

Averaged-DQN: Variance Reduction and Stabilization for Deep Reinforcement Learning

Oron Anschel

Joint work Nir Baram, and Nahum Shimkin

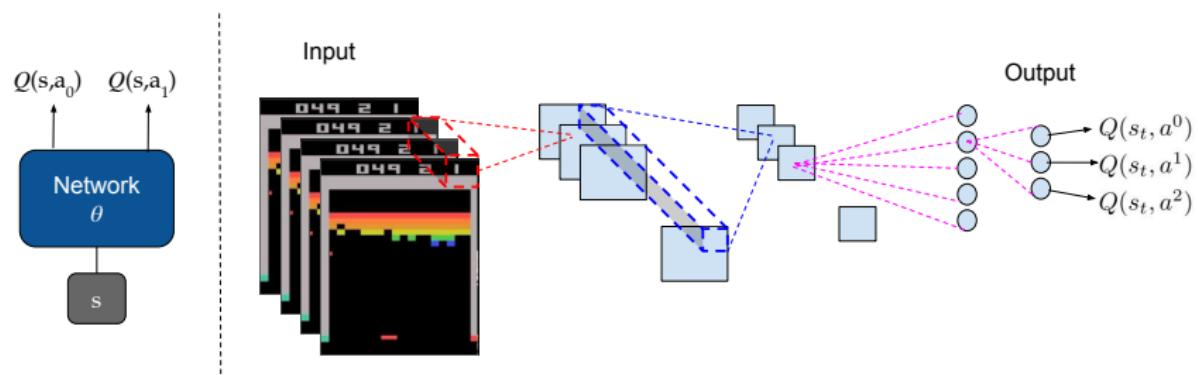
Department of Electrical Engineering
Technion - Israel Institute of Technology

Overview

- 1 Background & Motivation
- 2 Averaged-DQN
- 3 Overestimation and Approximation Errors
- 4 TAE Variance Reduction
- 5 Experiments

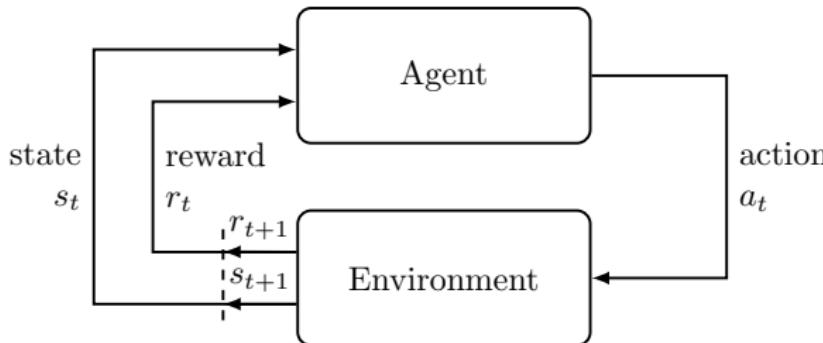
DQN

- Mnih, et al., (2013)
 - Target Network.
 - Experience Buffer (Lin et al., 1993).



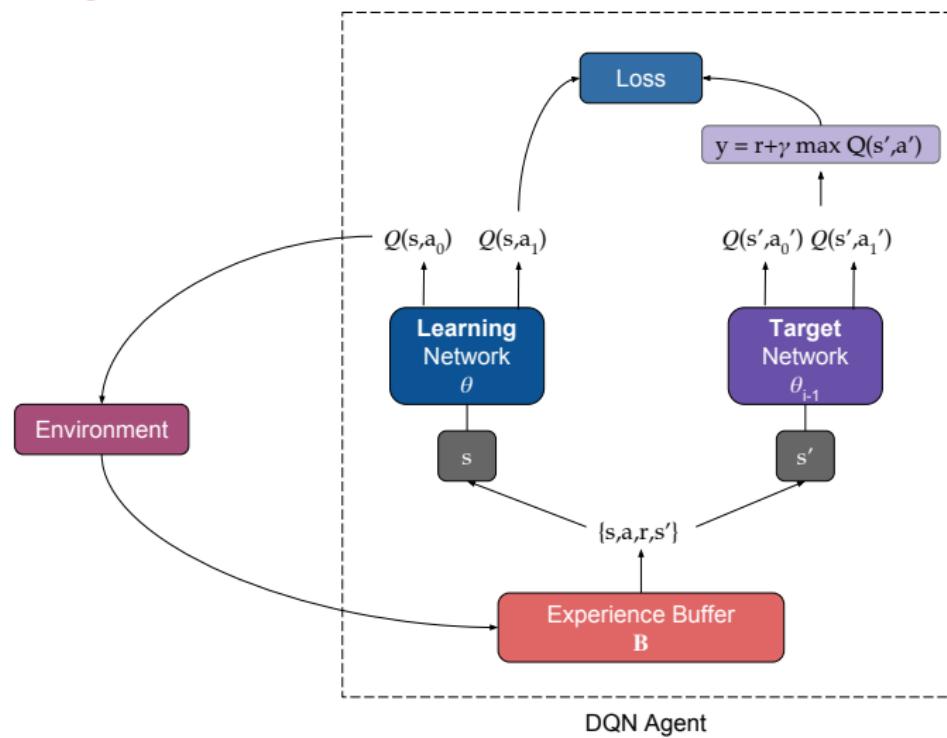
Reinforcement Learning

The **Agent-Environment** Interface



DQN Agent- Environment Interface

The **Agent-Environment** Interface



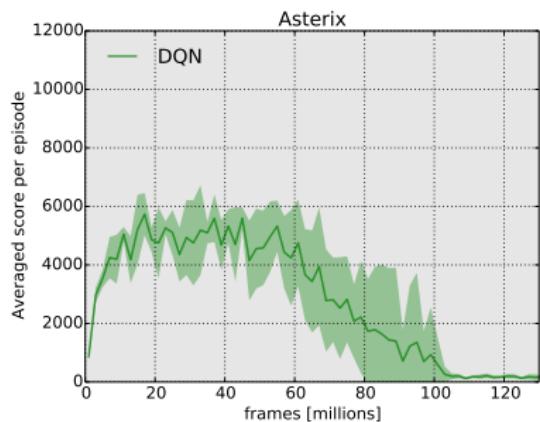
Motivational Example 1



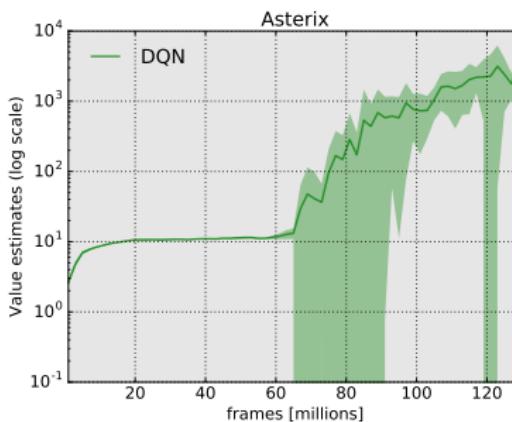
Example 1: **Instability and Variability**

- Upside: Super-human control.
- Downside: **Best model** evaluation methodology.

Motivational Example 2



(average score)



(value estimate)

Example 2:
Overestimation¹

¹H van Hasselt et al., 2015

- Overestimation → asymptotically sub-optimal policies.

Motivation Cont'd

- Are the two examples related?

Overview

- 1** Background & Motivation
- 2** Averaged-DQN
- 3** Overestimation and Approximation Errors
- 4** TAE Variance Reduction
- 5** Experiments

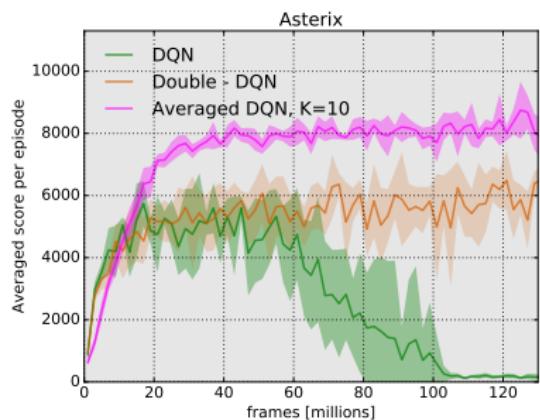
Motivational Example 1

Example 1:

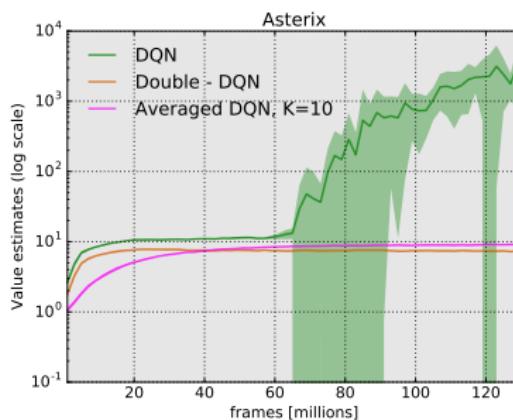
Instability and Variability

- Upside effects: **Stable training** and **better policies**.
- Bonus: Best model → Last model evaluation.

Motivational Example 2



(average score)



(value estimate)

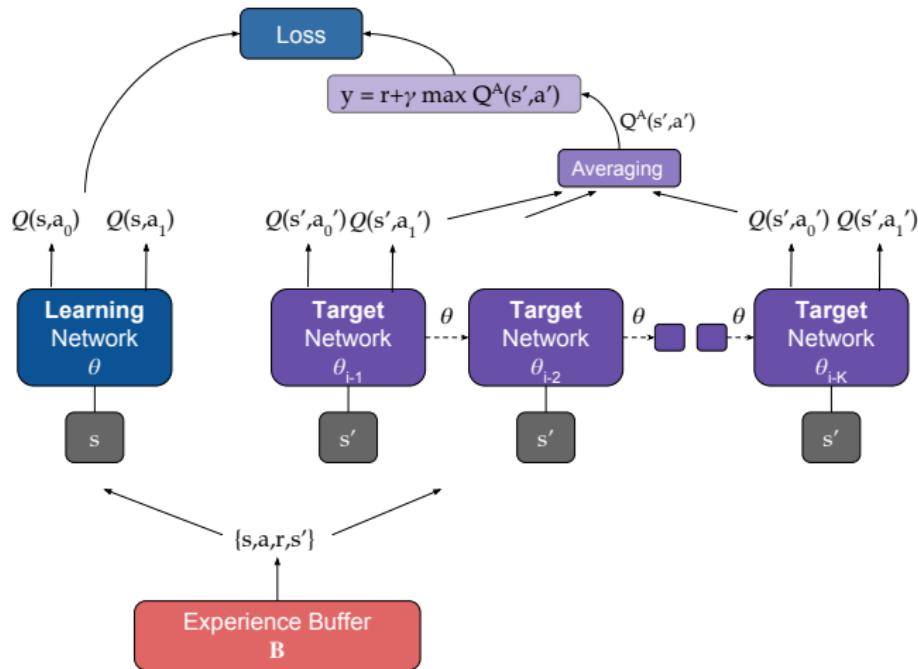
Example 2:
Overestimation¹

¹H van Hasselt et al., 2015

- Upside effect: Reduced overestimation.

Averaged DQN Agent

The Averaged-DQN Agent



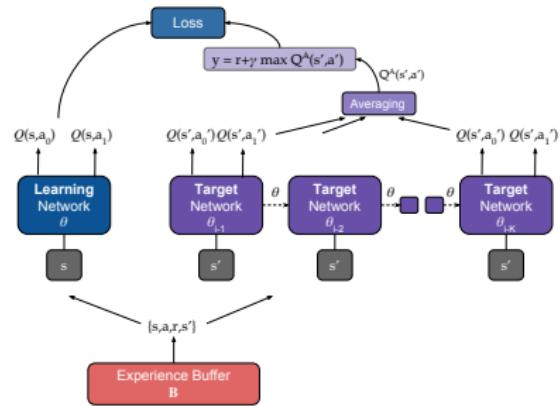
Averaged DQN

Algorithm Averaged DQN

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: Explore(\cdot), update \mathcal{B}
- 6: $Q_{i-1}^A(s, a) = \frac{1}{K} \sum_{k=1}^K Q(s, a; \theta_{i-k})$
- 7: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q_{i-1}^A(s', a') | s, a]$
- 8: $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

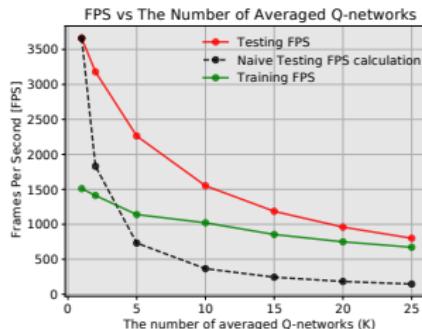
Output: $Q_N^A(s, a) = \frac{1}{K} \sum_{k=0}^{K-1} Q(s, a; \theta_{N-k})$

*In red: differences from DQN.



Computational Cost

- **Training:** $K - 1$ more forward passes (each minimization of the DQN loss).
Back-propagation updates → same as in DQN.
- **Testing:** $K - 1$ more forward passes (for each evaluation of the Q-function).



- *In general, the forward-passes through the K Q-networks can be done in parallel.
- *In practice, since the same state is forwarded, the GPU implicitly parallelize the computation.

Overview

- 1 Background & Motivation
- 2 Averaged-DQN
- 3 Overestimation and Approximation Errors
- 4 TAE Variance Reduction
- 5 Experiments

Target Approximation Error (TAE)

- TAE: $Z_{s,a}^i = Q(s, a; \theta_i) - y_{s,a}^i$
- The error in $Q(s, a; \theta_i)$ relative to $y_{s,a}^i$.

Algorithm DQN (Mnih, 2013)

```

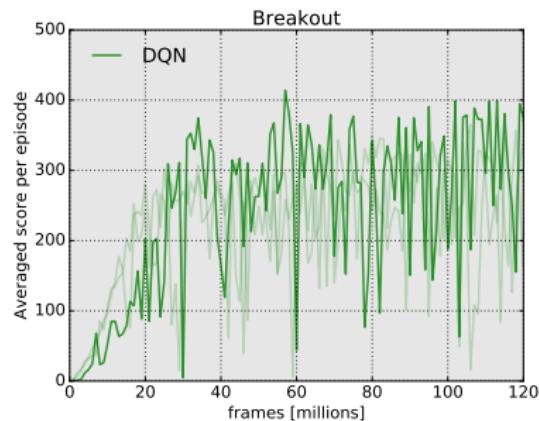
1: Initialize  $Q(s, a; \theta)$  with random weights  $\theta_0$ 
2: Initialize exploration procedure  $\text{Explore}(\cdot)$ 
3: Initialize Experience Replay (ER) buffer  $\mathcal{B}$ 
4: for  $i = 1, 2, \dots, N$  do
5:    $\text{Explore}(\cdot)$ , update  $\mathcal{B}$ 
6:    $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ 
7:    $\theta_i \approx \text{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$ 

```

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- Sub-optimality of θ_i .
- Model expressiveness power.
- Finite ER buffer \mathcal{B} .



Target Approximation Error (TAE)

- TAE: $Z_{s,a}^i = Q(s, a; \theta_i) - y_{s,a}^i$
- The error in $Q(s, a; \theta_i)$ relative to $y_{s,a}^i$.

Algorithm DQN (Mnih, 2013)

```

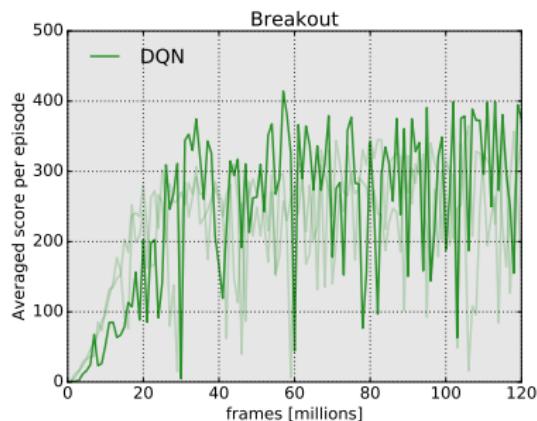
1: Initialize  $Q(s, a; \theta)$  with random weights  $\theta_0$ 
2: Initialize exploration procedure  $\text{Explore}(\cdot)$ 
3: Initialize Experience Replay (ER) buffer  $\mathcal{B}$ 
4: for  $i = 1, 2, \dots, N$  do
5:    $\text{Explore}(\cdot)$ , update  $\mathcal{B}$ 
6:    $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ 
7:    $\theta_i \approx \text{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$ 

```

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- Sub-optimality of θ_i .
- Model expressiveness power.
- Finite ER buffer \mathcal{B} .



Target Approximation Error (TAE)

- TAE: $Z_{s,a}^i = Q(s, a; \theta_i) - y_{s,a}^i$
- The error in $Q(s, a; \theta_i)$ relative to $y_{s,a}^i$.

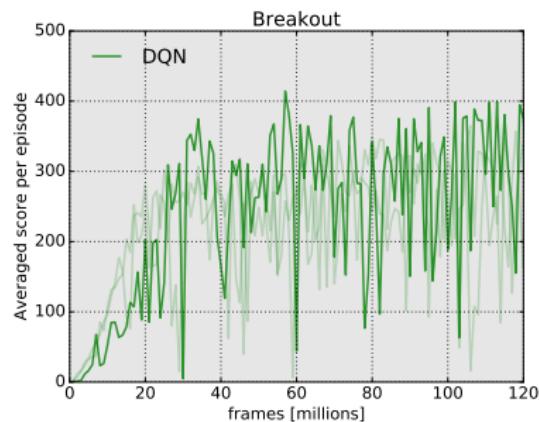
Algorithm DQN (Mnih, 2013)

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: Explore(\cdot), update \mathcal{B}
- 6: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- 7: $\theta_i \approx \underset{\theta}{\text{argmin}} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- Sub-optimality of θ_i .
- Model expressiveness power.
- Finite ER buffer \mathcal{B} .



Overestimation

- Overestimation: $R_{s,a}^i = y_{s,a}^i - \hat{y}_{s,a}^i$
- $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- $\hat{y}_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} (y_{s',a'}^{i-1}) | s, a]$ (without TAE)

Algorithm DQN (Mnih, 2013)

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: Explore(\cdot), update \mathcal{B}
- 6: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- 7: $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- TAE variance.
- Similar Q-values for different actions.

¹Thrun and Shwartz, 1993

Overestimation

- Overestimation: $R_{s,a}^i = y_{s,a}^i - \hat{y}_{s,a}^i$
- $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- $\hat{y}_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} (y_{s',a'}^{i-1}) | s, a]$ (without TAE)

Algorithm DQN (Mnih, 2013)

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: $\text{Explore}(\cdot)$, update \mathcal{B}
- 6: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- 7: $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- TAE variance.
- Similar Q-values for different actions.

¹Thrun and Shwartz, 1993

Overestimation

- Overestimation: $R_{s,a}^i = y_{s,a}^i - \hat{y}_{s,a}^i$
- $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- $\hat{y}_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} (y_{s',a'}^{i-1}) | s, a]$ (without TAE)

Algorithm DQN (Mnih, 2013)

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: $\text{Explore}(\cdot)$, update \mathcal{B}
- 6: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- 7: $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

- TAE variance.
- Similar Q-values for different actions.

¹Thrun and Shwartz, 1993

Overestimation

- Overestimation: $R_{s,a}^i = y_{s,a}^i - \hat{y}_{s,a}^i$
- $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- $\hat{y}_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} (y_{s',a'}^{i-1}) | s, a]$ (without TAE)

Algorithm DQN (Mnih, 2013)

- 1: Initialize $Q(s, a; \theta)$ with random weights θ_0
- 2: Initialize exploration procedure $\text{Explore}(\cdot)$
- 3: Initialize Experience Replay (ER) buffer \mathcal{B}
- 4: **for** $i = 1, 2, \dots, N$ **do**
- 5: $\text{Explore}(\cdot)$, update \mathcal{B}
- 6: $y_{s,a}^i = \mathbb{E}_{\mathcal{B}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$
- 7: $\theta_i \approx \operatorname{argmin}_{\theta} \mathbb{E}_{\mathcal{B}} [(y_{s,a}^i - Q(s, a; \theta))^2]$

Output: $Q^{\text{DQN}}(s, a; \theta_N)$

Error source

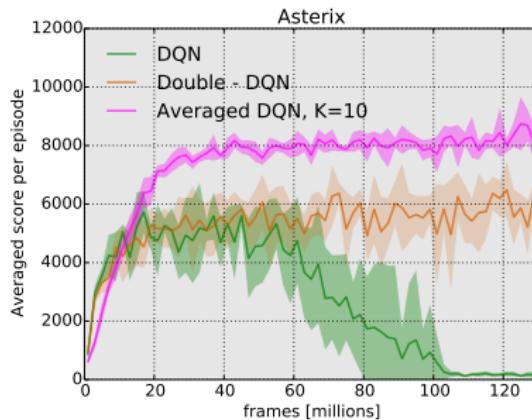
- TAE variance.
- Similar Q-values for different actions.

¹Thrun and Shwartz, 1993

Overestimation Effect

Overestimation effect:

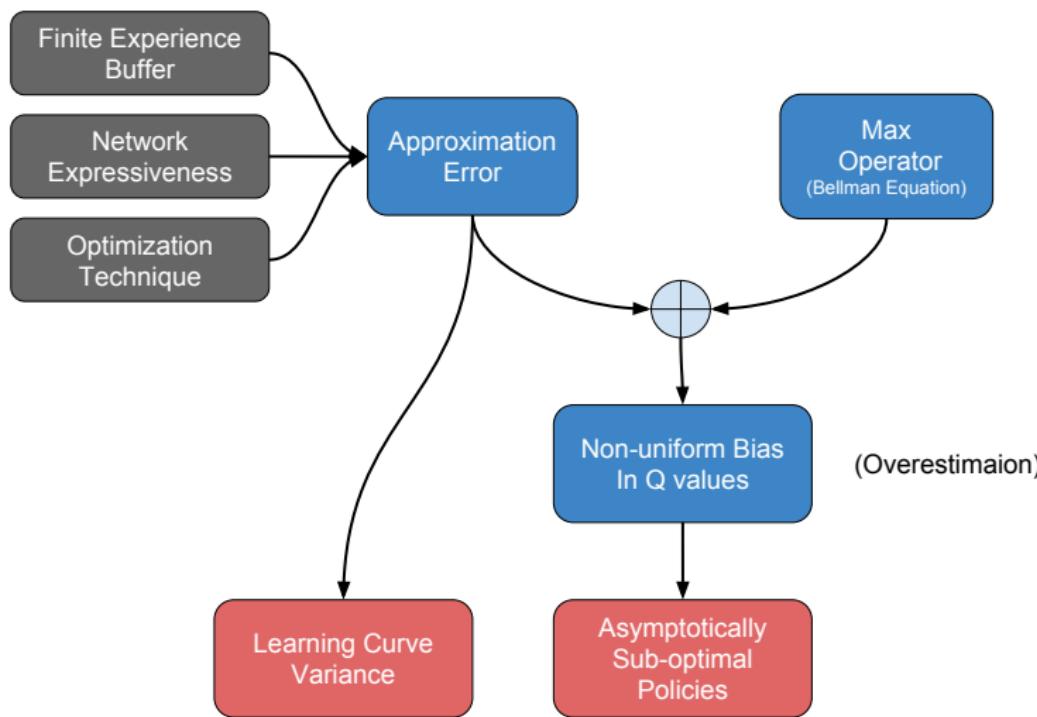
- 1 Non-uniform bias → Asymptotic sub-optimal policies^{1,2}.



¹Thrun and Shwartz, 1993

²H van Hasselt et al., 2015

TAE and Overestimation Summary



Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

- Bootstrapped-DQN (Osband, et al. 2016).

Other Options

Directly dealing with the TAE sources:

- Prioritized experience replay (Schaul, et al. 2015).
- Dueling network architectures (Wang, et al., 2015).
- Smaller SGD learning rate/ Better optimizers.

Changing the Bellmann Operator:

- Double-DQN (Hasselt, et al. 2015).

Ensemble Techniques:

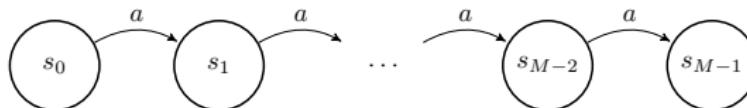
- Bootstrapped-DQN (Osband, et al. 2016).

Overview

- 1** Background & Motivation
- 2** Averaged-DQN
- 3** Overestimation and Approximation Errors
- 4** TAE Variance Reduction
- 5** Experiments

Averaged-DQN Variance Reduction

Unidirectional single-chain MDP



Result

For $K > 1, M > 1$:

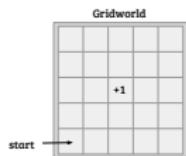
$$\begin{aligned}\text{Var}[Q_i^{\text{Averaged-DQN}}(s_0, a)] &< \text{Var}[Q_i^{\text{Ensemble-DQN}}(s_0, a)] \\ &= \frac{1}{K} \text{Var}[Q^{\text{DQN}}(s_0, a; \theta_i)]\end{aligned}$$

*Averaged-DQN averages over averages.

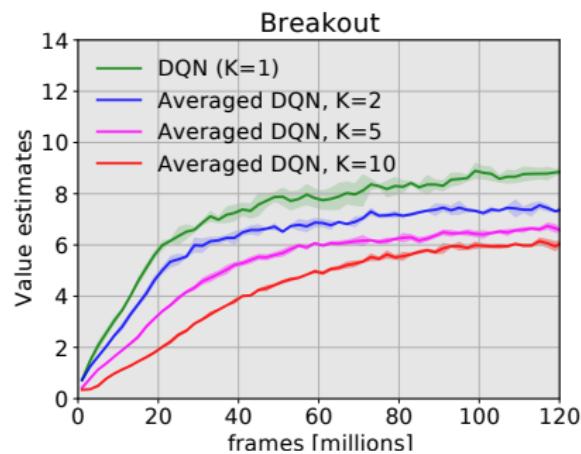
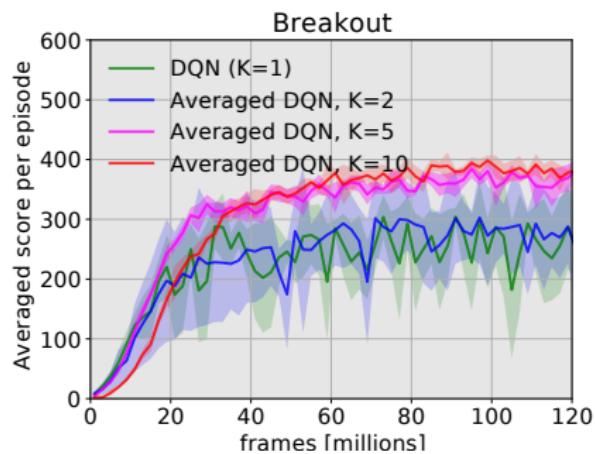
Overview

- 1 Background & Motivation
- 2 Averaged-DQN
- 3 Overestimation and Approximation Errors
- 4 TAE Variance Reduction
- 5 Experiments

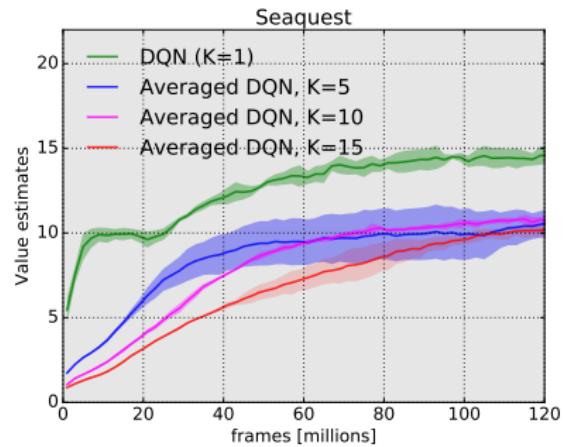
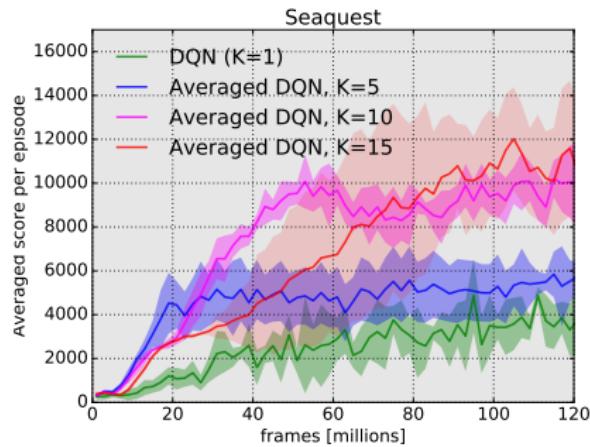
Gridworld Experiments

**DQN****Averaged-DQN (K=15)**

Arcade Learning Environment (Breakout)



Arcade Learning Environment (Seaquest)



Thank you!

oronanschel@campus.technion.ac.il

Talk to me at poster #99

Slides at <https://goo.gl/VQQ3Wj>