

GQ(λ) Quick Reference Guide

Adam White and Richard S. Sutton

July 7, 2014

This document should serve as a quick reference for the linear GQ(λ) off-policy learning algorithm. We refer the reader to Maei and Sutton (2010) for a more detailed explanation of the intuition behind the algorithm and convergence proofs. If you have questions or concerns about the content in this document or the attached java code please email adam.white@ualberta.ca.

1 Requirements and Setting

For each use of GQ(λ) you need to provide the following three *question functions*. (In the following \mathcal{S} and \mathcal{A} denote the sets of states and actions.)

- $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; target policy to be learned. Incidentally, if π is chosen as the greedy policy with respect to the learned value function, then the algorithm will implement a generalization of the Greedy-GQ algorithm (Maei, Szepesvari, Bhatnagar & Sutton 2010).
- $\gamma : \mathcal{S} \rightarrow [0, 1]$; termination function ($\gamma(s) = 1 - \beta(s)$ in GQ paper)
- $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$; reward function

In the earlier publications there was also specified a fourth question function, the terminal reward function $z : \mathcal{S} \rightarrow \mathbb{R}$ used to specify a final reward at termination. Since that time it has been recognized that this functionality can be included in the reward function, making use of the termination function (Modayil, White & Sutton 2014). For example, if one wanted only a terminal reward function $z(s)$ upon termination in state s , one would use a reward function of $r(s, a, s') = (1 - \gamma(s'))z(s')$. This completes the specification of the predictive question that you are seeking to answer using the GQ(λ) algorithm.

The specific approximate answer found by $\text{GQ}(\lambda)$ will depend upon the following four *answers functions* (these also must be provided):

- $b : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; behavior policy
- $I : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$; interest function (can set to 1 for all state-action pairs or indicate selected state-action pairs to be best approximated)
- $\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^n$; feature-vector function
- $\lambda : \mathcal{S} \rightarrow [0, 1]$; bootstrapping or eligibility-trace decay-rate function

The following data structures are internal to GQ :

- $\theta \in \mathbb{R}^n$; the learned weights of the linear approximation: $Q^\pi(s, a) = \theta^\top \phi(s, a) = \sum_{i=1}^n \theta_i \phi_i(s, a)$
- $w \in \mathbb{R}^n$; secondary set of learned weights
- $e \in \mathbb{R}^n$; eligibility trace vector

Parameters internal to GQ :

- α ; step-size parameter for learning θ
- $\eta \in [0, 1]$; relative step-size parameter for learning w ($\alpha\eta$)

2 Algorithm Specification

We can now specify $\text{GQ}(\lambda)$. Let w and e be initialized to zero and θ be initialized arbitrarily. Let the subscript t denote the current time step. Let ρ_t denote the ‘‘importance sampling’’ ratio:

$$\rho_t = \frac{\pi(S_t, A_t)}{b(S_t, A_t)}, \quad (1)$$

where $S_t \in \mathcal{S}$ and $A_t \in \mathcal{A}$ are the state and action occurring on time step t . Let $\bar{\phi}_t$ denote the expected next feature vector:

$$\bar{\phi}_t = \sum_{a \in \mathcal{A}} \pi(S_t, a) \phi(S_t, a) \quad (2)$$

The following equations fully specify $\text{GQ}(\lambda)$:

$$\delta_t = r(S_t, A_t, S_{t+1}) + \gamma(S_{t+1}) \theta_t^\top \bar{\phi}_{t+1} - \theta_t^\top \phi(S_t, A_t) \quad (3)$$

$$\theta_{t+1} = \theta_t + \alpha \left[\delta_t e_t - \gamma(S_{t+1})(1 - \lambda(S_{t+1}))(\mathbf{w}_t^\top e_t) \bar{\phi}_{t+1} \right] \quad (4)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha \eta [\delta_t e_t - (\mathbf{w}_t^\top \phi(S_t, A_t)) \phi(S_t, A_t)] \quad (5)$$

$$e_t = I(S_t) \phi(S_t, A_t) + \gamma(S_t) \lambda(S_t) \rho_t e_{t-1} \quad (6)$$

3 Pseudocode

The following pseudocode characterizes the algorithm and its use.

```
Initialize  $\theta$  arbitrarily and  $w = 0$ 
Repeat (for each episode):
  Initialize  $e = 0$ 
   $S \leftarrow$  initial state of episode
  Repeat (for each step of episode):
     $A \leftarrow$  action selected by policy  $b$  in state  $S$ 
    Take action  $A$ , observe next state,  $S'$ 
     $\bar{\phi} \leftarrow 0$ 
    For all  $a \in \mathcal{A}(s)$ :
       $\bar{\phi} \leftarrow \bar{\phi} + \pi(S', a)\phi(S', a)$ 
       $\rho = \frac{\pi(S,A)}{b(S,A)}$ 
      GQlearn( $\phi(S, A), \bar{\phi}, \lambda(S'), \gamma(S'), r(S, A, S'), \rho, I(S)$ )
     $S \leftarrow S'$ 
  until  $S'$  is terminal
```

```
GQLearn( $\phi, \bar{\phi}, \lambda, \gamma, R, \rho, I$ )
 $\delta \leftarrow R + \gamma\theta^\top \bar{\phi} - \theta^\top \phi$ 
 $e \leftarrow \rho e + I\phi$ 
 $\theta \leftarrow \theta + \alpha(\delta e - \gamma(1 - \lambda)(w^\top e)\bar{\phi})$ 
 $w \leftarrow w + \alpha\eta(\delta e - (w^\top \phi)\phi)$ 
 $e \leftarrow \gamma\lambda e$ 
```

4 Code

The files `GQlambda.java` and `GQlambda.cpp` contain implementations of the GQlearn function described in the pseudocode. We have excluded optimizations (e.g., binary features or efficient trace implementation) to ensure the code is simple and easy to understand. We leave it to the reader to provide environment code for interfacing to GQ(λ) (e.g., using RL-Glue).

5 References

- Maei, H. R., Szepesvari, Cs., Bhatnagar, S., Sutton, R. S. (2010). Toward Off-Policy Learning Control with Function Approximation. In *Proceedings of the 27th International Conference on Machine Learning*, Haifa, Israel.
- Maei, H. R. and Sutton, R. S. (2010). GQ(λ): A general gradient algorithm for temporal-difference prediction learning with eligibility traces. In *Proceedings of the Third Conference on Artificial General Intelligence 1*: 91–96.
- Modayil, J., White, A., Sutton, R. S. (2014). Multi-timescale nexting in a reinforcement learning robot. *Adaptive Behavior* 22(2):146–160.
- Sutton, R. S., Barto, A. G. (1998). Reinforcement Learning: An Introduction. MIT Press.