Introduction to Machine Learning



Reinforcement Learning







Reinforcement Learning

Markov Decision Process

$$(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma, s_0)$$

- S = set of possible states
- A = set of possible actions
- $T: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$: transition model
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function
- γ = discount factor

Reinforcement Learning

Markov Decision Process

$$(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma, s_0)$$

- S = set of possible states
- A = set of possible actions
- $T: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$: transition model
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function
- γ = discount factor



Reinforcement Learning

Markov Decision Process

$$(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \gamma, s_0)$$

- S = set of possible states
- $\mathcal{A} = \text{set of possible actions}$
- $T: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$: transition model
- $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$: reward function
- γ = discount factor



Goal in RL: find a "policy" $\pi: S \to A$ that maximizes reward

Review: Value of a policy

• Given an MDP and a policy $\pi : S \to A$ we can find the value of a policy by solving a system of linear equations.

Review: Value of a policy

• Given an MDP and a policy $\pi : S \to A$ we can find the value of a policy by solving a system of linear equations.



Review: Value of a policy

• Given an MDP and a policy $\pi : S \to A$ we can find the value of a policy by solving a system of linear equations.



R(rich, plant)=100 R(poor, plant)=10 R(rich, fallow)=0 R(poor, fallow)=0

$$V_{\pi}^{0}(s) = 0; V_{\pi}^{h}(s) = R(s, \pi(s)) + \sum_{s'} T(s, \pi(s), s') \cdot V_{\pi}^{h-1}(s')$$
$$V_{\pi}(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{\pi}(s')$$

Can use to evaluate which policy is better. How to compute best policy?

- *h*: horizon (e.g. how many growing seasons left)
- $V^h_{\pi}(s)$: value (expected reward) with policy π starting at s



- *h*: horizon (e.g. how many planting seasons)
- Q^h(s, a): expected reward of starting at s, making action a, and then making the "best" action for the h-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$



- h: horizon (e.g. how many planting seasons)
- Q^h(s, a): expected reward of starting at s, making action a, and then making the "best" action for the h-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$



- *h*: horizon (e.g. how many planting seasons)
- $Q^h(s,a)$: expected reward of starting at *s*, making action *a*, and then making the "best" action for the *h*-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$ $Q^1(\text{rich, plant}) = 100; Q^1(\text{rich, fallow}) = 0; Q^1(\text{poor, plant}) = 10; Q^1(\text{poor, fallow}) = 0$



- *h*: horizon (e.g. how many planting seasons)
- $Q^h(s,a)$: expected reward of starting at *s*, making action *a*, and then making the "best" action for the *h*-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$ $Q^1(\text{rich, plant}) = 100; Q^1(\text{rich, fallow}) = 0; Q^1(\text{poor, plant}) = 10; Q^1(\text{poor, fallow}) = 0$ $Q^2(\text{rich, plant}) = R(\text{rich, plant}) + T(\text{rich, plant, rich}) \max_{a'} Q^1(\text{poor, }a')$ $+ T(\text{rich, plant, poor)} \max_{a'} Q^1(\text{poor, }a')$



- *h*: horizon (e.g. how many planting seasons)
- Q^h(s, a): expected reward of starting at s, making action a, and then making the "best" action for the h-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$ $Q^1(\text{rich, plant}) = 100; Q^1(\text{rich, fallow}) = 0; Q^1(\text{poor, plant}) = 10; Q^1(\text{poor, fallow}) = 0$ $Q^2(\text{rich, plant}) = 100 + (0.1)(100)$ + (0.9)(10) = 119



- *h*: horizon (e.g. how many planting seasons)
- Q^h(s, a): expected reward of starting at s, making action a, and then making the "best" action for the h-1 steps left
- With Q, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$ $Q^1(\text{rich, plant}) = 100; Q^1(\text{rich, fallow}) = 0; Q^1(\text{poor, plant}) = 10; Q^1(\text{poor, fallow}) = 0$ $Q^2(\text{rich, plant}) = 119; Q^2(\text{rich, fallow}) = 91; Q^2(\text{poor, plant}) = 29; Q^2(\text{poor, fallow}) = 91$



- *h*: horizon (e.g. how many planting seasons)
- $Q^h(s,a)$: expected reward of starting at *s*, making action *a*, and then making the "best" action for the *h*-1 steps left
- With *Q*, can find **an optimal policy**: $\pi_h^*(s) = \arg \max_a Q^h(s, a)$ $Q^0(s, a) = 0; Q^h(s, a) = R(s, a) + \sum_{s'} T(s, a, s') \max_{a'} Q^{h-1}(s', a')$ $Q^1(\text{rich, plant}) = 100; Q^1(\text{rich, fallow}) = 0; Q^1(\text{poor, plant}) = 10; Q^1(\text{poor, fallow}) = 0$ $Q^2(\text{rich, plant}) = 119; Q^2(\text{rich, fallow}) = 91; Q^2(\text{poor, plant}) = 29; Q^2(\text{poor, fallow}) = 91$ What's best? Any *s*, $\pi_1^*(s) = \text{plant}; \pi_2^*(\text{rich}) = \text{plant}, \pi_2^*(\text{poor}) = \text{fallow}$



- What if I don't stop farming? Is there any optimal policy?
- **Theorem**. There exists a (stationary) optimal policy π^* . I.e., for every policy π and for every state $s \in S$, $V_{\pi^*}(s) \ge V_{\pi}(s)$



• What if I don't stop farming? Is there any optimal policy?

- **Theorem**. There exists a (stationary) optimal policy π^* . I.e., for every policy π and for every state $s \in S$, $V_{\pi^*}(s) \ge V_{\pi}(s)$
- $Q^*(s, a)$: expected reward if we make best actions in future
 - If we knew $Q^*(s,a)$, then: $\pi^*(s) = \arg \max_a Q^*(s,a)$
- Note: $Q^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q^*(s', a')$
 - Not linear in $Q^*(s,a)$, so not as easy to solve as $V_{\pi}(s)$

Finite-horizon value iteration:

$$Q^{0}(s,a) = 0 \quad Q^{1}(s,a) = R(s,a)$$
$$Q^{h}(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q^{h-1}(s',a')$$

Finite-horizon value iteration:

$$\begin{split} Q^0(s,a) &= 0 \ Q^1(s,a) = R(s,a) \\ Q^h(s,a) &= R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q^{h-1}(s',a') \\ \texttt{Infinite-Horizon-Value-Iteration} (\mathcal{S},\mathcal{A},T,R,\gamma,\epsilon) \\ \texttt{for} \text{ each state } s \in \mathcal{S} \text{ and each action } a \in \mathcal{A} \\ \texttt{Initialize } Q_{\text{old}}(s,a) &= 0 \end{split}$$

Finite-horizon value iteration:

 $\begin{array}{l} Q^0(s,a)=0 \ Q^1(s,a)=R(s,a) \\ Q^h(s,a)=R(s,a)+\gamma\sum_{s'}T(s,a,s')\max_{a'}Q^{h-1}(s',a') \end{array}$ Infinite-Horizon-Value-Iteration ($\mathcal{S},\mathcal{A},T,R,\gamma,\epsilon$) for each state $s\in\mathcal{S}$ and each action $a\in\mathcal{A}$ Initialize $Q_{\mathrm{old}}(s,a)=0$ while True

for each state
$$s \in S$$
 and each action $a \in A$
 $Q_{\text{new}}(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

Finite-horizon value iteration:

 $\begin{array}{l} Q^0(s,a)=0 \ Q^1(s,a)=R(s,a) \\ Q^h(s,a)=R(s,a)+\gamma\sum_{s'}T(s,a,s')\max_{a'}Q^{h-1}(s',a') \end{array}$ Infinite-Horizon-Value-Iteration ($\mathcal{S},\mathcal{A},T,R,\gamma,\epsilon$)
for each state $s\in\mathcal{S}$ and each action $a\in\mathcal{A}$ Initialize $Q_{\mathrm{old}}(s,a)=0$ while True
for each state $s\in\mathcal{S}$ and each patients $a\in\mathcal{A}$

for each state
$$s \in S$$
 and each action $a \in A$
 $Q_{\text{new}}(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$
if $\max_{s,a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$
return Q_{new}
 $Q_{\text{old}} = Q_{\text{new}}$

Finite-horizon value iteration:

 $Q^{0}(s,a) = 0 \ Q^{1}(s,a) = R(s,a)$ $Q^{h}(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q^{h-1}(s',a')$ Infinite-Horizon-Value-Iteration ($\mathcal{S}, \mathcal{A}, T, R, \gamma, \epsilon$) **for** each state $s \in \mathcal{S}$ and each action $a \in \mathcal{A}$ Initialize $Q_{\text{old}}(s, a) = 0$ while True **for** each state $s \in \mathcal{S}$ and each action $a \in \mathcal{A}$ $Q_{\text{new}}(s,a) = R(s,a) + \gamma \sum_{s'} T(s,a,s') \max_{a'} Q_{\text{old}}(s',a')$ if $\max_{s,a} |Q_{\text{old}}(s,a) - Q_{\text{new}}(s,a)| < \epsilon$ Issue: need to know reward and return Q_{new} transition functions! $Q_{\text{old}} = Q_{\text{new}}$

A more realistic scenario







Reinforcement Learning Overview

Goal in RL: find a "policy" $\pi: S \to A$ that maximizes reward in an **unknown** MDP.

Reinforcement Learning Overview

Goal in RL: find a "policy" $\pi: S \to A$ that maximizes reward in an unknown MDP.



https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html

Reinforcement Learning Overview

Goal in RL: find a "policy" $\pi: S \to A$ that maximizes reward in an **unknown** MDP.



https://spinningup.openai.com/en/latest/spinningup/rl_intro2.html

```
\begin{aligned} Q\text{-LEARNING}(\mathcal{S},\mathcal{A},s_0,\gamma,\alpha) \\ & \text{for } s\in\mathcal{S}, a\in\mathcal{A}: \\ & Q[s,a]=0 \end{aligned}
```

Initialize Q function to 0

```
Q-LEARNING(\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha)
    for s \in S, a \in A:
          Q[s, a] = 0
                                                      Initialize Q function to 0
    s = s_0 \parallel Or draw an s randomly from S
    while True:
                                                     "Act" in the environment
          a = select_action(s, Q)
          r, s' = execute(a)
                                              Receive reward, transition to s'
```

Intuition: sampling to estimate the transition model $T: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$

```
Q-LEARNING(\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha)
    for s \in S, a \in A:
          Q[s, a] = 0
                                                       Initialize Q function to 0
    s = s_0 \parallel Or draw an s randomly from S
    while True:
          a = select_action(s, Q)
                                                      "Act" in the environment
          r, s' = execute(a)
                                               Receive reward, transition to s'
          Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])
```

Update Q function for the sampled state and action

```
Q-LEARNING(\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha)
    for s \in S, a \in A:
           Q[s, a] = 0
                                                      Initialize Q function to 0
     s = s_0 \parallel Or draw an s randomly from S
     while True:
           a = select_action(s, Q)
                                                      "Act" in the environment
          r, s' = execute(a)
                                              Receive reward, transition to s'
           Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])
          s = s'
                                              Update Q function for the
Move onto next state
                                              sampled state and action
```

```
Q-LEARNING(S, A, s_0, \gamma, \alpha)

for s \in S, a \in A:

Q[s, a] = 0

s = s_0 // Or draw an s randomly from S

while True:

a = select\_action(s, Q)

r, s' = execute(a)

Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])
```

while True

s = s'

for each state $s \in S$ and each action $a \in A$ $Q_{\text{new}}(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

```
Q-\text{LEARNING}(S, \mathcal{A}, s_0, \gamma, \alpha)

for s \in S, a \in \mathcal{A}:

Q[s, a] = 0

s = s_0 \parallel \text{Or draw an s randomly from S}

while True:

a = \text{select}_action(s, Q)

r, s' = \text{execute}(a)

Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])

s = s'
```

while True

for each state $s \in S$ and each action $a \in A$ $Q_{\text{new}}(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

Since we don't know R and T, we estimate it by "sampling" in the environment

```
Q-LEARNING(S, A, s_0, \gamma, \alpha)

for s \in S, a \in A:

Q[s, a] = 0

s = s_0 // \text{ Or draw an s randomly from S}

while True:
```

$$a = select_action(s, Q)$$

$$r, s' = execute(a)$$

$$Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])$$

$$s = s'$$

Moving average

Q-LEARNING(S, A, s_0 , γ , α) for $s \in S$, $a \in A$: Q[s, a] = 0 $s = s_0 // \text{ Or draw an s randomly from S}$ while True:

$$a = select_action(s, Q)$$

$$r, s' = execute(a)$$

$$Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a'])$$

$$s = s'$$

Moving average

$$Q[s, a] = Q[s, a] - \alpha \left(Q[s, a] - (r + \gamma \max_{a'} Q[s', a']) \right)$$

Looks like gradient descent!

$$Q[s,a] = Q[s,a] - \alpha \left(Q[s,a] - (r + \gamma \max_{a'} Q[s',a']) \right)$$

Linear regression with just the bias term

$$y = \theta_0$$

$$L(\theta_0) = \frac{1}{2}(y - \theta_0)^2$$

$$\theta_0 = \theta_0 - \eta(\theta_0 - y)$$

$$Q[s,a] = Q[s,a] - \alpha \left(Q[s,a] - (r + \gamma \max_{a'} Q[s',a']) \right)$$

Linear regression with just the bias term

$$y = \theta_0$$
$$L(\theta_0) = \frac{1}{2}(y - \theta_0)^2$$
$$\theta_0 = \theta_0 - \eta(\theta_0 - y)$$

Current guess Target

$$Q[s,a] = Q[s,a] - \alpha \left(Q[s,a] - (r + \gamma \max_{a'} Q[s',a']) \right)$$

Linear regression with just the bias term

$$y = \theta_0$$
$$L(\theta_0) = \frac{1}{2}(y - \theta_0)^2$$
$$\theta_0 = \theta_0 - \eta(\theta_0 - y)$$

Learning rate Current guess Target

$$Q[s,a] = Q[s,a] - \alpha \left(Q[s,a] - (r + \gamma \max_{a'} Q[s',a']) \right)$$

Linear regression with just the bias term



Key difference: in Q-learning, the target itself is a function of what is being learned (and the reward)!

 $\begin{array}{l} Q\text{-LEARNING}(\$, \mathcal{A}, s_0, \gamma, \alpha) \\ \textbf{for } s \in \$, a \in \mathcal{A} : \\ Q[s, a] = 0 \\ s = s_0 \ \textit{//} \ Or \ draw \ an \ s \ randomly \ from \ \$ \\ \textbf{while True:} \\ a = \underbrace{\text{select_action}(s, Q)}_{r, s' = execute(a)} \\ Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a']) \\ s = s' \end{array}$

```
\begin{array}{l} Q\text{-LEARNING}(\$, \mathscr{A}, s_0, \gamma, \alpha) \\ \textbf{for } s \in \$, \alpha \in \mathscr{A} : \\ Q[s, \alpha] = 0 \\ s = s_0 \ /\!\!/ \ \text{Or draw an } s \text{ randomly from } \$ \\ \textbf{while True:} \\ a = \underbrace{\text{select\_action}(s, Q)}_{r, s' = execute(a)} \\ Q[s, \alpha] = (1 - \alpha)Q[s, \alpha] + \alpha(r + \gamma \max_{\alpha'} Q[s', \alpha']) \\ s = s' \end{array}
```

Argmax action from current policy?

$$arg \, max_{\mathfrak{a} \in \mathcal{A}} \, Q(s, \mathfrak{a})$$

Initially Q function is bad \Rightarrow argmax not a good idea.

 $\begin{array}{l} \text{Q-LEARNING}(\$, \mathcal{A}, s_0, \gamma, \alpha) \\ \text{for } s \in \$, a \in \mathcal{A} : \\ Q[s, a] = 0 \\ s = s_0 \ \text{// Or draw an } s \text{ randomly from } \$ \\ \text{while True:} \\ a = \begin{array}{l} \text{select_action(s, Q)} \\ r, s' = \text{execute}(a) \\ Q[s, a] = (1 - \alpha)Q[s, a] + \alpha(r + \gamma \max_{a'} Q[s', a']) \\ s = s' \end{array}$

ε-greedy strategy:

- with probability 1ϵ , choose $\arg \max_{a \in A} Q(s, a)$
- with probability ε , choose the action $a \in A$ uniformly at random

Exploitation

Exploration

- Q[s,a] is a scalar for each possible state and action. ("Tabular" Q-learning)
- What if states are high dimensional or continuous?



• What if actions are not discrete?



• Parameterize Q function to be an output from a neural network!









Q-LEARNING($\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha$)

- 1 **for** $s \in S$, $a \in A$:
- $2 \qquad Q[s,a] = 0$
- 3 $s = s_0 \parallel Or draw an s randomly from S$
- 4 while True:
- 5 $a = select_action(s, Q)$
- 6 r, s' = execute(a)
- 7 $Q[s,a] = (1-\alpha)Q[s,a] + \alpha(r + \gamma \max_{a'} Q[s',a'])$
- 8 s = s'

Q-LEARNING($\mathcal{S}, \mathcal{A}, s_0, \gamma, \alpha$)

- $s = s_0 \parallel Or draw$ an s randomly from S
- 4 while True:
- $a = select_action(s, Q)$
- r, s' = execute(a)
- $Q[s,a] = (1-\alpha)Q[s,a] + \alpha(r + \gamma \max_{a'} Q[s',a'])$
- s = s'

Q-LEARNING($(S, A, s_0, \gamma, \alpha)$) for $s \in S$, $a \in A$: Randomly initialize $Q_{\theta}(s, a)$ Q[s, a] = 02 3 $s = s_0 \parallel Or draw an s randomly from S$ while True: 4 5 $a = select_action(s, Q)$ 6 r, s' = execute(a) $\left| \mathbf{Q}[\mathbf{s},\mathbf{a}] = (1-\alpha)\mathbf{Q}[\mathbf{s},\mathbf{a}] + \alpha(\mathbf{r}+\gamma \max_{\mathbf{a}'}\mathbf{Q}[\mathbf{s}',\mathbf{a}']) \right|$ 7 8 $\overline{s = s'}$ Minimize regression loss $(Q_{\theta}(s,a) - (r + \gamma \max_{a'} Q_{\theta}(s',a')))^2$

Q-LEARNING($(S, A, s_0, \gamma, \alpha)$) for $s \in S$, $a \in A$: Randomly initialize $Q_{\theta}(s, a)$ Q[s, a] = 02 3 $s = s_0 \parallel Or draw an s randomly from S$ while True: 4 5 a = select action(s, Q)6 r, s' = execute(a) $\left| \mathbf{Q}[\mathbf{s},\mathbf{a}] = (1-\alpha)\mathbf{Q}[\mathbf{s},\mathbf{a}] + \alpha(\mathbf{r}+\gamma \max_{\mathbf{a}'}\mathbf{Q}[\mathbf{s}',\mathbf{a}']) \right|$ 7 8 $\overline{s = s'}$ Minimize regression loss Issue: instability arising from both guess and $(Q_{\theta}(s,a) - (r + \gamma \max_{a'} Q_{\theta}(s',a')))^2$ target being from a learned network.

FITTED-Q-LEARNING($\mathcal{A}, s_0, \gamma, \alpha, \varepsilon, m$)

$$s = s_0 \parallel Or draw an s randomly from S$$

 $\mathcal{D} = \{ \}$

initialize neural-network representation of Q

while True:

$$\begin{split} & \mathcal{D}_{new} = \text{experience from executing ε-greedy policy based on Q for m steps $$ $\mathcal{D} = \mathcal{D} \cup \mathcal{D}_{new}$ represented as (s, a, r, s') tuples $$ $D_{sup} = \{(x^{(i)}, y^{(i)})\}$ where $x^{(i)} = (s, a)$ and $y^{(i)} = r + \gamma \max_{a' \in \mathcal{A}} Q(s', a')$ $ for each tuple $(s, a, r, s')^{(i)} \in \mathcal{D}$ $ re-initialize neural-network representation of Q $$ $Q = supervised_NN_regression(D_{sup})$ $$ $$$

FITTED-Q-LEARNING($\mathcal{A}, s_0, \gamma, \alpha, \varepsilon, \mathfrak{m}$)

$$s = s_0 \parallel Or draw an s randomly from S $\mathcal{D} = \{ \}$$$

initialize neural-network representation of Q **While** True: Collect data from current policy

$$\begin{split} & \mathcal{D}_{new} = \text{experience from executing ε-greedy policy based on Q for m steps $$ $\mathcal{D} = \mathcal{D} \cup \mathcal{D}_{new}$ represented as (s, a, r, s') tuples $$ $D_{sup} = \{(x^{(i)}, y^{(i)})\}$ where $x^{(i)} = (s, a)$ and $y^{(i)} = r + \gamma \max_{a' \in \mathcal{A}} Q(s', a')$ $$ for each tuple $(s, a, r, s')^{(i)} \in \mathcal{D}$ $$ re-initialize neural-network representation of Q $$ $Q = supervised_NN_regression(D_{sup})$ $$ $$$$

FITTED-Q-LEARNING($\mathcal{A}, s_0, \gamma, \alpha, \varepsilon, m$)

$$s = s_0 \parallel Or draw an s randomly from S$$

 $\mathcal{D} = \{ \}$

initialize neural-network representation of Q **While** True: Collect data from current policy

$$\begin{split} & \mathcal{D}_{new} = experience \ from \ executing \ \varepsilon-greedy \ policy \ based \ on \ Q \ for \ m \ steps \\ & \mathcal{D} = \mathcal{D} \cup \mathcal{D}_{new} \ represented \ as \ (s, a, r, s') \ tuples \\ & D_{sup} = \{(x^{(i)}, y^{(i)})\} \ where \ x^{(i)} = (s, a) \ and \ y^{(i)} = r + \gamma \ max_{a' \in \mathcal{A}} \ Q(s', a') \\ & \text{for each tuple } (s, a, r, s')^{(i)} \in \mathcal{D} \\ & \text{re-initialize neural-network representation of } Q \\ & Q = \text{supervised}_NN_regression(D_{sup}) \end{aligned}$$

FITTED-Q-LEARNING($\mathcal{A}, s_0, \gamma, \alpha, \varepsilon, \mathfrak{m}$)

$$s = s_0 \parallel Or draw an s randomly from S$$

 $\mathcal{D} = \{ \}$

initialize neural-network representation of Q **While** True: Collect data from current policy

$$\begin{split} & \mathcal{D}_{new} = experience \ from \ executing \ \varepsilon-greedy \ policy \ based \ on \ Q \ for \ m \ steps \\ & \mathcal{D} = \mathcal{D} \cup \mathcal{D}_{new} \ represented \ as \ (s, a, r, s') \ tuples \\ & D_{sup} = \{(x^{(i)}, y^{(i)})\} \ where \ x^{(i)} = (s, a) \ and \ y^{(i)} = r + \gamma \ max_{a' \in \mathcal{A}} \ Q(s', a') \\ & \text{for each tuple } (s, a, r, s')^{(i)} \in \mathcal{D} \\ & \text{re-initialize neural-network representation of } Q \\ & Q = supervised_NN_regression(D_{sup}) \\ \end{split}$$

Train a neural network with regression loss!