# TWIN - DELAYED DDPG (TD3)

\# Bu algoritma modelsiz, gerçim-içi, politikadışı (off-policy) bir RL modelidir.
Bir TD3 ajanı, beklenen kümülatif uzun vadeli ödülü, en üst düzeye çıkaran
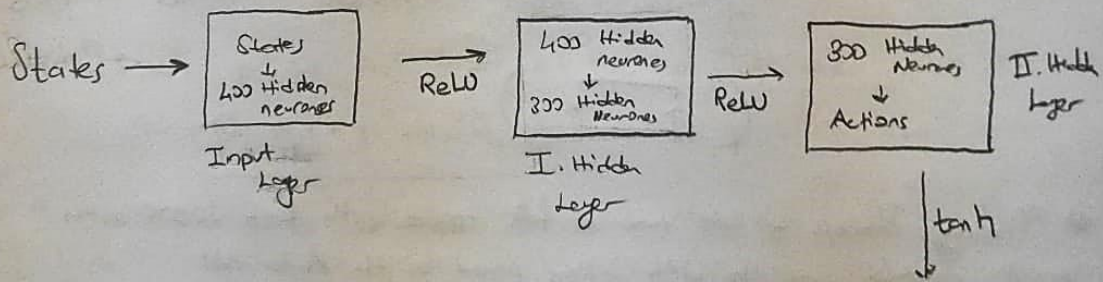optimal bir politikayı arayan actor-critic RL aracıdır.

## Initialization

**Step 1:** We initialize the Experience Reply memory, with a size of 20 000.
We will populated it with each new transition.

$$(S_1, S_1', a_1, r_1) \quad (S_2, S_2', a_2, r_2) \quad (S_3, S_3', a_3, r_3) \quad (S_4, S_4', a_4, r_4) \quad \cdots$$

first transition

— Experience Reply Memory —

**Step 2:** We build one neural network for the Actor model and one neural network for
the Actor target. (Yani 2 sinir ağı oluşturdık ve bunlar birbirinin aynısı)

States → [ States ↓ 400 Hidden neurones ] —ReLU→ [ 400 Hidden neurones ↓ 300 Hidden Neurones ] —ReLU→ [ 300 Hidden Neurones ↓ Actions ]

Input Layer          I. Hidden Layer          II. Hidden Layer

↓ tanh

Actions

\# Yukarıda ki sinir ağı modele yapanlar tarafında oluşturda bu model hem Actor model
hemde Actor target için yazılır.
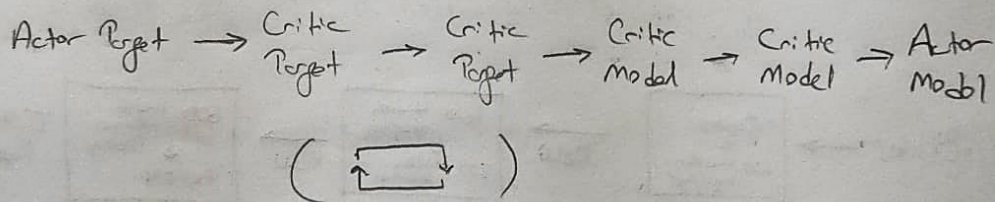
**Step 3 :** We build two neural networks for the two Critic models and two neural networks for the two Critic targets (4 adet kritik sinir ağ. iinn olacak.)

States → [States +Actions ↓ 400 Hidden neurons] ReLU → [400 Hidden Nrns ↓ 300 Hidden Neurons] ReLU → [300 Hidden Neurones ↓ Q-Value] → Q1 Value

Actions →

\# Yukarıda ki sinir ağı modeli üzerinde tachidan üretilen kritik modelidir. Bu Y.S.A. 4 adet olacaktır.

Ⅱ Elimizde toplam 6 adet Y.SA. oldu. Bunların ikisi Actor, dörde Critic'ler olmakta. Un-trendan iki Actor'ler politikaya öğrenecek, Critcler Q1 değerleri öğrenecek izleyeceğini yol su seçildi. (Bunlar kopu YSA)

Actor Target → Critic Target → Critic Target → Critic Model → Critic Model → Actor Model

( ⟲ )

\# Training Process – We run a full episode with first 10000 actions played randomly and then with actions played by the Actor model.

**Step 4 :** We sample a batch of transitions (s, s', a, r) from the memory. Then for each element of the batch:

\# Uygulama olarak 100 batch kullanacağız.

\# Bunlar hepsi ayrı batchler olacak. Yani her girişin batch ve 4 adet batch var.

| Last States | Next States | Actions | Rewards | |
|---|---|---|---|---|
| Last State 1 | Next State 1 | Action 1 | Reward 1 | I. transition |
| Last State 2 | Next State 2 | Action 2 | Reward 2 | II. transition |
| . . . | . . . | . . . | . . . | |
| Last State 100 | Next State 100 | Action 100 | Reward 100 | |

**Step 5:** From the next state $s'$, the Actor target plays the next action $a'$.

**Step 6:** We add Gaussian noise to the next action $a'$ and we clamp it in a range of values supported by the environment.
(Ajan farklı aksiyon keşfetmesi için garantiye alıyoruz.)

$$\tilde{a} \leftarrow \pi_{\phi'}(s') + \overset{\text{noise}}{\epsilon}, \quad \epsilon \sim clip(N(0, \tilde{\sigma}), -c, c)$$

$$\tilde{a} \leftarrow clip(\tilde{a})$$

**Step 7:** The two Critic targets take each the couple $(s', a')$ as input and return two Q-values $Q_{t_1}(s', a')$ and $Q_{t_2}(s', a')$ as outputs

**Step 8:** We keep the minimum of these two Q-values : $min(Q_{t_1}, Q_{t_2})$
It represents the approximated value of the next state.
(min alınır anlı aşırı tahminin önüne)

**Step 9:** We get the final target of the two Critic models, which is:

$$Q_t = r + \gamma * min(Q_{t_1}, Q_{t_2}) \quad \text{where } \gamma \text{ is the discount factor.}$$

\# Yukarıda formüle Q-func.'dan hatırla. Max değer geri gelen $Q_t$'leri incele.

**Step 10:** The two Critic models take each the couple $(s, a)$ as input and return two Q-values $Q_1(s, a)$ and $Q_2(s, a)$ as outputs.
(İki Critic model de iki farklı tehnik Q değeri çıkarıp ve bunu başlatıyoruz.)

**Step 11:** We compute the loss coming from the two Critic models:

$$Critic\ Loss = MSE\_Loss(Q_1(s,a), Q_t) + MSE\_Loss(Q_2(s,a), Q_t)$$

\# İki Critic model arasındaki kayıbı hesaplarız. 5

**Step 12 :** We backpropagate this Critic loss and update the parameters of the two Critic models with a SGD (Stochastic Gradient Descent) optimizer ("Adam" optimizers ile stokastik gradyen inis kullanede Critic kaybı azaltırız)

\# Training kısmının Q-learning parts bitti şimdi Policy Learning kısmına geçelim.

\* TD3'ün politika öğrenme bölümünün nihai amacı, her durumda (state) beklenen getiriyi en üst düzeye çıkaracak optimal aksiyonları gerçekleştirmek için Actor model'in (Ana politikası) optimal parametreleri bulmuş olduğu hatırlayın.

\* Q-value, bir kkhnın değeridir ve elbette bu var Q-value ne kadar artarsa, optimum beklenen değere o kadar yaklaşırız.

**Step 13 :** Once every two iterations, we update our Actor model by performing gradient ascent on the output of the first Critic Model:

$$\nabla_\phi J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s,a)\big|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s)$$

where $\phi$ and $\theta_1$ are resp. the weights of the Actor and the Critic.
( Her iki adımda, ilk Critic modelin çıktısında gradyen yükselişi ile Actor model'i güncelleriz.)

**Step 14 :** Still once every two iterations, we update the weights of the Actor target by Polyak averaging (gradient descent gibi bir optimizasyon yok.) :

$$\theta_i' \leftarrow \tau\theta_i + (1-\tau)\theta_i'$$
tek tek    tau

**Step 15 :** Still once every two iterations, we update the weights of the Critic target by Polyak averaging:

$$\phi' \leftarrow \tau\phi + (1-\tau)\phi'$$

-The End-
(O.B.)