

# MODEL-FREE RL WITH VALUE FUNCTION APPROXIMATION CONTINUED

## INTRODUCTION TO REINFORCEMENT LEARNING

**Bogdan Ivanyuk-Skulskyi, Dmytro Kuzmenko**

Department of Mathematics,  
National University of Kyiv-Mohyla Academy

March 13, 2023

# TABLE OF CONTENTS

- ▶ Model-free Function Approximation Convergence
  - **Policy Evaluation**
  - Model-free Control with Linear Function Approximation Convergence
  - Maximization bias
  - Double Q-learning
  - Double DQN

## CONVERGENCE GUARANTEES FOR LINEAR VALUE FUNCTION APPROXIMATION FOR POLICY EVALUATION

- Define the mean squared error of a linear value function approximation for a particular policy relative to the true value as

$$MSVE_{\mu}(w) = \sum_{s \in S} \mu(s) (V^{\pi}(s) - \hat{V}(s; w))^2$$

- where
  - $\mu(s)$ : probability of visiting state  $s$  under policy  $\pi$ . Note  $\sum_s \mu(s) = 1$
  - $\hat{V}^{\pi}(s; w) = x(s)^T w$ , a linear value function approximation
- Monte Carlo policy evaluation with VFA converges to the weights  $w_{MC}$  which has the minimum mean squared error possible with respect to the distribution  $\mu$ :

$$MSVE_{\mu}(w) = \min_w \sum_{s \in S} \mu(s) (V^{\pi}(s) - \hat{V}(s; w))^2$$

## CONVERGENCE GUARANTEES FOR TD LINEAR VFA FOR POLICY EVALUATION: PRELIMINARIES

- ▶ For infinite horizon, the Markov Chain defined by a MDP with a particular policy will eventually converge to a probability distribution over states  $d(s)$
- ▶  $d(s)$  is called the stationary distribution over states of  $\pi$
- ▶  $\sum_s d(s) = 1$
- ▶  $d(s)$  satisfies the following balance equation:

$$d(s') = \sum_s \sum_a \pi(a|s) \pi(s'|s, a) d(s)$$

## CONVERGENCE GUARANTEES FOR LINEAR VALUE FUNCTION APPROXIMATION FOR POLICY EVALUATION

- ▶ Define the mean squared error of a linear value function approximation for a particular policy relative to the true value given the distribution  $d$  as

$$MSVE_d(w) = \sum_{s \in S} d(s) (V^\pi(s) - \hat{V}(s; w))^2$$

- ▶ where
  - $d(s)$ : stationary distribution of  $\pi$  in the true decision process.
  - $\hat{V}^\pi(s; w) = x(s)^T w$ , a linear value function approximation
- ▶ TD(0) policy evaluation with VFA converges to weights  $w_{TD}$  which is within a constant factor of the min mean squared error possible given distribution  $d$ :

$$MSVE_d(w_{TD}) = \frac{1}{1 - \gamma} \min_w \sum_{s \in S} d(s) (V^\pi(s) - \hat{V}^\pi(s; w))^2$$

# TABLE OF CONTENTS

- ▶ Model-free Function Approximation Convergence
  - Policy Evaluation
  - **Model-free Control with Linear Function Approximation Convergence**
  - Maximization bias
  - Double Q-learning
  - Double DQN

## RECALL INCREMENTAL MODEL-FREE CONTROL APPROACHES

- ▶ Similar to policy evaluation, true state-action value function for a state is unknown and so substitute a target value
- ▶ In Monte Carlo methods, use a return  $G_t$  as a substitute target

$$\Delta w = \alpha(G_t - \hat{Q}(s_t, a_t; w)) \nabla_w \hat{Q}(s_t, a_t; w)$$

- ▶ For SARSA instead use a TD target  $r + \gamma \hat{Q}(s', a'; w)$  which leverages the current function approximation value

$$\Delta w = \alpha(r + \gamma \hat{Q}(s', a'; w) - \hat{Q}(s, a; w)) \nabla_w \hat{Q}(s, a; w)$$

- ▶ For Q-learning instead use a TD target  $r + \gamma \max_{a'} \hat{Q}(s', a'; w)$  which leverages the max of the current function approximation value

$$\Delta w = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; w) - \hat{Q}(s, a; w)) \nabla_w \hat{Q}(s, a; w)$$

## CONVERGENCE OF TD METHODS WITH VFA

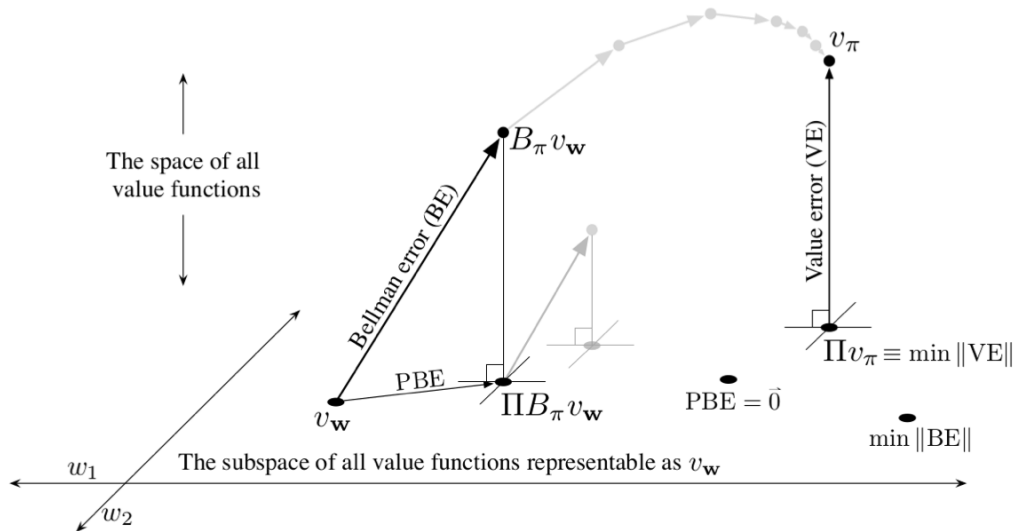
- ▶ Informally, updates involve doing an (approximate) Bellman backup followed by best trying to fit underlying value function to a particular feature representation
- ▶ Bellman operators are contractions, but value function approximation fitting can be an expansion



## ACTIVE AREA: OFF POLICY LEARNING WITH FUNCTION APPROXIMATION

- ▶ Extensive work in better TD-style algorithms with value function approximation, some with convergence guarantees: see Chp 11 SB

# VALUE FUNCTION APPROXIMATION



# TABLE OF CONTENTS

- ▶ Model-free Function Approximation Convergence
  - Policy Evaluation
  - Model-free Control with Linear Function Approximation Convergence
  - **Maximization bias**
  - Double Q-learning
  - Double DQN

## MAXIMIZATION BIAS PROOF

- ▶ Consider single-state MDP ( $|S| = 1$ ) with 2 actions, and both actions have 0-mean random rewards, ( $\mathbb{E}(r|a = a_1) = \mathbb{E}(r|a = a_2) = 0$ ).
- ▶ Then  $Q(s, a_1) = Q(s, a_2) = 0 = V(s)$
- ▶ Assume there are prior samples of taking action  $a_1$  and  $a_2$
- ▶ Let  $\hat{Q}(s, a_1), \hat{Q}(s, a_2)$  be the finite sample estimate of  $Q$
- ▶ Use an unbiased estimator for  $Q$ : e.g.  $\hat{Q}(s, a_1) = \frac{1}{n(s, a_1)} \sum_{i=1}^{n(s, a_1)} r_i(s, a_1)$
- ▶ Let  $\hat{\pi} = \operatorname{argmax}_a \hat{Q}(s, a)$  be the greedy policy w.r.t. the estimated  $\hat{Q}$
- ▶ Even though each estimate of the state-action values is unbiased, the estimate of  $\hat{\pi}$ 's value  $\hat{V}^{\hat{\pi}}$  can be biased:

$$\hat{V}^{\hat{\pi}}(s) = \mathbb{E}[\max \hat{Q}(s, a_1), \hat{Q}(s, a_2)] \geq \max[\mathbb{E}[\hat{Q}(s, a_1)], [\hat{Q}(s, a_2)]] = \max[0, 0] = V^{\pi},$$

where the inequality comes from Jensen's inequality.

# TABLE OF CONTENTS

- ▶ Model-free Function Approximation Convergence
  - Policy Evaluation
  - Model-free Control with Linear Function Approximation Convergence
  - Maximization bias
  - **Double Q-learning**
  - Double DQN

## DOUBLE Q-LEARNING

- ▶ The greedy policy w.r.t. estimated Q values can yield a maximization bias during finite-sample learning
- ▶ Avoid using max of estimates as estimate of max of true values
- ▶ Instead split samples and use to create two independent unbiased estimates of  $Q_1(s_1, a_i)$  and  $Q_2(s_1, a_i) \forall a$ .
  - Use one estimate to select max action:  $a^* = \operatorname{argmax}_a Q_1(s_1, a)$
  - Use other estimate to estimate value of  $a^* : Q_2(s, a^*)$
  - Yields unbiased estimate:  $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$

## DOUBLE Q-LEARNING

- ▶ The greedy policy w.r.t. estimated Q values can yield a maximization bias during finite-sample learning
- ▶ Avoid using max of estimates as estimate of max of true values
- ▶ Instead split samples and use to create two independent unbiased estimates of  $Q_1(s_1, a_i)$  and  $Q_2(s_1, a_i) \forall a$ .
  - Use one estimate to select max action:  $a^* = \operatorname{argmax}_a Q_1(s_1, a)$
  - Use other estimate to estimate value of  $a^* : Q_2(s, a^*)$
  - Yields unbiased estimate:  $\mathbb{E}(Q_2(s, a^*)) = Q(s, a^*)$
- ▶ Why is this an unbiased estimate of the max state-action value?  
Using independent samples to estimate the value
- ▶ If acting online, can alternate samples used to update  $Q_1$  and  $Q_2$ , using the other to select the action chosen

## DOUBLE Q-LEARNING

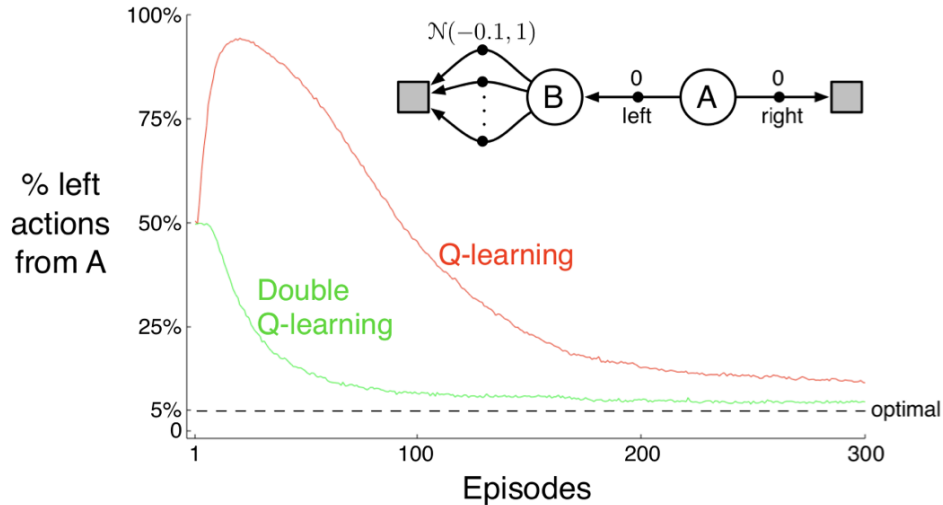
- ▶ Initialize  $Q_1(s, a)$  and  $Q_2(s, a)$ ,  $\forall s \in S, a \in A$   $t = 0$ , initial state  $s_t = s_0$
- ▶ loop
  - Select  $a_t$  using  $\epsilon$ -greedy  $\pi(s) = \operatorname{argmax}_a Q_1(s_t, a) + Q_2(s_t, a)$
  - Observe  $(r_t, s_{t+1})$
  - if (with prob. 0.5) then:
    - ▶  $Q_1(s_t, a_t) \leftarrow Q_1(s_t, a_t) + \alpha(r_t + \gamma Q_2(s_{t+1}, \operatorname{argmax}_a Q_1(s_{t+1}, a)) - Q_1(s_t, a_t))$
  - else
    - ▶  $Q_2(s_t, a_t) \leftarrow Q_2(s_t, a_t) + \alpha(r_t + \gamma Q_1(s_{t+1}, \operatorname{argmax}_a Q_2(s_{t+1}, a)) - Q_2(s_t, a_t))$
  - end if
  - $t = t + 1$

Compared to Q-learning, how does this change the: memory requirements, computation requirements per step, amount of data required?

Doubles the memory, same computation requirements, data requirements are subtle– might reduce amount of exploration needed due to lower bias



## DOUBLE Q-LEARNING (FIGURE 6.7 IN SUTTON AND BARTO 2018)



Due to the maximization bias, Q-learning spends much more time selecting suboptimal actions than double Q-learning.

# TABLE OF CONTENTS

- ▶ Model-free Function Approximation Convergence
  - Policy Evaluation
  - Model-free Control with Linear Function Approximation Convergence
  - Maximization bias
  - Double Q-learning
  - **Double DQN**

## RECALL DQN

- ▶ Deep Q-learning (DQN): Q-learning with deep neural networks and
  - Experience replay
  - Fixed Q-targets

$$\Delta w = \alpha(r + \gamma \max_{a'} \hat{Q}(s', a'; w) - \hat{Q}(s, a; w)) \nabla_w \hat{Q}(s, a; w)$$

## DOUBLE DQN

- ▶ Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- ▶ Extend double Q learning to DQN
- ▶ Current Q-network  $w$  is used to select actions
- ▶ Older Q-network  $w^-$  is used to evaluate actions

$$\Delta w = \alpha(r + \gamma \hat{Q}(\operatorname{argmax}_{a'} \hat{Q}(s', a'; w); w^-) - \hat{Q}(s, a; w))$$

## DOUBLE DQN

- ▶ Double DQN (Deep Reinforcement Learning with Double Q-Learning, Van Hasselt et al, AAAI 2016)
- ▶ Extend double Q learning to DQN
- ▶ Current Q-network  $w$  is used to select actions
- ▶ Older Q-network  $w^-$  is used to evaluate actions

$$\Delta w = \alpha(r + \gamma \hat{Q}(\operatorname{argmax}_{a'} \hat{Q}(s', a'; w); w^-) - \hat{Q}(s, a; w))$$

- ▶ In DQN the same weights  $w$  were used to choose the best action at  $s'$  and evaluate its value  $\hat{Q}(s', a'; w^-)$

# MODEL-FREE VALUE FUNCTION APPROXIMATION RL: WHAT YOU SHOULD KNOW

- ▶ Be able to derive weight update for generic function approximation for  $Q/V^\pi$
- ▶ Understand various (MC/SARSA/Q-learning) targets used when updating Q function
- ▶ Know what TD vs MC converge to for policy evaluation with a linear function approximator
- ▶ Be able to implement DQN
- ▶ Define the maximization bias and give one tool for alleviating it