
Faster Reinforcement Learning with Value Target Lower Bounding

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We show that an arbitrary lower bound of the maximum achievable value can be
2 used to improve the Bellman value target during value learning. In the tabular case,
3 value learning using the lower bounded Bellman operator converges to the same
4 optimal value as using the original Bellman operator, at a potentially faster speed.
5 In practice, discounted episodic return in episodic tasks and n-step bootstrapped
6 return in continuing tasks can serve as lower bounds to improve the value target.
7 We experiment on Atari games, FetchEnv tasks and a challenging physically
8 simulated car push and reach task. We see large gains in sample efficiency as
9 well as converged performance over common baselines such as TD3, SAC and
10 Hindsight Experience Replay (HER) in most tasks, and observe a reliable and
11 competitive performance against the stronger n-step methods such as td-lambda,
12 Retrace and optimality tightening. Prior works have already successfully applied a
13 special case of lower bounding (using episodic return), but are limited to a small
14 number of episodic tasks. To the best of our knowledge, we are the first to propose
15 the general method of value target lower bounding (with possibly bootstrapped
16 return), to demonstrate its optimality in theory, and effectiveness in a wide range
17 of tasks over many strong baselines.

18 1 Introduction

19 The value function is a key concept in dynamic programming approaches to Reinforcement Learning
20 (RL) (Bellman, 1957). It estimates the sum of all future rewards (usually time-discounted) of a given
21 state. In temporal difference (TD) learning, the value function is adjusted toward its Bellman target
22 which adds the reward of the current step with the (discounted) value of the next state (Sutton &
23 Barto, 2018). This forms the basis of many state of the art RL algorithms such as DQN (Mnih et al.,
24 2013), DDPG (Lillicrap et al., 2016), TD3 (Fujimoto et al., 2018), and SAC (Haarnoja et al., 2018).

25 The value of the next state is typically estimated using a “bootstrapped value” based on the value
26 function itself, which is being actively learned during training. The bootstrapped values can be
27 random and far from the optimal value, especially at the initial stage of training, or with sparse reward
28 tasks where rewards can only be achieved through a long sequence of actions. Consequently, the
29 Bellman value targets as well as the learned values are usually far away from the optimal value (the
30 value of the optimal policy).

31 Naturally, this leads to the following idea: If we can make the value target closer to the optimal value,
32 we may speedup TD learning. For example, we know that the optimal value is just the expected
33 discounted return of the optimal policy, which always upper bounds the expected return of any policy.
34 For episodic RL tasks, we could use the observed discounted return up to episode end from the
35 training trajectories to lower bound the value target. This makes the new value target closer to the
36 optimal value, when the empirical return is higher than the Bellman target.

Algorithm 1 Value iteration with value target lower bounding

Input: Finite MDP $p(s', r|s, a)$, convergence threshold θ , a lower bound $f(s)$ of the maximum achievable value $\bar{G}^v(s)$

Output: State value $v(s)$

$v(s) \leftarrow 0$

repeat

$\Delta \leftarrow 0$

$v_p(s) \leftarrow v(s)$

for each state s do

$\hat{v}(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a)[r + \gamma v_p(s')]$

$\hat{v}_f(s) \leftarrow \max(f(s), \hat{v}(s))$

$v(s) \leftarrow \hat{v}_f(s)$

$\Delta \leftarrow \max(\Delta, |v(s) - v_p(s)|)$

end for

until $\Delta < \theta$

37 The case for continuing or non-episodic tasks is less clear though. When a continuing task can
38 return negative rewards, any safe lower bound of the optimal value can be too low to be useful. One
39 could take the risk and use n-step bootstrapped return as a lower bound, which is unsafe because
40 bootstrapped return can overestimate and be greater than the optimal value. Can we still use them as
41 lower bounds to improve TD value targets?

42 2 Theoretical Results for the Tabular Case

43 Our results show that for the tabular case, arbitrary functions below a certain bootstrap bound can be
44 used to lower bound the value target to still converge to the same optimal value.

45 2.1 Background

46 In finite MDPs with a limited number of states and actions, a table can keep track of the value of
47 each state. Using dynamic programming algorithms such as value iteration, values are guaranteed to
48 converge to the optimum through Bellman updates (Chapter 4.4 (Sutton & Barto, 2018)).

49 The core of the value iteration algorithm (Algorithm 1) is the Bellman update of the value function,
50 $\mathcal{B}(v)$, where $v(s')$ is the bootstrapped value:

$$\mathcal{B}(v)(s) := \max_a \sum_{s', r} p(s', r|s, a)[r + \gamma v(s')] \quad (1)$$

51 It is well known that the Bellman operator, \mathcal{B} , is a contraction mapping over value functions (Denardo,
52 1967). That is, for any two value functions v_1 and v_2 , $\|\mathcal{B}(v_1) - \mathcal{B}(v_2)\|_\infty \leq \gamma \|v_1 - v_2\|_\infty$ for the
53 discount factor $\gamma \in [0, 1)$ and $\|x\|_\infty := \max_i |x_i|$ (the L_∞ norm). This guarantees that any value
54 function under the algorithm converges to the optimal value $\mathcal{B}^\infty(v) = v^*$.¹

55 2.2 Convergence of value target lower bounding

56 **Definition 2.1.** The expected n-step bootstrapped return for a given policy π and value function $v(s)$
57 is defined as the expected bootstrapped return of taking n steps according to policy π :

$$G_n^{\pi, v}(s_0) := \mathbb{E}^\pi \{r_1 + \dots + \gamma^{n-1} r_n + \gamma^n v(s_n)\} \quad (2)$$

58 Here, the step rewards r_i and the resulting n-th step state s_n are random variables, with the expectation
59 \mathbb{E}^π taken over all possible n-step trajectories under the policy π and the given MDP.

¹For the gist of the proof, see for example page 8 of <https://people.eecs.berkeley.edu/~pabbee1/cs287-fa09/lecture-notes/lecture5-2pp.pdf>

60 **Definition 2.2.** Given the current learned value function $v(s)$, policy class Π , the *maximum achievable*
 61 *value* of a state s is defined as:

$$\bar{G}^v(s) := \max_{\pi \in \Pi, n \in [1, +\infty)} G_n^{\pi, v}(s) \quad (3)$$

62 This is a more relaxed definition of maximum because for each state s , a different policy $\pi(s)$ and a
 63 different number of steps $n(s)$ can be used to achieve the maximum $\bar{G}^v(s)$. And the theorem below
 64 says any function not exceeding the maximum achievable value can be used to lower bound the value
 65 target, and still achieve the optimal value in convergence.

66 **Theorem 2.3.** *Under the same assumptions for Bellman value contraction, for any function f that*
 67 *lower bounds the maximum achievable value, i.e. $\forall s, f(s) \leq \bar{G}^v(s)$, if we define the lower bounded*
 68 *Bellman operator as $\mathcal{B}_f(v) := \max(\mathcal{B}(v), f)$, then $\mathcal{B}_f^\infty(v) = \mathcal{B}^\infty(v)$.*

69 Note, the value $v(s)$ and the bootstrapped value can be inaccurate, and even above the optimal value.
 70 As a consequence, when n is finite, the maximum achievable value $\bar{G}^v(s)$ (and f) can be above the
 71 maximum expected return (i.e. the optimal value). On the other hand, when n is sufficiently large,
 72 the effect of the bootstrap value $v(s_n)$ diminishes (see Equation 2), and the maximum achievable
 73 value becomes the maximum expected return (i.e. the optimal value). Therefore, $\forall s, \bar{G}^v(s)$ is no
 74 smaller than the optimal value $\mathcal{B}^\infty(v)(s)$. As a special case of the theorem, as long as f is below the
 75 optimal value, value target lower bounding converges correctly:

76 **Corollary 2.4.** *If function f lower bounds the optimal value, i.e. $\forall s, f(s) \leq \mathcal{B}^\infty(v)(s)$, then*
 77 *$\mathcal{B}_f^\infty(v) = \mathcal{B}^\infty(v)$.*

78 A few things to note about the proof of Theorem 2.3 (included in Appendix 1.1).

79 First, this only proves convergence, not contraction under the original $\|v_1 - v_2\|_\infty$ metric. In the
 80 case of the Bellman operator, contraction shows that $\forall v_1, v_2$ value functions, $\|\mathcal{B}(v_1) - \mathcal{B}(v_2)\|_\infty \leq$
 81 $\gamma \|v_1 - v_2\|_\infty$. Here, for value target lower bounding, what’s proved is convergence to v^* at a rate
 82 of γ , not contraction. There can be counter examples where the distance between v_1 and v_2 under
 83 one application of \mathcal{B}_f can increase in the original L_∞ metric space, even though v_1 and v_2 are both
 84 getting closer to v^* at a rate of γ . One difficulty caused by convergence instead of contraction is that
 85 the stopping criterion in Algorithm 1 ($\Delta < \theta$) no longer works, due to the inaccessible v^* during
 86 learning. In practice, this may not be a serious concern, as people often train algorithms for a fixed
 87 number of iterations or time steps.

88 Second, based on the proof, the new algorithm is at least as fast as the original. When the lower
 89 bound actually improves the value target, i.e. $f(s) > \mathcal{B}(v_1)(s)$, there is a chance for the convergence
 90 to be faster. Convergence is strictly faster when the lower bound f has an impact on the L_∞ distance
 91 between the current value and the optimal value, i.e. it increases the value target for the states where
 92 the differences between the current value and the optimal value are the largest.

93 Third, the lower bound function doesn’t have to be static during training. As long as there is a single
 94 f during each training update, convergence is preserved.

95 The following sections detail how to compute lower bounds of the maximum achievable value
 96 (Section 3), how to integrate the lower bounds into state of the art RL algorithms (Section 4), and
 97 provide an illustration of how this method may benefit value learning in practice (Section 4.3).

98 3 Example Lower Bound Functions

99 We show a few cases where lower bound functions can be readily obtained from the training
 100 experience. Future work may investigate alternatives.

101 3.1 Episodic tasks

102 In episodic tasks, discounted return is accumulated up to the last step of an episode. In this case, we
 103 can wait until an episode ends, and compute future discounted returns of all time steps up to the end
 104 of the episode. This episodic return is a lower bound of the optimal value when the environment is

105 deterministic, because the reward sequence can be repeated using the same sequence of actions². To
 106 make training efficient, we can compute and store such discounted returns into the replay buffer for
 107 each time step, and simply read them out during training, which adds very little computation to the
 108 baseline one-step TD computation.

$$f(s_0) = \sum_{i=0, \dots, \infty} \gamma^i r(s_i, a_i) \quad (4)$$

109 We call this variant “lb-DR”, short for lower bounding with discounted return.

110 3.1.1 Episodic with hindsight relabeled goals

111 In goal conditioned tasks, one helpful technique is hindsight goal relabeling (Andrychowicz et al.,
 112 2017). It takes a future state that is d time steps away from the current state as the hindsight / relabeled
 113 goal for the current state. When the goal is reached, a reward of 0 is given, otherwise a -1 reward is
 114 given for each time step.

115 In this case, we know it took d steps to reach the hindsight goal, so the discounted future return is:

$$\begin{aligned} f(s_0) &= \sum_{i=0, \dots, d-1} -1\gamma^i \\ &= -1(1 - \gamma^d)/(1 - \gamma) \end{aligned} \quad (5)$$

116 This calculation can be done on the fly as hindsight relabeling happens, requiring no extra space and
 117 very little computation.

118 We call this variant “lb-GD”, short for lower bounding with goal distance based return.

119 Additionally, we can also apply lb-DR and lb-GD together, with discounted episodic return (lb-DR)
 120 on the original experience and goal distance based return (lb-GD) on the hindsight experience, giving
 121 the “lb-DR+GD” variant, which was used in Fujita et al. (2020).

122 3.2 In general (including non-episodic tasks)

123 If the task is continuing, without an episode end³, discounted return needs to be accumulated all the
 124 way to infinity. When rewards are always non-negative, one can still use the accumulated discounted
 125 reward of the future n -steps to lower bound the value. But accumulated n -step discounted reward is no
 126 longer a lower bound when rewards can be negative, in which case, the more general lower bounding
 127 with bootstrapped value can be used: given a trajectory of training experience $\tau := \langle s_0, \dots, s_n \rangle$:

$$G_n^v(\tau) := r_1 + \gamma r_2 + \dots + \gamma^{n-1} r_n + \gamma^n v(s_n) \quad (6)$$

128 Assuming the rewards and the state s_n can be repeated with the same action sequence, $G_n^v(\tau)$ lower
 129 bounds the maximum achievable value $\bar{G}^v(s_0)$ (Equation 3).

130 Two variations are possible: Given a trajectory of length n ,

131 1. compute $v(s_i)$ for all $i \in [1, n]$ and take the maximum of all $G_i^v(\tau)$ to obtain a tighter lower
 132 bound. We call this variant “lb-b- n step”:

$$f(s_0) = \max_{i \in [1, n]} G_i^v(\tau) \quad (7)$$

133 2. only evaluate v on the last (n th) step and use the n th-step bootstrapped return as the lower
 134 bound, which involves less compute but results in a looser bound. (When n is large enough,
 135 this becomes the lb-DR variant.) We call this variant “lb-b- n step-only”.

$$f(s_0) = G_n^v(\tau) \quad (8)$$

²Note that the behavior policy can be stochastic, as long as the policy class contains the optimal policy, value learning will converge to the optimal.

³Chapter 3.3 of Sutton & Barto (2018) has more details on episodic vs continuing tasks.

136 4 Integration into RL algorithms

137 4.1 Background

138 The value target lower bounds can be readily plugged into RL algorithms that regress value to a
139 target, e.g. DQN, DDPG or SAC.

140 In these algorithms, the action value $q(s, a)$ is learned through a squared loss with the target value y .
141 In one step TD return, for a batch \mathbf{B} of experience $\{s, a \rightarrow r, s'\}$, the loss is:

$$\mathcal{L}_q := \sum_{(s,a,r,s') \in \mathbf{B}} |q(s, a) - y|^2 \quad (9)$$

142 In one step TD return, y is the one step TD return $\hat{q}(s, a, r, s')$:

$$\hat{q}(s, a, r, s') := r(s, a) + \gamma q'(s', \mu'(s')) \quad (10)$$

143 Here, q' and μ' are the bootstrap value and policy functions, typically following the value and policy
144 functions in a delayed schedule during training. (They are also called “target value” and “target
145 policy”, and are very different from the “value target” y in this paper.)

146 4.2 Value target lower bounding

147 With lower bounding, we replace the value target y with the lower bounded target:

$$y \leftarrow \max(f, \hat{q}(s, a, r, s')) = \max(f, r + \gamma q'(s', \mu'(s'))) \quad (11)$$

148 This way of lower bounding the value target is the same as was done by Fujita et al. (2020) (confirmed
149 via personal communication), but is subtly and importantly different from lower bounding the q
150 value directly (Oh et al., 2018; Tang, 2020): $q(s, a) \leftarrow \max(f, q(s, a))$, which stays overestimated if
151 $q(s, a)$ initially overestimates.

152 To the best of our knowledge, value target lower bounding with bootstrapped values is a novel
153 contribution of this work.

154 4.3 An Illustrative Example

155 Figure 1 includes a fairly general example showing how value target lower bounding would improve
156 value learning. Suppose we enhance an off policy algorithm such as DDPG with value target lower
157 bounding (lb-DR), when there is no training experience hitting the target state, no meaningful training
158 happens for the baseline or lb-DR. However, when there is one trajectory hitting the target state, all
159 states along the trajectory will soon be propagated with meaningful return, and nearby states will also
160 enjoy faster learning. As the state space becomes larger and the time horizon longer, a successful
161 trajectory will speed up learning quite a bit.

162 5 Experiments

163 The goal is to demonstrate the sample efficiency of lower bounding the value target over baseline such
164 as DDPG, TD3, SAC and HER. Because the lower bounded value target can now look potentially
165 many steps into the future, we suspect it to be best suited for long horizon, sparse reward tasks.
166 Hence, we choose to experiment on a sampled subset of Atari games, the goal conditioned FetchEnv
167 tasks and the harder goal conditioned Pioneer Push and Reach tasks. See details of the experiment
168 setup in Appendix 1.2.

169 5.1 Baselines

170 Baselines include DDPG (Lillicrap et al., 2016), TD3 (Fujimoto et al., 2018), SAC (Haarnoja et al.,
171 2018) and HER (Andrychowicz et al., 2017; Plappert et al., 2018). Implementations are based on

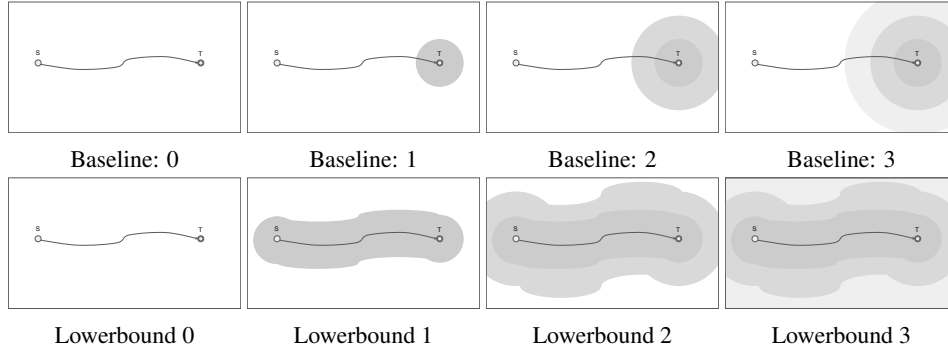


Figure 1: Illustration of value target lower bounding speeding up value learning as training progresses from stages 0 to 3. The task is to navigate in the state space from start state S to end state T , with sparse reward 1 at T and 0 elsewhere. The curve from S to T denotes a training experience that reaches the target. The shaded areas denote roughly states whose value has been significantly improved during training up to that stage.

172 open sourced repositories, and baseline performance is verified against published results under similar
 173 settings. The Appendix 1.6 and 1.5 include results on more baselines such as DDQN (van Hasselt
 174 et al., 2015), td-lambda (Sutton & Barto, 2018) and Retrace (Munos et al., 2016).

175 5.2 Hyperparameters

176 Value target lower bounding is applied on top of these baselines without any additional hyperparameter
 177 (Section 4). The only hyperparameters come from the baselines. These hyperparameters follow
 178 published work as much as possible. When baseline hyperparameters need to be tuned for an
 179 environment, e.g. Atari games or Pioneer tasks, we search for the best performance in total episode
 180 reward averaged across all tasks for that environment on one set of random seeds, then the optimal
 181 hyperparameters are fixed and evaluated on a separate set of random seeds never seen during
 182 development. Value target lower bounding simply uses the the parameter values optimal for the
 183 baselines. Details are in Appendix 1.3.

184 5.3 Results

185 We report results on both episodic and continuing/non-episodic tasks. We report evaluation per-
 186 formance averaged across several runs of the algorithms (five for the less stable Atari games and
 187 three for the others). Each run uses a random seed never seen during development. Due to space
 188 constraints, the main paper only reports performance aggregated across all tasks for each environment.
 189 During each run, we take one task and one random seed, run baseline and treatment algorithms, and
 190 record whether treatment agent evaluates strictly above the baseline agent as training progresses. We
 191 average across all the runs of the same environment, and plot the fraction of times where treatment is
 192 above baseline and the standard deviation of that fraction in Figure 2. Appendix 1.4 contains per task
 193 evaluation curves.

194 Overall, value target lower bounding is a simple, effective, efficient, carefree (no hyperparameter)
 195 and theoretically justified approach. Although the example lower bounds are limited to deterministic
 196 environments, the theory is generally applicable to stochastic environments. A similar prior work
 197 to compare would be Hindsight Experience Replay (HER) (Andrychowicz et al., 2017), which is
 198 simple, effective, efficient, and also limited to deterministic environments (Blier & Ollivier, 2021).
 199 However, unlike our work, HER relies on the task being goal conditioned with full knowledge of the
 200 reward function, has one hyperparameter to tune (the proportion of hindsight experience), and is not
 201 justified in theory for stochastic environments. Our work shows further significant gains on top of
 202 HER on hard continuous control tasks.

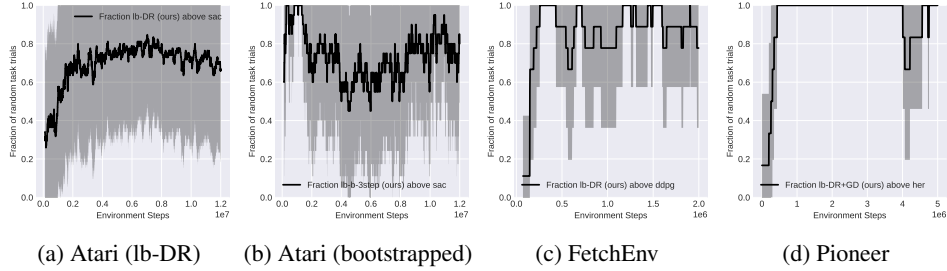


Figure 2: Aggregated evaluation performance: The fraction of times where treatment performs strictly above baseline, plotted along the number of time steps used for training. The solid curve is the sample mean of the fraction across all runs, and the shaded area is \pm one standard deviation. We use, for Atari (lb-DR), 85 runs – 17 games each with 5 seeds, for Atari (bootstrapped), 20 runs – 4 games 5 seeds, for FetchEnv, 9 runs – 3 tasks 3 seeds, and for Pioneer, 6 runs – 2 tasks 3 seeds.

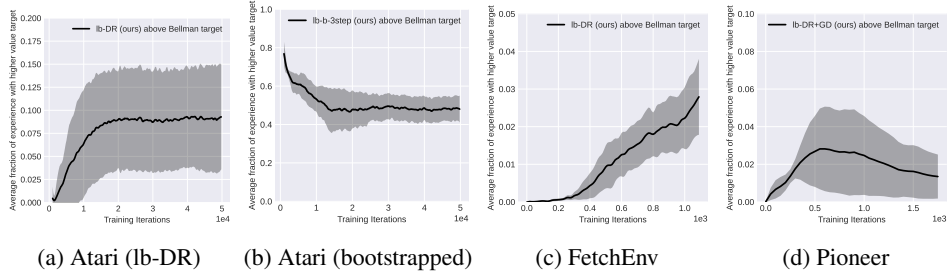


Figure 3: Average fraction of training experience where lower bounding improves Bellman value target, plotted along the number of iterations of training. The solid curve is the average of the fraction across all runs – number of tasks times number of seeds, and the shaded area is \pm one standard deviation. Other setups are the same as Figure 2.

203 **5.3.1 Lower bounding vs baselines SAC/DDPG/HER**

204 Figure 2 compares the lower bounding treatment with SAC, DDPG or HER baseline on 17 sampled
 205 Atari games, the FetchEnv tasks and the Pioneer tasks. For all the environments, value target lower
 206 bounding is not only more sample efficient, but also enjoys a higher converged performance. After
 207 training starts, it quickly gains ground and outperforms the baseline for 70% to 100% of the runs. It
 208 keeps that advantage even at the end of the training, outperforming baseline in converged performance.
 209 These plots show how frequently treatment is above baseline, but is insensitive to the magnitude of
 210 change.

211 Appendix 1.4 shows the magnitude of change with total episode return plotted for each task, often
 212 with large gains in sample efficiency and sometimes much higher converged performance. Among all
 213 the 22 tasks, only one task (Atari Breakout) shows lb-DR underperforming the baseline.

214 Investigations show that the loss on Atari Breakout is likely due to the mismatch between training
 215 objective and evaluation metric. During training, raw rewards are clipped to $[-1, 1]$ and step discounted
 216 at $\gamma = 0.99$ to compute value, while in evaluation, total reward is the unclipped and undiscounted
 217 cumulative sum of episode rewards. The discount and clipping together severely penalizes large
 218 rewards earned later in the episode, which is what’s happening for Breakout, because hitting a top
 219 layer block produces a reward of 7 while hitting a bottom layer block produces 1. When we use
 220 non-clipped rewards or a higher γ in training, the lower bounding method performs much better in
 221 total reward. Note, this train-test discrepancy as well as an additional training bias (Thomas, 2014) is
 222 likely present in all of the prior works using policy gradient methods on Atari games.

223 **5.3.2 Value target improvement**

224 The lb-DR method is mostly effective, but is it really due to improvements to the value targets?
 225 Figure 3 looks at the fraction of training experience where lower bounding actually improves the

226 Bellman value target over the course of training. Overall, improved value target roughly coincides
227 with performance gains. Appendix 1.4 shows the plots per each task.

228 The Appendix also has more results, comparing with baselines such as n-step methods, DDQN and
229 optimality tightening, and more analyses such as ablations and robustness to hyperparameter choices.

230 6 Related Work

231 Prior works (Fujita et al., 2020; Hoppe & Toussaint, 2020; He et al., 2017; Oh et al., 2018; Tang,
232 2020) employed several different ways of computing future returns and using that as a lower bound
233 to improve value learning. It is quite easy to introduce biases and inefficiencies into the process and
234 end up with a suboptimal or inefficient algorithm. Our work is the first to propose the general form of
235 value target lower bounding (possibly with bootstrapping), to show its convergence to the optimal
236 value in the tabular case, and to demonstrate its effectiveness in illustrative examples and extensive
237 experiments on a wide range of tasks.

238 Fujita et al. (2020)’s method is similar to a special case (the lb-DR+GD variant) of the general
239 method. They used it as a part of a large system and showed that it improved sample efficiency for a
240 robotic grasping task. Hoppe & Toussaint (2020) also bounded the value target. But instead of using
241 empirical return, they used a simplified MDP with a subset of actions. Although without theoretical
242 proof and only experimented on a limited set of robotic manipulation tasks, both works show that
243 value target lower bounding increased sample efficiency. This work, in addition to the theory and the
244 more general method, shows that lower bounding improves both sample efficiency and converged
245 performance in a wide range of tasks.

246 He et al. (2017) used empirical return with bootstrap to improve value learning. They formulated value
247 learning as a constrained optimization problem with the empirical bootstrapped value being the lower
248 (and upper) constraints of the value function. In their experiments, the Lagrangian multiplier was
249 fixed, which would likely lead to suboptimal solutions. Our lb-b-nstep method also uses bootstrapped
250 value. But we lower bound the value target directly, which is simpler, more efficient, and likely more
251 optimal. Our work points out that for episodic tasks, even more efficient and effective methods like
252 lb-DR exist. Appendix 1.5 offers more discussion and results related to this.

253 Our work is subtly but importantly different from the prior works on lower bound Q learning or Self
254 Imitation Learning (SIL) (Oh et al., 2018; Tang, 2020). SIL uses empirical return R to lower bound
255 the value function itself (instead of the *value target*). This is achieved by adding an off policy value
256 loss during on-policy (AC or PPO) training ($L_{value}^{sil} = \frac{1}{2}|v(s) - \max(v(s), R)|^2$). When the value
257 function overestimates, the SIL value loss becomes zero, and keeps overestimating. Mixing the SIL
258 loss with the loss from the baseline algorithms probably helped to correct the overestimation, but no
259 theoretical guarantee was given. In evaluation, SIL was often compared to on-policy Actor Critic or
260 PPO baselines, so it was not clear how much of the gain was due to lower bounding and how much
261 due to off-policy value learning. In this work, we bound the Bellman value target (Equation 11), so
262 overestimates are automatically corrected via Bellman updates, and convergence is guaranteed in the
263 tabular case. We also use off-policy algorithms as baselines for a cleaner comparison.

264 N-step return methods such as td-lambda (Sutton & Barto, 2018) and Retrace (Munos et al., 2016)
265 also look a few steps ahead, but to obtain more accurate value of the behavior policy. Traditionally,
266 this requires careful off-policy correction, and the value can still be far from the optimal value due
267 to the often suboptimal behavior. This work shows that value target lower bounding efficiently and
268 effectively looks ahead much further without the need for off-policy correction, due to aiming at the
269 optimal value. Appendix 1.5 has more detailed observations and discussions.

270 Planning methods can look into the future to achieve higher value targets and better control. Examples
271 include Monte Carlo Tree Search (MCTS) (Schrittwieser et al., 2019; Ye et al., 2021) and Model
272 Predictive Control (MPC) or receding horizon planning with raw actions (Chua et al., 2018; Hafner
273 et al., 2019; Zhang et al., 2022), options (Silver & Ciosek, 2012), or subgoals (Nasiriany et al.,
274 2019; Nair & Finn, 2020; Chane-Sane et al., 2021). Planning methods use either a dynamics model
275 together with the learned value or just the learned value (in the case of goal conditioned tasks)
276 (Nasiriany et al., 2019) to improve policy or value estimates. Planning typically happens during roll
277 out (Nasiriany et al., 2019), but can also be used to improve the value target, as in Reanalyze of
278 MuZero (Schrittwieser et al., 2019; Ye et al., 2021). During value improvement, if planning takes

279 the maximum over a set of possible future values (e.g. from different trajectories as in the case
280 of MPC), and if this set includes the one step Bellman value target, then the planner is essentially
281 using alternative trajectories and their values to lower bound the Bellman value target. In this sense,
282 the theory developed here can potentially justify and improve Reanalyze. In general, planning is
283 orthogonal to value target lower bounding, and typically requires additional components and a lot
284 more compute than the basic TD learning does. Therefore, we leave it to future work to explore the
285 synergy between the two.

286 Interestingly, it is common practice to lower and upper bound the returns to the possible region,
287 e.g. Andrychowicz et al. (2017) bounds value between $[-\frac{1}{1-\gamma}, 0]$. Similar to lower bounding with
288 episodic return (Section 3.1), such strict bounds of the actual value can be thought of as admissible
289 heuristics (bounds) used during search of the optimal solution (Russell & Norvig, 2020). What’s new
290 in this work is that lower bounding with bootstrapped values (which can overestimate the value) still
291 converges to the optimal value.

292 Kumar et al. (2020) (DisCor) also recognized that bootstrapped value targets can be inaccurate. This
293 bias impacts learning adversely under function approximation. DisCor uses distribution correction to
294 sample experience with accurate bootstrap targets more frequently, while value target lower bounding
295 aims to directly reduce the bias.

296 While in theory using empirical return to lower bound the value target is only correct for deterministic
297 environments, in practice, it seems as long as the environment is not heavily impacted by random fluc-
298 tuations, they still perform well. In fact, with function approximation, the agents cannot distinguish
299 between two slightly different states, making the problem partially observable (Sutton & Barto, 2018)
300 and appear slightly random. Prior methods such as SIL (Oh et al., 2018), Optimality Tightening (He
301 et al., 2017), and even Hindsight relabeling (Andrychowicz et al., 2017) and MuZero (Schrittwieser
302 et al., 2019) require the environment to be deterministic. Despite this theoretical limitation, the lower
303 bounding methods and the prior methods can still be very useful, outperforming baselines often by
304 large margins and when deploying to the real world (Fujita et al., 2020).

305 7 Conclusions

306 We propose a general form of lower bounding the value target using possibly bootstrapped return. In
307 theory, value target lower bounding converges to the same optimal solution as the original Bellman
308 operator. In practice, several ways of finding value lower bounds are examined.

309 For episodic tasks, discounted episodic return is an efficient and effective method involving very
310 little extra computation. Precomputing the episodic return and storing it into the replay buffer
311 allows efficient lower bound computation. It achieves much higher sample efficiency and converged
312 performance than one-step baselines such as SAC, DDPG or TD3 in most tasks, and is competitive
313 among n-step baselines. Simple goal distance based return uses even less compute and achieves large
314 gains in certain long horizon tasks over Hindsight relabeling (HER).

315 For non-episodic tasks or in general, lower bounding with n-step bootstrapped return outperforms
316 one-step baselines and is a strong competitor to the n-step methods such as (truncated) td-lambda and
317 Retrace.

318 7.1 Future Work

319 There are probably better ways of finding value lower bounds that improve training even more. One
320 direction may be to use planning (e.g. Monte Carlo Tree Search, the Cross Entropy Method or using
321 subgoals) to achieve tighter lower bounds given a model of the task.

322 Estimating value lower bound for stochastic tasks may be possible, e.g. by learning a reward
323 function and a dynamics model and using imagined rollouts to obtain bootstrapped returns without
324 overestimation.

325 Other ways of bounding the value target, e.g. upper bounding (He et al., 2017), may be worth
326 investigating as well, e.g. to reduce overestimation in regions of poor reward.

327 **References**

- 328 Andrychowicz, M., Crow, D., Ray, A., Schneider, J., Fong, R., Welinder, P., McGrew, B.,
329 Tobin, J., Abbeel, P., and Zaremba, W. Hindsight experience replay. In Guyon, I., von
330 Luxburg, U., Bengio, S., Wallach, H. M., Fergus, R., Vishwanathan, S. V. N., and Gar-
331 nett, R. (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference
332 on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA,
333 USA*, pp. 5048–5058, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/
334 453fadbd8a1a3af50a9df4df899537b5-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/453fadbd8a1a3af50a9df4df899537b5-Abstract.html).
- 335 Bellman, R. *Dynamic Programming*. Princeton Univ. Press, Princeton, NJ, USA, 1957. ISBN
336 0-691-07951-X.
- 337 Blier, L. and Ollivier, Y. Unbiased methods for multi-goal reinforcement learning. *CoRR*,
338 abs/2106.08863, 2021. URL <https://arxiv.org/abs/2106.08863>.
- 339 Chane-Sane, E., Schmid, C., and Laptev, I. Goal-conditioned reinforcement learning with imagined
340 subgoals. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on
341 Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 1430–1440.
342 PMLR, 18–24 Jul 2021. URL [https://proceedings.mlr.press/v139/chane-sane21a.
343 html](https://proceedings.mlr.press/v139/chane-sane21a.html).
- 344 Chua, K., Calandra, R., McAllister, R., and Levine, S. Deep reinforcement learning in a handful of
345 trials using probabilistic dynamics models. In *Proceedings of the 32nd International Conference
346 on Neural Information Processing Systems, NIPS’18*, pp. 4759–4770, Red Hook, NY, USA, 2018.
347 Curran Associates Inc.
- 348 Denardo, E. V. Contraction mappings in the theory underlying dynamic programming. *SIAM Review*,
349 9(2):165–177, 1967. ISSN 00361445. URL <http://www.jstor.org/stable/2027440>.
- 350 Fujimoto, S., van Hoof, H., and Meger, D. Addressing function approximation error in actor-critic
351 methods. *CoRR*, 2018. URL <http://arxiv.org/abs/1802.09477>.
- 352 Fujita, Y., Uenishi, K., Ummadisingu, A., Nagarajan, P., Masuda, S., and Castro, M. Distributed rein-
353 forcement learning of targeted grasping with active vision for mobile manipulators. In *IEEE/RSJ
354 International Conference on Intelligent Robots and Systems (IROS)*, pp. 9712–9719, Oct 2020.
- 355 Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy
356 deep reinforcement learning with a stochastic actor. In Dy, J. and Krause, A. (eds.), *Proceedings of
357 the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine
358 Learning Research*, pp. 1861–1870. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.
359 press/v80/haarnoja18b.html](https://proceedings.mlr.press/v80/haarnoja18b.html).
- 360 Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., and Davidson, J. Learning latent
361 dynamics for planning from pixels. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of
362 the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine
363 Learning Research*, pp. 2555–2565. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.
364 press/v97/hafner19a.html](https://proceedings.mlr.press/v97/hafner19a.html).
- 365 He, F. S., Liu, Y., Schwing, A. G., and Peng, J. Learning to play in a day: Faster deep reinforcement
366 learning by optimality tightening. In *5th International Conference on Learning Representations,
367 ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net,
368 2017. URL <https://openreview.net/forum?id=rJ8Je4c1g>.
- 369 Hoppe, S. and Toussaint, M. Qgraph-bounded q-learning: Stabilizing model-free off-policy deep
370 reinforcement learning. *CoRR*, abs/2007.07582, 2020. URL [https://arxiv.org/abs/2007.
371 07582](https://arxiv.org/abs/2007.07582).
- 372 Kumar, A., Gupta, A., and Levine, S. Discor: Corrective feedback in reinforcement learning via
373 distribution correction. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin,
374 H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18560–18572.
375 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/
376 d7f426ccbc6db7e235c57958c21c5dfa-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/d7f426ccbc6db7e235c57958c21c5dfa-Paper.pdf).

377 Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D.
378 Continuous control with deep reinforcement learning. In Bengio, Y. and LeCun, Y. (eds.), *4th*
379 *International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May*
380 *2-4, 2016, Conference Track Proceedings*, 2016. URL <http://arxiv.org/abs/1509.02971>.

381 Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller,
382 M. A. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013. URL
383 <http://arxiv.org/abs/1312.5602>.

384 Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. G. Safe and efficient off-policy
385 reinforcement learning. In Lee, D. D., Sugiyama, M., von Luxburg, U., Guyon, I., and
386 Garnett, R. (eds.), *Advances in Neural Information Processing Systems 29: Annual Con-*
387 *ference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona,*
388 *Spain*, pp. 1046–1054, 2016. URL [https://proceedings.neurips.cc/paper/2016/hash/](https://proceedings.neurips.cc/paper/2016/hash/c3992e9a68c5ae12bd18488bc579b30d-Abstract.html)
389 [c3992e9a68c5ae12bd18488bc579b30d-Abstract.html](https://proceedings.neurips.cc/paper/2016/hash/c3992e9a68c5ae12bd18488bc579b30d-Abstract.html).

390 Nair, S. and Finn, C. Hierarchical foresight: Self-supervised learning of long-horizon tasks via
391 visual subgoal generation. In *International Conference on Learning Representations*, 2020. URL
392 <https://openreview.net/forum?id=H1gzR2VKDH>.

393 Nasiriany, S., Pong, V., Lin, S., and Levine, S. Planning with goal-conditioned policies. *Advances in*
394 *Neural Information Processing Systems*, 2019.

395 Oh, J., Guo, Y., Singh, S., and Lee, H. Self-imitation learning. *CoRR*, abs/1806.05635, 2018. URL
396 <http://arxiv.org/abs/1806.05635>.

397 Plappert, M., Andrychowicz, M., Ray, A., McGrew, B., Baker, B., Powell, G., Schneider, J., Tobin,
398 J., Chociej, M., Welinder, P., Kumar, V., and Zaremba, W. Multi-goal reinforcement learning:
399 Challenging robotics environments and request for research. *CoRR*, abs/1802.09464, 2018. URL
400 <http://arxiv.org/abs/1802.09464>.

401 Russell, S. J. and Norvig, P. *Artificial Intelligence: A Modern Approach (4th Edition)*. Pearson, 2020.
402 ISBN 9780134610993. URL <http://aima.cs.berkeley.edu/>.

403 Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart,
404 E., Hassabis, D., Graepel, T., Lillicrap, T. P., and Silver, D. Mastering atari, go, chess and shogi by
405 planning with a learned model. *CoRR*, abs/1911.08265, 2019. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1911.08265)
406 [1911.08265](http://arxiv.org/abs/1911.08265).

407 Silver, D. and Ciosek, K. Compositional Planning Using Optimal Option Models. In *Proceedings of*
408 *the 29th International Conference on Machine Learning*, pp. 165. icml.cc / Omnipress, 2012.

409 Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. A Bradford Book,
410 Cambridge, MA, USA, 2018. ISBN 0262039249.

411 Tang, Y. Self-imitation learning via generalized lower bound q-learning. *CoRR*, abs/2006.07442,
412 2020. URL <https://arxiv.org/abs/2006.07442>.

413 Thomas, P. Bias in natural actor-critic algorithms. In Xing, E. P. and Jebara, T. (eds.), *Proceedings of*
414 *the 31st International Conference on Machine Learning*, volume 32 of *Proceedings of Machine*
415 *Learning Research*, pp. 441–448, Beijing, China, 22–24 Jun 2014. PMLR. URL [https://](https://proceedings.mlr.press/v32/thomas14.html)
416 proceedings.mlr.press/v32/thomas14.html.

417 van Hasselt, H., Guez, A., and Silver, D. Deep reinforcement learning with double q-learning. *CoRR*,
418 abs/1509.06461, 2015. URL <http://arxiv.org/abs/1509.06461>.

419 Ye, W., Liu, S., Kurutach, T., Abbeel, P., and Gao, Y. Mastering atari games with limited data. *CoRR*,
420 abs/2111.00210, 2021. URL <https://arxiv.org/abs/2111.00210>.

421 Zhang, H., Xu, W., and Yu, H. Generative planning for temporally coordinated exploration in
422 reinforcement learning. In *International Conference on Learning Representations*, 2022. URL
423 <https://openreview.net/forum?id=YZHES8wIdE>.

424 **Checklist**

- 425 1. For all authors...
- 426 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
427 contributions and scope? [Yes]
- 428 (b) Did you describe the limitations of your work? [Yes] Regarding the implementations
429 being theoretically limited to deterministic MDPs, see Section 3.1 and the last paragraph
430 of Related Works (Section 6) for details, and Appendix 1.7 for an example. The
431 theorems, however, do not assume a deterministic MDP, and follow standard Bellman
432 value contraction assumptions. Experiments use function approximation on continuous
433 observations or action spaces, which already violate the theoretical assumptions, but
434 the lower bounding methods still show large gains over the baselines.
- 435 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Appendix
436 1.9
- 437 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
438 them? [Yes]
- 439 2. If you are including theoretical results...
- 440 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 441 (b) Did you include complete proofs of all theoretical results? [Yes] Appendix 1.1
- 442 3. If you ran experiments...
- 443 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
444 perimental results (either in the supplemental material or as a URL)? [No] Omitted
445 due to anonymity. Code is already in a github repository, and will be released upon
446 publication.
- 447 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
448 were chosen)? [Yes] Appendix 1.2 and 1.3.
- 449 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
450 ments multiple times)? [Yes] Done for all plots
- 451 (d) Did you include the total amount of compute and the type of resources used (e.g., type
452 of GPUs, internal cluster, or cloud provider)? [No] Treatments only cost a bit more
453 compute than the baselines.
- 454 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 455 (a) If your work uses existing assets, did you cite the creators? [Yes] Except for the open
456 source code repository, to keep anonymity.
- 457 (b) Did you mention the license of the assets? [No]
- 458 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 459 (d) Did you discuss whether and how consent was obtained from people whose data you’re
460 using/curating? [N/A]
- 461 (e) Did you discuss whether the data you are using/curating contains personally identifiable
462 information or offensive content? [N/A]
- 463 5. If you used crowdsourcing or conducted research with human subjects...
- 464 (a) Did you include the full text of instructions given to participants and screenshots, if
465 applicable? [N/A]
- 466 (b) Did you describe any potential participant risks, with links to Institutional Review
467 Board (IRB) approvals, if applicable? [N/A]
- 468 (c) Did you include the estimated hourly wage paid to participants and the total amount
469 spent on participant compensation? [N/A]