisso di MC/TD/TD( $\lambda$ )

# LS Updates

---

LSMC	$\mathbf{w} = \left( \sum_{t=1}^T \mathbf{x}(S_t) \mathbf{x}(S_t)^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) G_t$
LSTD	$\mathbf{w} = \left( \sum_{t=1}^T \mathbf{x}(S_t) (\mathbf{x}(S_t) - \gamma \mathbf{x}(S_{t+1}))^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) R_{t+1}$
LSTD( $\lambda$ )	$\mathbf{w} = \left( \sum_{t=1}^T E_t (\mathbf{x}(S_t) - \gamma \mathbf{x}(S_{t+1}))^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t) R_{t+1}$

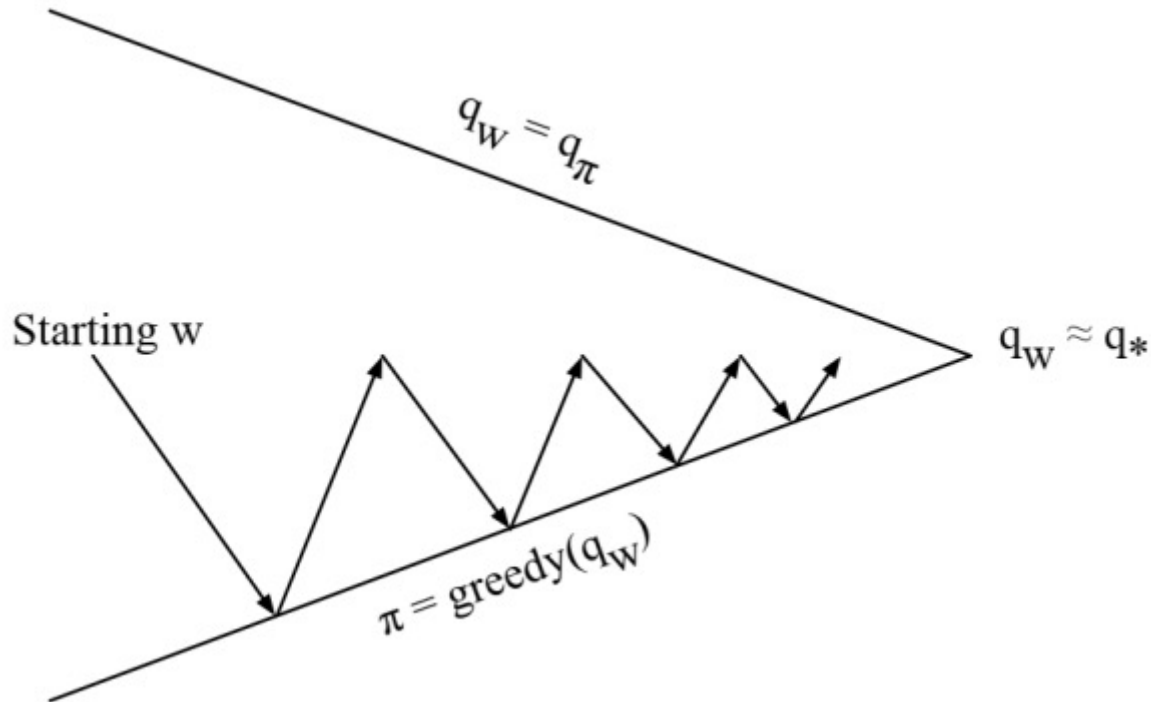
# Convergenza degli Algoritmi di Linear Least Squares Prediction

---

On/Off-Policy	Algorithm	Table Lookup	Linear	Non-Linear
On-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	✓	X
	LSTD	✓	✓	-
Off-Policy	MC	✓	✓	✓
	LSMC	✓	✓	-
	TD	✓	X	X
	LSTD	✓	✓	-

# Least Squares Control

---



- ▶ Valutazione della Policy – Policy evaluation attraverso **least squares Q-learning**
- ▶ Miglioramento della Policy – Greedy policy improvement



# Least Squares Action-Value Function Approximation

---

- ▶ Approssimare  $q_\pi(s, a)$  usando combinazioni lineari delle feature  $\mathbf{x}(s, a)$

$$\hat{q}(s, a; \mathbf{w}) = \mathbf{x}(s, a)^T \mathbf{w} \approx q_\pi(s, a)$$

- ▶ Minimizzare il least squares error tra  $\hat{q}(s, a; \mathbf{w})$  e  $q_\pi(s, a)$ 
  - ▶ Dall'esperienza generata usando la policy  $\pi$
  - ▶ Composta da coppie  $\langle (stato, azione), valore \rangle$

$$\mathcal{D} = \{ \langle (s_1, a_1), v_1^\pi \rangle, \dots, \langle (s_T, a_T), v_T^\pi \rangle \}$$

# Least Squares Control

---

- ▶ Per la policy evaluation, vogliamo utilizzare in modo efficiente tutta l'esperienza acquisita
- ▶ Per il control, vogliamo anche migliorare la policy
- ▶ Questa esperienza è generata da molte policy
- ▶ Quindi, per valutare  $q_{\pi}(s, a)$  dobbiamo apprendere off-policy
- ▶ Utilizziamo la stessa idea del Q-learning:
  - ▶ Utilizziamo l'esperienza generata dalla vecchia policy  $S_t, A_t, R_{t+1}, S_{t+1} \sim \pi_{old}$
  - ▶ Consideriamo un'azione alternativa  $A' \sim \pi_{new}(S_{t+1})$
  - ▶ Aggiorniamo  $\hat{q}(S_t, A_t; \mathbf{w})$  rispetto al valore dell'azione alternativa  $R_{t+1} + \gamma \hat{q}(S_{t+1}, A'; \mathbf{w})$

# Least Squares Q-Learning

---

- ▶ Si consideri il seguente linear Q-learning update

$$\begin{aligned}\delta_t &= R_{t+1} + \gamma \hat{q}(S_{t+1}, \pi(S_{t+1}); \mathbf{w}) - \hat{q}(S_t, A_t; \mathbf{w}) \\ \Delta \mathbf{w}_t &= \alpha \delta_t \mathbf{x}(S_t, A_t)\end{aligned}$$

- ▶ Algoritmo LSTDQ: risoluzione per aggiornamento totale uguale a zero

$$\begin{aligned}\Delta \mathbf{w}_t &= 0 \\ \mathbf{w} &= \left( \sum_{t=1}^T \mathbf{x}(S_t, A_t) (\mathbf{x}(S_t, A_t) - \gamma \mathbf{x}(S_{t+1}, \pi(S_{t+1})))^T \right)^{-1} \sum_{t=1}^T \mathbf{x}(S_t, A_t) R_{t+1}\end{aligned}$$

# Convergenza degli Algoritmi di Control

---

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

(✓) = si aggira intorno a una value-function quasi ottimale

# Take home messages

---

- ▶ Value function approximation: scalare RL a problemi di dimensioni reali
- ▶ Utilizzare MC, TD e TD( $\lambda$ ) per generare campioni per l'addestramento degli (action) value approximator
  - ▶ Stochastic Gradient Descent (incrementale)
  - ▶ Least Squares (batch)
- ▶ Experience replay dei dati memorizzati e aggiornamenti con target dei parametri meno recenti (DQN)
- ▶ Linear least squares fornisce una soluzione in forma chiusa
- ▶ La convergenza dell'apprendimento potrebbe essere complicata