
UNDERSTANDING BELLMAN EQUATIONS

Shuo Liu

Computer Science

Northeastern University

shuo.liu2@northeastern.edu

ABSTRACT

This article covers Bellman Equation, its RL specializations, and generalizations.

1 BELLMAN EQUATION AND OPERATOR

1.1 DEFINITION

The Bellman Equation and optimal Bellman Equation for V-values are, Sutton & Barto (2018),¹

$$\begin{aligned} V^\pi(s) &\doteq \mathbb{E}_{a \sim \pi(\cdot|s)} [Q^\pi(s, a)], \\ &= \mathbb{E}_{a \sim \pi(\cdot|s)} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')]], \\ V^*(s) &\doteq \max_a [Q^*(s, a)], \\ &= \max_a [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^*(s')]]. \end{aligned} \tag{1}$$

and the Bellman Equation and optimal Bellman Equation for Q-values are,

$$\begin{aligned} Q^\pi(s, a) &\doteq R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')], \\ &= R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [\mathbb{E}_{a' \sim \pi(a'|s')} Q^\pi(s', a')], \\ Q^*(s, a) &\doteq R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^*(s')], \\ &= R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [\max_{a'} Q^*(s', a')]. \end{aligned} \tag{2}$$

where $V^\pi(s)$ and $Q^\pi(s, a)$ are value representations following policy π , e.g., vectors and functions.

$$\tilde{\pi}(s) \doteq \operatorname{argmax}_a Q^\pi(s, a). \tag{3}$$

Bellman Equations establish relations between states and succeeding states, which can be applied as updating rules for value prediction. A succinct representation is to define the Bellman Equation as a unary mathematical operator. The V-value Bellman and optimal Bellman Operators are,

$$\begin{aligned} (\mathcal{T}^\pi \circ V^\pi)(s) &\doteq \mathbb{E}_{a \sim \pi(\cdot|s)} [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')]], \\ (\mathcal{T}^* \circ V^\pi)(s) &\doteq \max_a [R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V^\pi(s')]]. \end{aligned} \tag{4}$$

The Bellman and optimal Bellman Operators \mathcal{T}^π for Q-values are,

$$\begin{aligned} (\mathcal{T}^\pi \circ Q^\pi)(s, a) &\doteq R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [\mathbb{E}_{a' \sim \pi(a'|s')} Q^\pi(s', a')], \\ (\mathcal{T}^* \circ Q^\pi)(s, a) &\doteq R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [\max_{a'} Q^\pi(s', a')]. \end{aligned} \tag{5}$$

Curse of Dimension Why do we mostly use MDP (where the future evolution is independent of its history) and hence Bellman Equations to model RL problems? Bellman (1957) coined the “curse of dimension”, which describes the exponential increase in the state space size as dimensionality grows, making calculations extremely complex. Breaking this curse often requires altering the problem or its constraints, though complete solutions are not always achievable.

For convenience, we use Q-value as the representative in the following parts of this article.

¹The Bellman Equations shown in the main article are stochastic, the deterministic one can be derived as, $V(s) = \max_a \{R(s, a) + \gamma V(s')\}$, where $s' \leftarrow T(s, a)$ is a deterministic succeeding state of s following a .

1.2 IMPORTANT PROPERTIES

Proposition 1 (γ -contraction). *Given any $Q, Q' \mapsto \mathbb{R}^{|S| \times |A|}$, Bellman Operators are γ -contraction Operators in L^∞ norm,*

$$\begin{aligned} \|\mathcal{T}^\pi \circ Q - \mathcal{T}^\pi \circ Q'\|_\infty &\leq \gamma \|Q - Q'\|_\infty, \\ \text{and } \|\mathcal{T}^* \circ Q - \mathcal{T}^* \circ Q'\|_\infty &\leq \gamma \|Q - Q'\|_\infty. \end{aligned} \quad (6)$$

Corollary 1 (Fixed-point Iteration). *For any $Q^0 \mapsto \mathbb{R}^{|S| \times |A|}$, after $k \rightarrow \infty$ iterations of Bellman transformation, $Q^{\pi, \infty} \doteq \lim_{k \rightarrow \infty} (\mathcal{T}^\pi)^k \circ Q^0$, or $Q^{*, \infty} \doteq \lim_{k \rightarrow \infty} (\mathcal{T}^*)^k \circ Q^0$, according to Banach's Fixed Point Theorem,*

$$Q^{\pi, \infty} = Q^{*, \infty} = Q^*, \quad \text{which **uniquely** satisfies } \mathcal{T}^\pi \circ Q^* = Q^*, \text{ or } \mathcal{T}^* \circ Q^* = Q^*. \quad (7)$$

Theorem 1 (Fundamental theorem). *Any memoryless policy that is greedy to Q^* (**deterministically** maximizes) is optimal Szepesvári (2010),*

$$\tilde{\pi}^* \doteq \operatorname{argmax}_a Q^* = \pi^*. \quad (8)$$

Proposition 2 (Monotone). *Bellman Operators are monotonic. For any Q -values $Q, Q' \mapsto \mathbb{R}^{|S| \times |A|}$,*

$$\begin{aligned} (Q \leq Q') &\Leftrightarrow (\mathcal{T}^\pi \circ Q \leq \mathcal{T}^\pi \circ Q'), \\ (Q \leq Q') &\Leftrightarrow (\mathcal{T}^* \circ Q \leq \mathcal{T}^* \circ Q'). \end{aligned} \quad (9)$$

2 BELLMAN BACKUP FOR PLANNING

2.1 DYNAMIC PROGRAMMING

According to the Fundamental Theorem, we can find π^* efficiently once having access to Q^* , without the need to find the policy whose Q-function **dominates** the others' brute-force-ly. To avoid the Bellman Curse of Dimensionality, we can apply Dynamic Programming (DP) methods to solve MDPs by keeping track of Q-values during calculations, thanks to Bellman recursions.

Value Iteration Value iteration (so-called backward induction) involves iteratively applying \mathcal{T}^* to arbitrarily initialized values Q^0 until convergence. According to Corollary 1 and Theorem 1, value iteration converges to Q^* as $k \rightarrow \infty$, then an optimal policy π^* can be derived by greedifying Q^* .

Policy Iteration Policy iteration starts with an arbitrary policy π^0 and values Q^0 . In each iterative step k , $Q^{\pi, k}$ is calculated by applying Bellman Operator $\mathcal{T}^{\pi, k}$ that follows current policy π^k to $Q^{\pi, k-1}$ from the last iteration, and then π^{k+1} is derived from greedifying $Q^{\pi, k}$. This process is repeated until convergence, and policy iteration can produce optimal policy after sufficient iterations.

3 BELLMAN RESIDUAL FOR LEARNING

3.1 TD-LEARNING WITH LOOK-UP TABLE

When the transition model is unavailable (model-free), we can use the residuals (RHS minus LHS) of the Bellman Equations as learning objective,

$$\begin{aligned} (\mathcal{B}^\pi \circ Q)(s, a) &\doteq r + \gamma Q(s', \pi(s')) - Q(s, a), \\ (\mathcal{B}^* \circ Q)(s, a) &\doteq r + \gamma \max_{a'} Q(s', a') - Q(s, a). \end{aligned} \quad (10)$$

Assuming that our sampling and parameter updating roughly follow the true state distribution $\mu(s)$, the expectation of Bellman residual will be closed to zero at the optima. This approach is often called temporal difference (TD) learning.

TD-learning In TD-learning with learning rate α , the update rule for Q-values is,

$$Q(s, a) \leftarrow Q(s, a) + \alpha(B^\pi \circ Q)(s, a). \quad (11)$$

According to Stochastic Approximation Theorem, let k be the visitation times of state-action pair, and learning rates $0 \leq \alpha^k < 1$ satisfies $\forall(s, a), \sum_{k=1}^{\infty} \alpha^k(s, a) = \infty, \sum_{k=1}^{\infty} [\alpha^k(s, a)]^2 < \infty$. Following TD-learning updates, $Q^{\pi, k}(s, a)$ converges to $Q^*(s, a)$ as $k \rightarrow \infty$ (Jaakkola et al. (1994)).

Q-learning In Q-learning that relies on optimal Bellman Equation, the Q-value update is,

$$Q(s, a) \leftarrow Q(s, a) + \alpha(B^* \circ Q)(s, a). \quad (12)$$

According to Stochastic Approximation Theorem, let k be the visitation times of state-action pair, and learning rates $0 \leq \alpha^k < 1$ satisfies $\forall(s, a), \sum_{k=1}^{\infty} \alpha^k(s, a) = \infty, \sum_{k=1}^{\infty} [\alpha^k(s, a)]^2 < \infty$. Following Q-learning updates, $Q^{*, k}(s, a)$ converges to $Q^*(s, a)$ as $k \rightarrow \infty$ (Watkins & Dayan (1992)). The deep version of Q-learning algorithm, Deep Q-Network (DQN), is shown in Appendix.

However, the nice property of convergence only holds in the tabular case and cannot be extended to a function approximation as discussed later.

3.2 TD-LEARNING WITH FUNCTION APPROXIMATION

To introduce generalization to the value function, we represent the approximated Q-value in a parameterized functional form. Our goal is to minimize the mean squared value error,

$$\mathcal{L}(\theta) = \frac{1}{2} \sum_{s \in \mathcal{S}} \mu(s) \left[Q^{\text{target}} - Q_\theta(s, a) \right]^2, \quad (13)$$

where Q^{target} is the ground truth and Q_θ is the prediction. Just like TD-learning, the Bellman residual can be applied for the value function approximation.

Semi gradient for Bellman Residual Similar to stochastic gradient methods with unbiased target estimators, if we use the Bellman Equation to get target Q-value Q^{target} , but here we just ignore its potential gradient change, the gradient ascent for Bellman residual is,

$$\begin{aligned} \Delta_{\text{semi}} \theta &= -\frac{1}{2} \alpha \nabla_\theta \left[Q^{\text{target}} - Q_\theta(s, a) \right]^2 \\ &= \alpha \left[Q^{\text{target}} - Q_\theta(s, a) \right] \nabla_\theta Q_\theta(s, a), \text{ where } Q^{\text{target}} = r + \gamma Q_\theta(s', a') \end{aligned} \quad (14)$$

Since we neglects a part of the gradient of Q^{target} , it is called semi gradient for Bellman residual (θ in red). Though semi-gradient methods are fast and simple, they could have divergence issue, e.g., Baird's counter-example (the star problem).

Full Gradient for Bellman Residual The full Bellman residual gradient should include all gradient components, including the gradient of the target estimation,

$$\begin{aligned} \Delta_{\text{full}} \theta &= -\frac{1}{2} \alpha \nabla_\theta \left[r + \gamma Q_\theta(s', a') - Q_\theta(s, a) \right]^2 \\ &= -\alpha \left[r + \gamma Q_\theta(s', a') - Q_\theta(s, a) \right] \left[\gamma \nabla_\theta Q_\theta(s', a') - \nabla_\theta Q_\theta(s, a) \right]. \end{aligned} \quad (15)$$

If the approximation system is general enough and the value functions are continuous, the full Bellman residual gradient is guaranteed to converge to the optima. However, this is at the sacrifice of learning speed, as illustrated by the hall problem.

Hybrid Gradient for Bellman Residual In contrast to Figure 1a where Δ_{semi} boosts Δ_{full} , Figure 1b represents the case where the semi gradient may diverge. Baird (1995) combined these 2 methods: to keep stable, Δ_B should stay in the same direction as Δ_{full} (above the perpendicular axis); meanwhile, Δ_B should stay as close as possible to Δ_{semi} to increase learning speed.

$$\begin{aligned} \Delta_B \theta &= (1 - \omega) \cdot \Delta_{\text{semi}} \theta + \omega \cdot \Delta_{\text{full}} \theta, \\ &= -\alpha \left[r + \gamma Q_\theta(s', a') - Q_\theta(s, a) \right] \left[\omega \gamma \nabla_\theta Q_\theta(s', a') - \nabla_\theta Q_\theta(s, a) \right], \\ \text{s.t., } \Delta_B \theta \cdot \Delta_{\text{full}} \theta &\geq 0 \Leftrightarrow \omega \geq \frac{\Delta_{\text{semi}} \theta \cdot \Delta_{\text{full}} \theta}{\Delta_{\text{semi}} \theta \cdot \Delta_{\text{full}} \theta - \Delta_{\text{full}} \theta \cdot \Delta_{\text{full}} \theta}. \end{aligned} \quad (16)$$

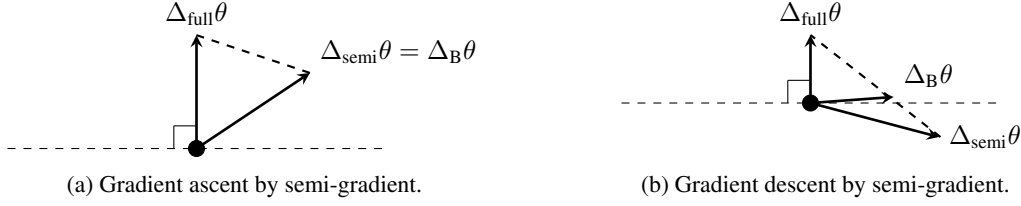


Figure 1: Epoch-wise gradient vectors for Bellman residual gradients.

4 BELLMAN GENERALIZATIONS

4.1 SOFT BELLMAN EQUATION

The Soft Bellman Equation incorporates an entropy regularization term, encouraging exploration by penalizing overly deterministic policies. It is useful where balancing exploration and exploitation is critical, such as environments with sparse rewards or high uncertainty.

$$Q(s, a) = R(s, a) + \gamma \mathbb{E}_{a' \sim \pi(\cdot|s'), s' \sim P(\cdot|s, a)} [Q(s', a') - \lambda \log \pi(\cdot|s')], \quad (17)$$

where λ controls the weight of the entropy term $\mathcal{H}(\cdot|s)$.

4.2 CONTINUOUS-TIME BELLMAN EQUATION

Hamilton-Jacobi-Bellman (HJB) Equation is the continuous-time analogue of the Bellman Equation and used in problems involving continuous state and action spaces. Given system dynamics $f(s, a)$,

$$\frac{\partial Q(s, a)}{\partial t} + \max_a [R(s, a) + \nabla_s Q(s, a)^T f(s, a)] = 0. \quad (18)$$

4.3 DISTRIBUTIONAL BELLMAN EQUATION

The Distributional Bellman Equation models the return as a distribution \mathcal{Q} , which is particularly useful in risk-sensitive scenarios with specific demands of quantifying the variability of returns,

$$\mathcal{Q}(s, a) \stackrel{\text{D}}{=} R(s, a) + \gamma \mathcal{Q}(s', a'), \quad (19)$$

where \mathcal{Q} is distributional expected returns, and $\stackrel{\text{D}}{=}$ denotes equality in distribution.

4.4 MULTI-AGENT BELLMAN EQUATION

The Multi-Agent (MA) Bellman Equation generalizes the Bellman Equation to cooperative MA systems, where multiple agents collaborate to achieve a common goal.

Joint Bellman Equation The Joint MA Bellman Equation models the group of agents as an entity, considering their collective actions $\mathbf{a} = (a_1, \dots, a_N)$ and expected returns Q ,

$$Q(s, \mathbf{a}) = R(s, \mathbf{a}) + \gamma \mathbb{E}_{\mathbf{a}' \sim \pi(\cdot|s'), s' \sim P(\cdot|s, \mathbf{a})} [Q(s', \mathbf{a}')]. \quad (20)$$

Individual Bellman Equation The Individual Bellman Equation computes Q_i for each agent, assuming that agents' execution and rewards are independent (but observing a shared reward $R(s, \mathbf{a})$),

$$Q_i(s, a_i) = R(s, \mathbf{a}) + \gamma \mathbb{E}_{a'_i \sim \pi_i(\cdot|s'), s' \sim P(\cdot|s, a_i)} [Q_i(s', a'_i)]. \quad (21)$$

4.5 MULTI-TASK BELLMAN EQUATION

The Generalized Bellman Equation extends the Bellman Equation to include additional objectives. It is particularly useful in multi-task learning, constrained RL or hierarchical RL with additional objectives $\phi(s, a)$ or constraints $\phi(s, a) \leftarrow \lambda c(s, a)$ alongside the primary goal,

$$Q(s, a) = R(s, a) + \phi(s, a) + \gamma \mathbb{E}_{a' \sim \pi(\cdot|s'), s' \sim P(\cdot|s, a)} [Q(s', a')]. \quad (22)$$

REFERENCES

- Leemon C. Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pp. 30–37. Elsevier, 1995.
- Richard Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, 1957.
- Thomas Jaakkola, Michael I. Jordan, and Satinder P. Singh. On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation*, 6(6):1185–1201, 1994. doi: 10.1162/neco.1994.6.6.1185.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2nd edition, 2018.
- Csaba Szepesvári. *Algorithms for Reinforcement Learning*, volume 4 of *Synthesis Lectures on Artificial Intelligence and Machine Learning*. Morgan & Claypool Publishers, 2010. ISBN 9781608454921. doi: 10.2200/S00268ED1V01Y201005AIM009.
- Christopher J. C. H. Watkins and Peter Dayan. Q-learning. *Machine Learning*, 8(3–4):279–292, 1992. doi: 10.1007/BF00992698.

A DEEP Q-NETWORK

Deep Q-Network (DQN) approximates the Q-function with a target network to stabilize training.

Algorithm 1 Deep Q-learning with Experience Replay

```
1: Initialize: parameter  $\theta$  for Q-network  $Q_\theta$ , target Q-network  $\theta' \leftarrow \theta$ , replay memory  $\mathcal{D}$ 
2: for each episode do
3:   Choose  $s_0$  from  $\mathcal{S}$  as initial state
4:   for  $t = 0$  to  $T - 1$  do
5:     Take action  $a_t$  according to policy derived by  $Q$  (e.g.  $\epsilon$ -greedy)
6:     Observe reward  $r_{t+1}$  and next state  $s_{t+1}$ 
7:     Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  into  $\mathcal{D}$ 
8:     for each mini-batch of  $N$  transitions  $\{s_i, a_i, r_{i+1}, s_{i+1}\}$  do
9:        $\delta_i \leftarrow r_{i+1} + \underbrace{\gamma \max_{a'} Q_{\theta'}(s_{i+1}, a') - Q_\theta(s_i, a_i)}_{\text{target Q prediction}}$ 
10:       $\theta \leftarrow \theta + \frac{1}{N} \sum_i \alpha \delta_i \nabla_\theta Q_\theta(s_i, a_i)$ 
11:    end for
12:     $\theta' \leftarrow \tau \theta' + (1 - \tau) \theta$ 
13:  end for
14: end for
15: return value function  $Q$ 
```
