Approximate Value Iteration and Variants

Chen-Yu Wei

Value Iteration

For
$$k = 1, 2, ...$$

 $\forall s, a, \qquad Q_k(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{k-1}(s', a')$
unknown
unknown

Idea: In each iteration, use multiple samples to estimate the right-hand side.

Value Iteration with Samples
$$\underbrace{\sum_{i, A_{i}, Y_{i}} \underbrace{\sum_{k_{i}, Y_{i}$$

Recall: Contextual Bandits with Regression

 (x_i, a_i, r_i)



Train \hat{R} such that $\hat{R}(x_i, a_i) \approx r_i$

Value Iteration with Regression

 $(s_i, a_i, r_i, \frac{s'_i}{s'_i})$



Train Q_k such that $Q_k(s_i, a_i) \approx r_i + \gamma \max_{a'} Q_{k-1}(s'_i, a')$

This is just one iteration of Value Iteration

Value Iteration with Samples

For
$$k = 1, 2, ..., N$$
:
For $i = 1, 2, ..., N$:
Choose action $a_i \sim EG(Q_{\theta_k}(s_i, \cdot))$
Receive reward $r_i \sim R(s_i, a_i)$ and $s'_i \sim P(\cdot | s_i, a_i)$
 $s_{i+1} = s'_i$ if episode continues, $s_{i+1} \sim \rho$ if episode ends
 $\theta \leftarrow \theta_k$
For $m = 1, 2, ..., M$:
Randomly pick an i (or a mini-batch) from $\{1, 2, ..., N\}$
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a')\right)^2$
 $\theta_{k+1} \leftarrow \theta$
Perform one iteration of Value Iteration

2nd for-loop: trying to find $\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^{N} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \underset{a'}{\max} Q_{\theta_k}(s'_i, a') \right)^2$

It is Valid to Reuse Samples



The algorithm in the previous slide only use \mathcal{D}_k to train θ_{k+1} .

However, as the reward function *R* and transition *P* remains unchanged, it is valid (actually, even better) to reuse samples:



Benefits of Reusing Samples

- Improving data efficiency
 - Every sample is used multiple times in training just like we usually go through multiple epochs for supervised learning tasks.
- The buffer \mathcal{B} will consist of a wider range of state-actions
 - It allows better approximation of

$$\forall s, a, \qquad Q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_k(s', a')$$

Value Iteration with Reused Samples (= Deep Q-Learning or DQN)



Another Popular Implementation

```
Initialize \mathcal{B} = \{\} \in \mathsf{Replay buffer}
For k = 1, 2, ...
   For i = 1, 2, ..., N:
          Choose action a_i \sim EG(Q_\theta(s_i, \cdot))
          Receive reward r_i \sim R(s_i, a_i) and s'_i \sim P(\cdot | s_i, a_i)
          s_{i+1} = s'_i if episode continues, s_{i+1} \sim \rho if episode ends
          Insert (s_i, a_i, r_i, s'_i) to \mathcal{B}
    For m = 1, 2, ..., M:
```

HW4 task

Randomly pick an i (or a mini-batch) from \mathcal{B}

When Does DQN Succeed?

DQN tries to approximate Value Iteration by solving

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \sum_{i \in \mathcal{B}} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \underset{a'}{\max} Q_{\theta_k}(s'_i, a') \right)^2 \tag{1}$$

The true Value Iteration:

$$\forall s, a, \qquad Q_{k+1}(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_k(s', a')$$
(2)

Under what conditions can (1) well approximate (2)?

- *B* should contain a wide range of state-action pairs (a challenge of **exploration**)
- $Q_{\theta_{k+1}}(s, a)$ should recover $R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$ well for all state-actions (a challenge of function approximation, or generalization)

1. Exploration in MDPs (Not Easy)



Environment:

- Fixed-horizon MDP with episode length *H*
- Initial state at 0
- A single rewarding state at state *H*
- Actions: Go LEFT or RIGHT

Suppose we perform DQN with ϵ -greedy with random initialization \Rightarrow On average, we need 2^{*H*} episodes to see the reward (before that, we won't make any meaningful update and will just do random walk around state 0)



Key issue:

- The ε-greedy strategy (or BE, IGW) performs action-space exploration but not state-space exploration.
- This problem becomes more severe when the reward signal is **sparse**.
- To solve this, we usually require the exploration bonus (a form of reward shaping) technique will be covered much later.

At this point (for the discussion of DQN), we pretend that EG, BE, or IGW will lead to sufficient exploration over the state space.

1. Exploration in MDPs (Not Easy)

Classic sparse-reward environments:



Mountain Car



Montezuma's Revenge

2. Function Approximation

To make DQN well approximate VI, we need

$$\forall s, a \qquad Q_{\theta_{k+1}}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$$

(ϵ -approximate) Bellman Completeness an assumption both on the MDP and the function expressiveness

$$\forall \theta', \exists \theta \quad \forall s, a, \qquad \left| Q_{\theta}(s, a) - \left(R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta'}(s', a') \right) \right| \le \epsilon$$

This allows us to quantify the regression error in each iteration.

2. Function Approximation

In HW1 you have shown



In value-based contextual bandits, the requirement / assumption for function approximation is

$$\exists \theta \quad \forall x, a \qquad R_{\theta}(x, a) \approx R(x, a)$$

In value-based MDPs, the requirement / assumption for function approximation is

$$V_{\theta'}, \exists \theta \forall s, a \qquad Q_{\theta}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta'}(s', a') \qquad \beta \subset a \text{ supposed}$$

Analysis of DQN assuming sufficient exploration and Bellman Completeness

Recall the analysis for the exact Value Iteration:



Completing the Analysis of VI (1st Step)

$$\begin{aligned} & \bigcup_{k} (s, \alpha) = \mathcal{R}(s, \alpha) + \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k,1}(s', \alpha') \\ & \bigcirc_{k} (s, \alpha) = \mathcal{R}(s, \alpha) + \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k-2}(s', \alpha') \\ & \bigcirc_{k-1}(s, \alpha) = \mathcal{R}(s, \alpha) + \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k-2}(s', \alpha') \\ & \bigcirc_{k} (s, \alpha) - \mathcal{O}_{k-1}(s, \alpha) = \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k-1}(s', \alpha') - \max_{a'} \mathcal{O}_{k-2}(s', \alpha') \\ & \square_{k} f(x) - \max_{x} g(x) \\ & = \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & = \mathcal{Y} \sum_{s'} \mathcal{P}(s'|s, \alpha) \max_{a'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigcirc_{k} (s, \alpha) - \mathcal{O}_{k}(s, \alpha) \\ & \bigcirc_{k} (s, \alpha) - \mathcal{O}_{k-1}(s, \beta) \\ & \bigcirc_{k} (s, \alpha) - \mathcal{O}_{k-1}(s, \beta) \\ & \bigcirc_{k} (s, \alpha) - \mathcal{O}_{k-1}(s, \beta) \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-2}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha') - \mathcal{O}_{k-1}(s', \alpha') \\ & \bigvee_{s', n'} \mathcal{O}_{k-1}(s', \alpha')$$

Analysis of DQN assuming sufficient exploration and Bellman Completeness

$$\frac{m_{M}}{M_{M}} \underbrace{Q_{k}(5,5) - Q_{k,1}(5,c)}_{M_{k,1}} \leq \int_{m_{k}}^{m_{k}} \underbrace{Q_{k,1}(5,c) - Q_{k,2}(5,c)}_{M_{k,2}} + \underbrace{b_{k}(5,c) - b_{k,1}(5,c)}_{K_{k,2}} \leq \int_{m_{k,1}}^{m_{k}} \underbrace{Q_{k-2}(5,c) - Q_{k-2}(5,c)}_{K_{k}} + \frac{b_{k}(5,c) - b_{k}(5,c)}{K_{k}} + \frac{b_{k}(5,c) - b_{k}(5,c$$

DQN can be offline

Let \mathcal{B} consists of (s, a, r, s') tuples collected by a mixture of **arbitrary policies.**

For k = 1, 2, ... $\theta \leftarrow \theta_k$ For m = 1, 2, ..., M: Randomly pick an *i* (or a mini-batch) from \mathcal{B} $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta_k}(s'_i, a') \right)^2$ $\theta_{k+1} \leftarrow \theta$

Data collection

Again, its success relies on 1) \mathcal{B} contains data with sufficiently wide range of state-actions, 2) Bellman completeness.

The same theoretical analysis applies.

Handling the Non-Ideal Case

When DQN cannot well-approximate VI

In practice,

- We may not be able to collect sufficiently wide range of state-actions
- Bellman completeness may not hold

In either case, we may not have

$$\forall s, a \quad Q_{\theta_{k+1}}(s, a) \approx R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta_k}(s', a')$$

This makes our previous analysis based on VI fails.

When DQN cannot well-approximate VI

In this case, $Q_{\theta_k}(s, a)$ tends to **overestimate** $Q^*(s, a)$, and the greedy policy $\hat{\pi}(s) = \underset{a}{\operatorname{argmax}} Q_{\theta_k}(s, a)$ could be very bad.

5,9

When DQN cannot well-approximate VI

Such "seeking the error" behavior is due to "bootstrapping"

• An issue only in MDP but not in bandits

To prevent overestimation, two strategies are

- Double Q-learning: decorrelating the choice of argmax action and the error of the value function
- Conservative Q-learning: being conservative

Double DQN (v1)



Double DQN (v1)



Double DQN (v2)



Hado van Hasselt, Arthur Guez, David Silver. Deep Reinforcement Learning with Double Q-learning. 2015.

Conservative Q-learning (CQL)

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \sum_{i \in \mathcal{B}} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \underset{a'}{\max} Q_{\theta_k}(s'_i, a') \right)^2 + \alpha \sum_{i \in \mathcal{B}} \left(\log \left(\sum_{a} \exp(Q_{\theta}(s_i, a)) \right) - Q_{\theta}(s_i, a_i) \right)$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_{i \in \mathcal{B}} \left(Q_{\theta}(s_{i}, a_{i}) - r_{i} - \gamma \underset{a'}{\max} Q_{\theta_{k}}(s_{i}', a') \right)^{2} + \alpha \sum_{i \in \mathcal{B}} \left(\max_{\mu} \sum_{a} \mu(a|s_{i}) Q_{\theta}(s_{i}, a) - Q_{\theta}(s_{i}, a_{i}) + \operatorname{KL}(\mu(\cdot|s_{i}), \operatorname{Unif}) \right)$$

Aviral Kumar, Aurick Zhou, George Tucker, Sergey Levine Conservative Q-Learning for Offline Reinforcement Learning. 2020.

Comparison

- Double-Q: make the $\underset{a}{\operatorname{argmax}} Q_{\theta}(s, a)$ choice decoupled from θ
- Conservative-Q: mitigate the overestimation of $\max_{a} Q_{\theta}(s_i, a)$ over $Q_{\theta}(s_i, a_i)$

Summary for DQN

- Motivation: approximating Value Iteration using samples and function approximation
- Standard elements: target network, replay buffer
- Work as desired when both of the following conditions hold:
 - The learner is able to obtain exploratory data (online or offline)
 - Neural network is sufficiently expressive: Bellman completeness
- When the conditions above do not hold
 - Tends to overestimate *Q* values and suggest arbitrary actions
- Solutions
 - Double Q-learning
 - Conservative Q-learning

Improvements on DQN

- Dueling DDQN
- Prioritized replay
- Distributional DQN

• ...





Other Variants that Fail

An Unstable Variant
DQN without target network
For
$$k = 1, 2, ...$$

Randomly pick an i (or a mini-batch) from \mathcal{B}
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(\underline{s}_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(\underline{s}'_{i}, a') \right)^{2}$
 $\overline{\theta} \leftarrow \theta$
cf. DQN with target network
For $k = 1, 2, ...$
Randomly pick an i (or a mini-batch) from \mathcal{B}
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(\underline{s}_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(\underline{s}'_{i}, a') \right)^{2}$
f. DQN with target network
For $k = 1, 2, ...$
Randomly pick an i (or a mini-batch) from \mathcal{B}
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(\underline{s}_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(\underline{s}'_{i}, a') \right)^{2}$
 $\overline{\theta} \leftarrow (1 - \tau)\overline{\theta} + \tau \theta$
For $k = 1, 2, ...$
Randomly pick an i (or a mini-batch) from \mathcal{B}
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(\underline{s}_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(\underline{s}'_{i}, a') \right)^{2}$
 $\overline{\theta} \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(\underline{s}_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(\underline{s}'_{i}, a') \right)^{2}$
 $\overline{\theta} \leftarrow \theta$

An Unstable Variant

Could diverge for off-policy problems (e.g. Deep Q-learning) even when exploration assumption and Bellman completeness hold.

Off-policy: Training Q^{π} using data collected from some other policy $\pi' \neq \pi$ or some distillation very different from the distribution induced by π

$$r = 1$$

$$s_{1}$$

$$r = 1$$

$$s_{2}$$

$$\phi(s_{1}, a) = (1, 0)$$

$$\phi(s_{2}, a) = (2, 1)$$

$$r = 0$$

$$r = 0$$

$$s_{2}$$

 $Q_{\mathbf{A}}(s,s) = \phi(s,a)^T p$



Simplified from the "Baird's counterexample" (see Sutton and Barto Section 11.2)

The Effect of Target Network

Let
$$KN = 10000$$

For $k = 1, 2, ..., K$
 $\theta_k \leftarrow \theta$
For $i = 1, ..., N$:
Sample $(s, a, r, s') \sim$ Uniform $\{(\underline{s_1, a, 1, s_2}), (\underline{s_2, a, 0, s_2})\}$
 $\theta \leftarrow \theta - \alpha \left(\phi(s, a)^{\mathsf{T}}\theta - r - \gamma \phi(s', a)^{\mathsf{T}}\theta_k\right)\phi(s, a)$
 $\theta_{k+1} \leftarrow \theta$

The Effect of Target Network









Evolution of $\theta(0)$ and $\theta(1)$ over time $\theta(0)$ over time 200 - — θ(1) over time — θ_{tar}(0) over time θ_{tar}(1) over time 100 Valt -100 --200 ò 20000 40000 60000 80000 100000 Iteration





N=190













N=800

An Unstable Variant

- This type of algorithm (i.e., perform regression without fixing the target) is still useful in certain cases:
 - **On-policy** problems where Q^{π} is trained using data collected from π will see this later
 - Off-policy Q-learning without function approximation this is the very original "Qlearning" algorithm being proposed by Chris Watkins in 1989.
- However, in Deep Q-learning, this variant is generally less stable and less recommended.

A Biased Variant

DQN without stop gradient

For k = 1, 2, ...

Randomly pick an i (or a mini-batch) from \mathcal{B}

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta}(s_i', a') \right)^2$$

Residual gradient

cf. standard DQN

For k = 1, 2, ...Randomly pick an i (or a mini-batch) from \mathcal{B} $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\overline{\theta}}(s'_i, a') \right)^2$ $\overline{\theta} \leftarrow (1 - \tau)\overline{\theta} + \tau \theta$

For
$$k = 1, 2, ...$$

 $\theta \leftarrow \overline{\theta}$
For $m = 1, ..., M$:
Randomly pick an i (or a mini-batch) from \mathcal{B}
 $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\overline{\theta}}(s'_i, a') \right)^2$
 $\overline{\theta} \leftarrow \theta$

A Biased Variant

This variant will converge (as it is similar to standard SGD), but the solution it converges to could be undesirable.

The underlying loss function of this algorithm is

$$\text{Loss}^{\text{RG}}(\theta) = \sum_{i \in \mathcal{B}} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta}(s'_i, a') \right)^2$$

cf. the desired Bellman error that we truly want to minimize is

$$BE(\theta) = \sum_{i \in \mathcal{B}} \left(Q_{\theta}(s_i, a_i) - R(s_i, a_i) - \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta}(s', a') \right)^2$$

Recall the Bellman equation: $\forall s, a \quad Q(s, a) = R(s, a) + \gamma \sum_{s'} P(s'|s, a) \max_{a'} Q_{\theta}(s', a')$

Lemma 1. For random variables X, Y, let $\mu_X \triangleq \mathbb{E}_Y[Y|X]$. Then

$$\mathbb{E}[(X - Y)^2] = \mathbb{E}_X[(X - \mu_X)]^2 + \mathbb{E}_{X,Y}[(\mu_X - Y)^2].$$

Proof.

$$\mathbb{E}_{X,Y}[(X-Y)^2] = \mathbb{E}_{X,Y}[(X-\mu_X+\mu_X-Y)^2] = \mathbb{E}_X[(X-\mu_X)]^2 + 2\mathbb{E}_{X,Y}[(X-\mu_X)(\mu_X-Y)] + \mathbb{E}_{X,Y}[(\mu_X-Y)^2].$$

The second term above is zero because

$$\mathbb{E}_{X,Y}[(X - \mu_X)(\mu_X - Y)] = \mathbb{E}_X \left[\mathbb{E}_Y \left[(X - \mu_X)(\mu_X - Y) \mid X \right] \right]$$
$$= \mathbb{E}_X \left[(X - \mu_X) \mathbb{E}_Y \left[(\mu_X - Y) \mid X \right] \right]$$
$$= \mathbb{E}_X \left[(X - \mu_X) 0 \right]$$
$$= 0.$$

A Biased Variant

Assuming μ is the state-action distribution in \mathcal{B} , we have

$$\mathbb{E}[\operatorname{Loss}^{\operatorname{RG}}(\theta)] = \mathbb{E}_{(s,a)\sim\mu} \mathbb{E}_{r\sim R(s,a)} \mathbb{E}_{s'\sim P(\cdot|s,a)} \left[\left(Q_{\theta}(s,a) - r - \gamma \max_{a'} Q_{\theta}(s',a') \right)^{2} \right]$$

$$= \mathbb{E}_{(s,a)\sim\mu} \left[\left(Q_{\theta}(s,a) - R(s,a) - \gamma \mathbb{E}_{s'\sim P(\cdot|s,a)} \max_{a'} Q_{\theta}(s',a') \right)^{2} \right]$$

$$+ \mathbb{E}_{(s,a)\sim\mu} \mathbb{E}_{r\sim R(s,a)} \left[(r - R(s,a))^{2} \right]$$

$$+ \mathbb{E}_{(s,a)\sim\mu} \mathbb{E}_{\tilde{s}\sim P(\cdot|s,a)} \left[\left(\max_{a'} Q_{\theta}(\tilde{s},a') - \mathbb{E}_{s'\sim P(\cdot|s,a)} \left[\max_{a'} Q_{\theta}(s',a') \right] \right)^{2} \right]$$

$$(3)$$

where the last equality is by Lemma 1. Term (1) is equal to $BE(\theta)$, which is the loss we truly want to minimize. Term (3) is an undesired term that is non-zero when the transition is non-deterministic.

A Biased Variant

- Theoretically, if the transition is deterministic, residual gradient is optimizing a correct loss function without bias.
 - Still not recommended: for other reasons, it converges slower and has worse performance when there is no sufficient 1) data coverage or 2) powerful function approximation.
 - It deviates from the idea of Value Iteration / dynamic programming.

Fujimoto et al. Why should I trust you, bellman? the bellman error is a poor replacement for value error. 2022. Saleh and Jiang. Deterministic Bellman Residual Minimization

Summary

For k = 1, 2, ... $\theta \leftarrow \overline{\theta}$ For m = 1, ..., M: Randomly pick an *i* (or a mini-batch) from \mathcal{B} $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_{i}, a_{i}) - r_{i} - \gamma \max_{a'} Q_{\overline{\theta}}(s'_{i}, a') \right)^{2}$ $\overline{\theta} \leftarrow \theta$ $\overline{\theta} \leftarrow \theta$ $\overline{\theta}$ is updated at the same speed as θ

For k = 1, 2, ...Randomly pick an *i* (or a mini-batch) from \mathcal{B} $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\overline{\theta}}(s'_i, a') \right)^2$ $\overline{\theta} \leftarrow (1 - \tau)\overline{\theta} + \tau \theta$

DQN

Not taking gradient on the $Q_{\overline{\theta}}(s', a')$ term $\overline{\theta}$ is updated slower than θ

For
$$k = 1, 2, ...$$

Randomly pick an *i* (or a mini-batch) from \mathcal{B} $\theta \leftarrow \theta - \alpha \nabla_{\theta} \left(Q_{\theta}(s_i, a_i) - r_i - \gamma \max_{a'} Q_{\theta}(s'_i, a') \right)^2$

Taking gradient on the $Q_{\overline{\theta}}(s', a')$ term

Not Recommended