# Lecture 28

## CS 59000-RL

### Agenda
- Competetive games
- Bellman Residual minimization
- Deep Q Network (DQN)
  { DDGN
  { BDQN
- Wrap up

---

So far single agent MDP. $\Pi \to \eta(\Pi) \to \max_{\Pi} \eta(\Pi)$

what about two player game ?, e.g. competing
against
each other?

Two players _ play chess

- Play ping pong
- competing for the same resources
- competing for a same radio channel.
- racing
- :

$\to$ Competetive Markov Decision processes.

- zero sum games.

Two players at each time step make
their decision abter observing $X_t$

$$P1 \rightarrow A^1 \sim \Pi^1(X_t) \qquad\qquad r_t^1$$

$$P2 \rightarrow A^2 \sim \Pi^2(X_2) \quad\Rightarrow X_{t+1}, \qquad r_t^2$$

$\Rightarrow$ zero sum game $r_t^1 + r_t^2 = 1$

$$\Rightarrow r_t = r_t^1 = -r_t^2$$

$$\gamma^1(\Pi^1, \Pi^2) \Rightarrow \text{objective of } P1$$

$$\gamma^2(\Pi^1, \Pi^2) \Rightarrow \text{objective of } P2$$

$$\eta(\Pi^1, \Pi^2) = \eta^1(\Pi^1, \Pi^2) = -\eta^2(\Pi^1, \Pi^2) = E_{\Pi_1, \Pi_2}\left[\sum r_t\right]$$

$P1$; first player aims to maximize $\eta^1$ and respectively $\eta$

$P2$; second player aims to maximize $\eta^2$, ergo minimize $\eta$

$$\max_{\Pi_1} \eta(\Pi^1, \Pi^2), \qquad \min_{\Pi_2} \eta(\Pi^1, \Pi^2)$$
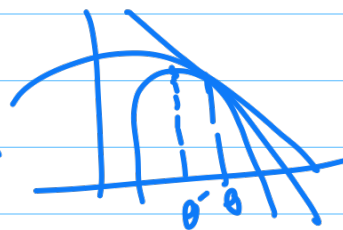
How we develope policy gradient method for
such setting?

How we did it for single player MDP?

$$\eta(\Pi_\theta) \Rightarrow \theta' = \theta + \underset{\Delta\theta}{\arg\max}\ \Delta\theta^T \nabla\eta(\theta) - \frac{1}{2\alpha}\|\Delta\theta\|^2$$

linear approximation of the objective

$$\Rightarrow \theta' = \theta + \alpha \nabla\eta(\theta)$$

Tylor expansion



$\theta' \to \Pi^1_{\theta_1} ; \theta^2 \to \Pi^2_{\theta_2}$

Naive: $\theta^1 \leftarrow \theta^1 + \underset{\Delta\theta^1}{\arg\max}\ \Delta\theta^{1^T} \nabla_{\theta^1}\eta(\theta^1,\theta^2) - \frac{1}{2\alpha}\|\Delta\theta^1\|^2$

$\theta^2 \leftarrow \theta^2 + \arg\min\ \Delta\theta^{2^T} \nabla_{\theta^2}\eta(\theta^1,\theta^2) + \frac{1}{2\alpha}\|\Delta\theta^2\|^2$

Approximate $\eta(\theta^1,\theta^2)$ such that the local game is linear in each player's parameter.

$$\eta(\theta^1+\Delta\theta^1, \theta^2+\Delta\theta^2) = \eta(\theta^1,\theta^2)$$

$$+\ \Delta\theta^1 \nabla_{\theta_1}\eta + \Delta\theta^2 \nabla_{\theta^2}\eta + \Delta\theta_1^T \nabla_{\theta_1\theta_2}\eta \Delta\theta_2$$

$$+\ \Delta\theta^{2^T} \nabla_{\theta_2\theta_1}\eta \Delta\theta^1 + \text{other}$$

$\Rightarrow$ This is a bilinear approximation of the game

$$\to \theta^1 \leftarrow \theta^1 + \underset{\Delta\theta^1}{\text{argmax}}\ \Delta\theta^{1T}\nabla_{\theta^1}\eta + \Delta\theta^1 \nabla_{\theta^1\theta^2}\eta\ \Delta\theta_2 - \frac{1}{2\alpha}\|\Delta\theta_1\|^2$$

$$\to \theta^2 \leftarrow \theta^2 + \underset{\Delta\theta^2}{\text{argmin}}\ \Delta\theta^{2T}\nabla_{\theta^2}\eta + \Delta\theta^2 \nabla_{\theta^2\theta^1}\eta\ \Delta\theta_1 + \frac{1}{2\alpha}\|\Delta\theta_2\|^2$$

bilinear local game $\to$
has close form solution, which is the
Nash equilibrium of this bilinear game

$$\Rightarrow \theta^1 \leftarrow \theta^1 + \alpha\left(I + \alpha^2 \nabla_{\theta^1\theta^2}\eta\ \nabla_{\theta^2\theta^1}\eta\right)^{-1}\left(\nabla_{\theta^1}\eta - \alpha\nabla_{\theta^1\theta^2}\nabla_{\theta^2}\eta\right)$$

$$\theta^2 \leftarrow \theta^2 - \alpha\left(I + \alpha^2\nabla_{\theta^2\theta^1}\eta\nabla_{\theta^1\theta^2}\eta\right)^{-1}\left(\nabla_{\theta^2}\eta + \alpha\nabla_{\theta^2\theta^1}\nabla_{\theta^1}\eta\right)$$

It is called competetive policy optimization principle
(CoPO)          Prajapt et al 2620

---

Bellman residual minimization

we turn the problem of learning Q function
to somewhat clasical regression problem.

$\to$ follow a policy $\to$ $x_t, a_t, r_t, x_{t+1}$

$$\Rightarrow \left\|f(x_t, a_t) - \left(r_t + \gamma \max_{a'} f(x_{t+1}, a')\right)\right\|^2$$

$f(x) \to y$   $\|f(x) - y\|^2 \to r_t + \gamma f(x_{t+1})$

Notes:> in general the learned $f$ is
a biased estimate of $Q$.
↳ well behaved in Linear $Q$

This very objective function is the same used

in Deep $Q$ Network (DQN)

→ To tackle Atari games
→ To tackle Go

↳ Boosts of the field Deep RL

---

# DQN
## Initialize $Q$

1) run epsilon greedy policy
$$\begin{cases} \text{argmax } Q(n_t, a) & 1-\varepsilon \\ \text{uniform random} & \varepsilon \end{cases}$$

2) $x_t, a_t, r_t, n_{t+1}$

3) update $Q \Rightarrow \left\| Q(n_t, a_t) - r_t - \gamma \max_{a'} Q(n_{t+1}, a') \right\|^2$

4) using gradient descent

Improve. $Q^{target}$

For line 3 we we $\|Q(x_b, a_b) - r_b - \gamma \max_{a'} Q^{target}(x_{b+1}, a')\|^2$

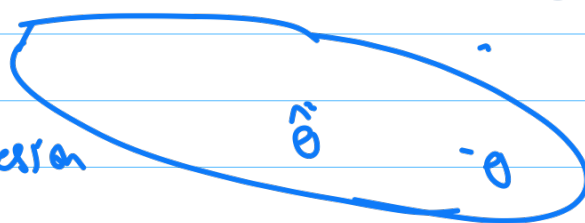line 5) update $\underline{Q^{target}} = \underline{Q}$ once in a while.

Mnih 2015

---

In linear bandit $\rightarrow$ Thompson sampling

$$\theta^T \phi(x, a) \qquad \theta \sim Normal(\underset{mean}{\underbrace{m}}, \underset{Cov}{\underbrace{b}})$$

Bayesian linear regression

$\hat{\hat{\theta}}$     $-\theta$

---

$$x \rightarrow W \phi(x) \rightarrow Q(x, a)$$

learn $\phi$ using gradient descent

learn $W$ using Bayesian linear regression.
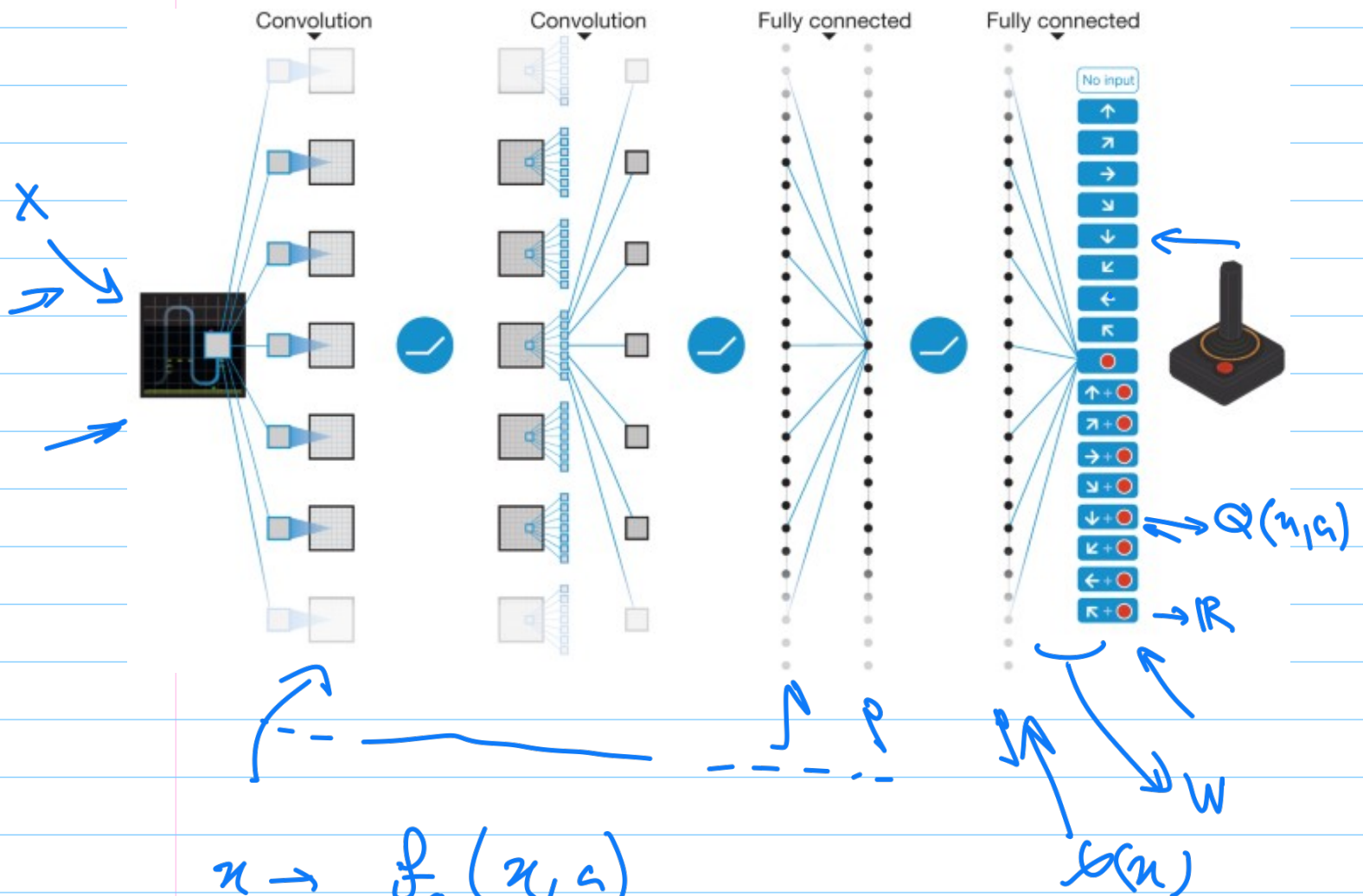
For line 1) we do Thompson sampling.

3) $\left\|\left(W_{mean} \phi(x_b)\right) - r_i - \max_{a'} Q^{target}(x_{b+1}, a')\right\|^2$

learn $\phi$

# This is called Bayessian DQN (BDQN)

---

- Measure theory
- Bandits.
- MDP
- PCMDP
- Control setting
- Model based
- Model free Value based
- Model free Policy based
- Competetisu optimization
- Deep RL
- off policy learning

$$x \to f_\theta(x, a)$$

$$f_\theta(x, a) = \left( W^T \phi(x) \right)_a$$

..... Infinite recursion

Level 4 recursion

Level 3 recursion

Level 2 recursion

Level 1 recursion

Player 1

Player 2

TRCoPO vs TRCoPO