Model Predictive Control and Reinforcement Learning - On-Policy Control with Function Approximation -

Joschka Boedecker and Moritz Diehl

University Freiburg

October 6, 2022



#### Lecture Overview



- 1 Function Approximation in Reinforcement Learning
- 2 Linear Methods
- **3** On-policy Control with Function Approximation
- 4 Off-policy Learning
- 5 Problems of Off-policy Learning with Function Approximation
- 6 Deep Q-learning
- 7 DDPG
- 8 TD3



Slide contents are partially based on *Reinforcement Learning: An Introduction* by Sutton and Barto and the Reinforcement Learning lecture by David Silver.

- Up to this point, we represented all elements of our RL systems by tables (value functions, models and policies)
- ▶ If the state and action spaces are very large or infinite, this is not a feasible solution
- We can apply function approximation to find a more compact representation of RL components and to generalize over states and actions
- Reinforcement Learning with function approximation comes with new issues that do not arise in Supervised Learning – such as non-stationarity, bootstrapping and delayed targets

• Here: we estimate value-functions  $v_{\pi}(\cdot)$  and  $q_{\pi}(\cdot, \cdot)$  by function approximators  $\hat{v}(\cdot, \mathbf{w})$  and  $\hat{q}(\cdot, \cdot, \mathbf{w})$ , parameterized by weights  $\mathbf{w}$ 



But we can also represent models or policies

We can use different types of function approximators:

- Linear combinations of features
- Neural networks
- Decision trees
- Gaussian processes
- Nearest neighbor methods
- ▶ ...

Here: We focus on differentiable FAs and update the weights via gradient descent.



We want to update our weights w.r.t. the Mean Squared Value Error of our prediction:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{1}{2}\alpha \nabla [v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2$$
$$= \mathbf{w}_t + \alpha [v_{\pi}(S_t) - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t)$$

However, we don't have  $v_{\pi}(S_t)$ .



#### Gradient MC

 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [\mathbf{G}_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$ 

#### Semi-gradient TD(0)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha[R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$$

Why are bootstrapping methods, defined this way, called semi-gradient methods?



 $\mathbf{w} \leftarrow \mathbf{w} + \alpha [\mathbf{G}_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$ 

#### Semi-gradient TD(0)

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$$

Why are bootstrapping methods, defined this way, called *semi-gradient methods*? They take into account the effects of changing  $\mathbf{w}$  w.r.t. the prediction, but not w.r.t. the target!

#### Linear Methods



- ▶ Represent state s by feature vector  $\mathbf{x}(s) = (x^1(s), x^2(s), \dots, x^d(s))^\top$
- These features can also be non-linear functions/combinations of state dimensions
- Linear methods approximate the value function by a linear combination of these features

$$\hat{v}(s, \mathbf{w}) = \mathbf{w}^{\top} \mathbf{x}(s) = \sum_{i=1}^{d} w^{i} x^{i}(s)$$

- ▶ Therefore,  $\nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w}) = \mathbf{x}(s)$
- Gradient MC prediction converges under linear FA
- On-policy linear semi-gradient TD(0) is stable
- Unfortunately, this does not hold for non-linear FA

# Fixed point of on-policy linear semi-gradient TD



$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \left( R_{t+1} + \gamma \mathbf{w}_t^\top \mathbf{x}_{t+1} - \mathbf{w}_t^\top \mathbf{x}_t \right) \mathbf{x}_t$$
$$= \mathbf{w}_t + \alpha \left( R_{t+1} \mathbf{x}_t - \mathbf{x}_t (\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top \mathbf{w}_t \right)$$

The expected next weight vector can thus be written:

$$\mathbb{E}[\mathbf{w}_{t+1}|\mathbf{w}_t] = \mathbf{w}_t + \alpha(\mathbf{b} - \mathbf{A}\mathbf{w}_t),$$

where  $\mathbf{b} = \mathbb{E}[R_{t+1}\mathbf{x}_t]$  and  $\mathbf{A} = \mathbb{E}[\mathbf{x}_t(\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]$ 

▶ If the system converges, it has to converge to the *fixed point*:

$$\mathbf{w}_{\mathsf{TD}} = \mathbf{A}^{-1}\mathbf{b}$$

## **Coarse Coding**



Divide the state space in circles/tiles/shapes and check in which some state is inside. This is a binary representation of the location of a state and leads to generalization.







Broad generalization Asymmetric generalization







Log stripes

Diagonal stripes



# On-policy Control with Function Approximation



- ▶ Again, up to this point we discussed Policy Evaluation based on state value functions
- ▶ In order to apply FA in control, we parameterize the action-value function

#### Semi-gradient SARSA

 $\mathbf{w} \leftarrow \mathbf{w} + \alpha[R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \mathbf{w}) - \hat{q}(S_t, A_t, \mathbf{w})] \nabla \hat{q}(S_t, A_t, \mathbf{w})$ 



#### Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q}: \mathbb{S} \times \mathcal{A} \times \mathbb{R}^d \to \mathbb{R}$ Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$ Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ ) Loop for each episode:  $S, A \leftarrow \text{initial state and action of episode (e.g., <math>\varepsilon$ -greedy) Loop for each step of episode: Take action A, observe R, S'If S' is terminal:  $\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$ Go to next episode Choose A' as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedv)  $\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$  $S \leftarrow S'$  $A \leftarrow A'$ 

# Off-policy Learning



- We want to learn the optimal policy, but we have to account for the problem of maintaining exploration
- We call the (optimal) policy to be learned the *target policy*  $\pi$  and the policy used to generate behaviour the *behaviour policy* b
- We say that learning is from data off the target policy thus, those methods are referred to as off-policy learning

## Importance Sampling



- Weight returns according to the relative probability of target and behaviour policy
- ▶ Define state-transition probabilities p(s'|s, a) as  $p(s'|s, a) = \Pr{S_t = s'|S_{t-1} = s, A_{t-1} = a} = \sum_{r \in \mathcal{R}} p(s', r|s, a)$
- The probability of the subsequent trajectory under any policy  $\pi$ , starting in  $S_t$ , then is:

$$\Pr\{A_t, S_{t+1}, A_{t+1}, \dots, S_T | S_t, A_{t:T-1} \sim \pi\}$$
  
=  $\pi(A_t | S_t) p(S_{t+1} | S_t, A_t) \pi(A_{t+1} | S_{t+1}) \cdots p(S_T | S_{T-1}, A_{T-1})$   
=  $\prod_{k=t}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$ 

## Importance Sampling



The relative probability therefore is:

Definition: Importance Sampling Ratio

$$\rho_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k|S_k) p(S_{k+1}|S_k, A_k)} = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k)}{\prod_{k=t}^{T-1} b(A_k|S_k)}$$

The expectation of the returns by b is  $\mathbb{E}[G_t|S_t = s] = v_b(s)$ . However, we want to estimate the expectation under  $\pi$ . Given the importance sampling ratio, we can transform the MC returns by b to yield the expectation under  $\pi$ :

$$\mathbb{E}[\rho_{t:T-1}G_t|S_t=s] = v_{\pi}(s).$$

Importance Sampling can come with a vast increase in variance.



To use importance sampling with function approximation, replace the update to an array to an update to weight vector  $\mathbf{w}$ , and correct it with the importance sampling weight.

#### Off-policy MC Prediction

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \rho_{t:T-1}[G_t - \hat{v}(S_t, \mathbf{w})] \nabla \hat{v}(S_t, \mathbf{w})$$

#### Semi-gradient Off-policy TD(0)

$$\begin{split} \mathbf{w} &\leftarrow \mathbf{w} + \alpha \rho_t \delta_t \nabla \hat{v}(S_t, \mathbf{w}) \\ \text{where } \delta_t &= R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \end{split}$$

# Baird's Counterexample





The reward is 0 for all transitions, hence  $v_{\pi}(s) = 0$ . This could be exactly approximated by  $\mathbf{w} = \mathbf{0}$ .

MPC and RL - Lecture 9

#### J. Boedecker and M. Diehl, University Freiburg

# Baird's Counterexample



#### Semi-gradient DP

$$\mathbf{w} \leftarrow \mathbf{w} + \frac{\alpha}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (\mathbb{E}[R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) | S_t = s] - \hat{v}(s, \mathbf{w})) \nabla \hat{v}(s, \mathbf{w})$$



### The Deadly Triad



The combination of

- Function Approximation,
- Bootstrapping and
- Off-policy Learning

is known as the Deadly Triad, since it can lead to stability issues and divergence.

# Neural Fitted-Q Iteration (NFQ) [Riedmiller 2005]



- Model-free off-policy RL algorithm that works on continuous state and discrete action spaces
- Q-function is represented by a multi-layer perceptron
- ► One of the first approaches that combined RL with ANNs, predecessor of DQN

# Neural Fitted-Q Iteration (NFQ) [Riedmiller 2005]

```
for iteration i = 1, ..., N do
       sample trajectory with \epsilon-greedy exploration and add to memory D
       initialize network weights randomly
       generate pattern set P = \{(x_i, y_i) | j = 1..|D|\} with
      x_j = (s_j, a_j) \text{ and } y_j = \begin{cases} r_j & \text{if } s_j \text{ is terminal} \\ r_j + \gamma \max_{a'} Q(s_{j+1}, a', \mathbf{w}_i) & \text{else} \end{cases}
       for iteration k = 1, ..., K do
              Fit weights according to:
 L(\mathbf{w}_{i}) = \frac{1}{|D|} \sum_{j=1}^{|D|} (y_{j} - Q(x_{j}, \mathbf{w}_{i}))^{2}
       end
end
```

Algorithm 1: NFQ



DQN provides a stable solution to deep RL:

- Use experience replay (as in NFQ)
- Sample minibatches (as opposed to Full Batch in NFQ)
- Freeze target Q-networks (no target networks in NFQ)
- > Optional: Clip rewards or normalize network adaptively to sensible range

To remove correlations, build data set from agent's own experience

- Take action  $a_t$  according to  $\epsilon$ -greedy policy
- Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  in replay memory D
- **>** Sample random mini-batch of transitions (s, a, r, s') from D
- Optimize MSE between Q-network and Q-learning targets, e.g.

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

To avoid oscillations, fix parameters used in Q-learning target

 $\blacktriangleright$  Compute Q-learning targets w.r.t. old, fixed parameters  $\mathbf{w}^-$ 

$$r + \gamma \operatorname*{arg\,max}_{a'} Q(s', a', \mathbf{w}^-)$$

Optimize MSE between Q-network and Q-learning targets

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

- $\blacktriangleright$  Periodically update fixed parameters  $\mathbf{w}^- \leftarrow \mathbf{w}$ 
  - $\blacktriangleright$  hard update: update target network every N steps
  - slow update: slowly update weights of target network every step by

$$\mathbf{w}^- \leftarrow (1-\tau)\mathbf{w}^- + \tau \mathbf{w}$$

# Deep Q-Networks (DQN)

```
Initialize replay memory D to capacity N
Initialize action-value function Q with random weights
for episode i = 1, ..., M do
       for t = 1, ..., T do
                select action a_t \epsilon-greedily
                Store transition (s_t, a_t, s_{t+1}, r_t) in D
               Sample minibatch of transitions (s_j, a_j, r_j, s_{j+1}) from D
               Set y_j = \begin{cases} r_j & \text{if } s_{j+1} \text{ is terminal} \\ r_j + \gamma \max_{a'} Q(s_{j+1}, a', \mathbf{w}^-) & \text{else} \end{cases}
                Update the parameters of Q according to:
                      \nabla \mathbf{w}_i L_i(\mathbf{w}_i) = \mathbb{E}_{s,a,s,r \sim D}[(r + \gamma \max_{i} Q(s', a', \mathbf{w}_i) - Q(s, a, \mathbf{w}_i)) \nabla_{\mathbf{w}_i} Q(s, a, \mathbf{w}_i)]
                  Update target network
        end
end
```



## Deep Q-Networks: Reinforcement Learning in Atari



# Deep Q-Networks: Reinforcement Learning in Atari

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





				DQN
	Q-Learning	Q-Learning	Q-Learning	Q-learning
			+ Replay	+ Replay
		+ Target Q		+ Target Q
Breakout	3	10	241	317
Enduro	29	142	831	1006
River Raid	1453	2868	4103	7447
Seaquest	276	1003	831	2894
Space Invaders	302	373	826	1089



- DDPG is an actor-critic method (Continuous DQN)
- ▶ Recall the DQN-target:  $y_j = r_j + \gamma \max_a Q(s_{j+1}, a, \mathbf{w}^-)$
- In case of continuous actions, the maximization step is not trivial
- Therefore, we approximate deterministic actor μ representing the arg max<sub>a</sub> Q(s<sub>j+1</sub>, a, w) by a neural network and update its parameters following the *Deterministic Policy Gradient Theorem*:

$$\nabla_{\theta} J \approx \frac{1}{N} \sum_{j} \nabla_{a} Q(s_{j}, a, \mathbf{w})|_{a = \mu(s_{j})} \nabla_{\theta} \mu(s_{j}, \boldsymbol{\theta})$$

• Exploration by adding Gaussian noise to the output of  $\mu$ 



► The Q-function is fitted to the adapted TD-target:

$$y_j = r_j + \gamma Q(s_{j+1}, \mu(s_{j+1}, \boldsymbol{\theta}^-), \mathbf{w}^-)$$

▶ The parameters of target networks  $\mu(\cdot, \theta^-)$  and  $Q(\cdot, \cdot, \mathbf{w}^-)$  are then adjusted with a soft update

$$\mathbf{w}^- \leftarrow (1-\tau)\mathbf{w}^- + \tau \mathbf{w} \text{ and } \boldsymbol{\theta}^- \leftarrow (1-\tau)\boldsymbol{\theta}^- + \tau \boldsymbol{\theta}$$

with  $\tau \in [0,1]$ 

- DDPG is very popular and builds the basis for more SOTA actor-critic algorithms
- However, it can be quite unstable and sensitive to its hyperparameters

```
Initialize replay memory D to capacity N
Initialize critic Q and actor \mu with random weights
for episode i = 1, ..., M do
        for t = 1, ..., T do
                select action a_t = \mu(s_t, \theta) + \epsilon, where \epsilon \sim \mathcal{N}(0, \sigma)
                Store transition (s_t, a_t, s_{t+1}, r_t) in D
                Sample minibatch of transitions (s_j, a_j, r_j, s_{j+1}) from D
               Set y_j = \begin{cases} r_j & \text{if } s_{j+1} \text{ is terminal} \\ r_j + \gamma \ Q(s_{j+1}, \mu(s_{j+1}, \boldsymbol{\theta}^-), \mathbf{w}^-) & \text{else} \end{cases}
                Update the parameters of Q according to the TD-error
                Update the parameters of \mu according to:
                                              \nabla_{\theta} J \approx \frac{1}{N} \sum_{i} \nabla_{a} Q(s_{i}, a, \mathbf{w})|_{a = \mu(s_{j})} \nabla_{\theta} \mu(s_{j}, \boldsymbol{\theta})
                Adjust the parameters of the target networks via a soft update
        end
end
```



## **Overestimation** Bias

- In all control algorithms so far, the target policy is created by the maximization of a value-function
- We thus consider the maximum over estimated values as an estimate of the maximum value
- This can lead to the so-called overestimation bias



## Overestimation Bias



▶ Recall the Q-learning target:  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$ 

▶ Imagine two random variables X<sub>1</sub> and X<sub>2</sub>:

 $\mathbb{E}[\max(X_1, X_2)] \ge \max(\mathbb{E}[X_1], \mathbb{E}[X_2])$ 

▶  $Q(S_{t+1}, a)$  is not perfect – it can be *noisy*:

$$\max_{a} Q(S_{t+1}, a) = \overbrace{Q(S_{t+1}, \operatorname*{arg\,max}_{a} Q(S_{t+1}, a))}^{\text{value comes from } Q}$$

If the noise in these is decorrelated, the problem goes away!

## Double Q-learning



#### Double Q-learning, for estimating $Q_1 \approx Q_2 \approx q_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$ Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ , for all  $s \in S^+$ ,  $a \in \mathcal{A}(s)$ , such that  $Q(terminal, \cdot) = 0$ Loop for each episode: Initialize SLoop for each step of episode: Choose A from S using the policy  $\varepsilon$ -greedy in  $Q_1 + Q_2$ Take action A, observe R, S'With 0.5 probabilility:  $Q_1(S,A) \leftarrow Q_1(S,A) + \alpha \Big( R + \gamma Q_2 \big( S', \operatorname{arg\,max}_a Q_1(S',a) \big) - Q_1(S,A) \Big)$ else:  $Q_2(S,A) \leftarrow Q_2(S,A) + \alpha \Big( R + \gamma Q_1 \big( S', \operatorname{arg\,max}_a Q_2(S',a) \big) - Q_2(S,A) \Big)$  $S \leftarrow S'$ until S is terminal

#### Double Q-learning







TD3 adds three adjustments to vanilla DDPG

- Clipped Double Q-Learning
- Delayed Policy Updates
- Target-policy smoothing

# TD3: Clipped Double Q-Learning



- In order to alleviate the overestimation bias (which is also present in actor-critic methods), TD3 learns two approximations of the action-value function
- It the takes the minimum of both predictions as the second part of the TD-target:

$$y_j = r_j + \gamma \min_{i \in \{1,2\}} Q(s_{j+1}, \mu(s_{j+1}), \mathbf{w}_i^-)$$



## TD3: Delayed Policy Updates



- Due to the mutual dependency between actor and critic updates...
  - values can diverge when the policy leads to overestimation and
  - the policy will lead to bad regions of the state-action space when the value estimates lack in (relative) accuracy
- Therefore, policy updates on states where the value-function has a high prediction error can cause divergent behaviour
- We already know how to compensate for that: target networks
- $\blacktriangleright$  Freeze target and policy networks between d updates of the value function
- ► This is called a *Delayed Policy Update*



- ► Target-policy Smoothing adds Gaussian noise to the next action in target calculation
- It transforms the Q-update towards an Expected SARSA update fitting the value of a small area around the target-action:

$$y_j = r_j + \gamma \min_{i \in \{1,2\}} Q(s_{j+1}, \mu(s_{j+1}) + \mathsf{clip}(\epsilon, -c, c), \mathbf{w}_i^-),$$

where  $\epsilon \sim \mathcal{N}(0,\sigma)$ 

#### TD3: Ablation



Table 2. Average return over the last 10 evaluations over 10 trials of 1 million time steps, comparing ablation over delayed policy updates (DP), target policy smoothing (TPS), Clipped Double Q-learning (CDQ) and our architecture, hyper-parameters and exploration (AHE). Maximum value for each task is bolded.

Method	HCheetah	Hopper	Walker2d	Ant
TD3	9532.99	3304.75	4565.24	4185.06
DDPG	3162.50	1731.94	1520.90	816.35
AHE	8401.02	1061.77	2362.13	564.07
AHE + DP	7588.64	1465.11	2459.53	896.13
AHE + TPS	9023.40	907.56	2961.36	872.17
AHE + CDQ	6470.20	1134.14	3979.21	3818.71
TD3 - DP	9590.65	2407.42	4695.50	3754.26
TD3 - TPS	8987.69	2392.59	4033.67	4155.24
TD3 - CDQ	9792.80	1837.32	2579.39	849.75
DQ-AC	9433.87	1773.71	3100.45	2445.97
DDQN-AC	10306.90	2155.75	3116.81	1092.18