



SBEED (out of the Deadly Triad):

Convergent Reinforcement Learning with Nonlinear Function Approximation

Bo Dai^{1→2}, Albert Shaw¹, Lihong Li², Lin Xiao³, Niao He⁴,
Zhen Liu¹, Jianshu Chen⁵, Le Song¹

¹Gatech, ²Google Brain, ³Microsoft Research, ⁴UIUC, ⁵Tencent AI

Main Results

Background

- This talk: $\text{RL} \approx \text{solving Bellman equation based on data}$
- Approx. dynamic programming (DQN, ...) inherently **unstable**
- Remained challenging for decades — “**deadly triad**” (Sutton 15)

Contributions

1. Bellman equation reformulated as a saddle-point problem
2. First provably convergent ADP algorithm (**SBEED**) with general nonlinear function approximation
3. Empirical validation in simulated robotics tasks (Mujoco)

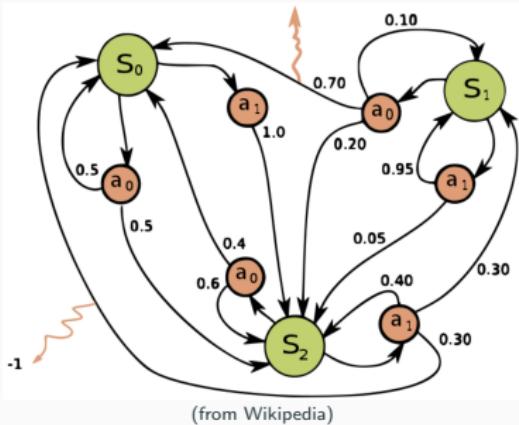
Outline

- **Background**
- Saddle-point Reformulation of Bellman Equation
- SBEED Learning
- Conclusions

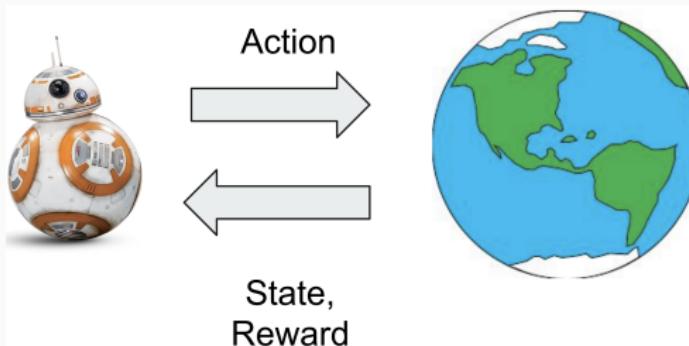
Markov Decision Process (MDP)

$$M = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$$

- Set of states \mathcal{S}
- Set of actions \mathcal{A}
- Transition probabilities $P(s'|s, a)$
- Immediate reward $R(s, a)$
- Discount factor $\gamma \in (0, 1)$



(from Wikipedia)



Bellman Equation and Dynamic Programming

Optimal value function V^* satisfies Bellman equation

$$\forall s \in \mathcal{S}, V^*(s) = \underbrace{\max_{a \in \mathcal{A}} (R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)}[V^*(s')])}_{\mathcal{T}V^*(s)}$$

Well-known facts of Bellman operator \mathcal{T} :

- \mathcal{T} is γ -contraction: $\|\mathcal{T}V_1 - \mathcal{T}V_2\|_\infty \leq \gamma \|V_1 - V_2\|_\infty$
- Hence, $V, \mathcal{T}V, \mathcal{T}^2V, \mathcal{T}^3V, \dots \rightarrow V^*$ ("fixed point")
- Mathematical foundation of value iteration, TD(λ), Q-learning, etc.
in the exact (\approx finite-MDP) case

When Bellman Meets Gauss: Approximate DP

In practice, V^* is often approximated

- Eg: least-squares fit on linear models or neural networks, ...
- Composing \mathcal{T} and approximation loses contraction
- Many known divergent examples
Baird (93), Boyan & Moore (95), Tsitsiklis & Van Roy (96), ...
- Limited positive theory or algorithms

Gordon (96), Tsitsiklis & Van Roy (97), Lagoudakis & Parr (03), Sutton et al. (08, 09), Maei et al. (10), ...

Functional Approximations and Dynamic Programming

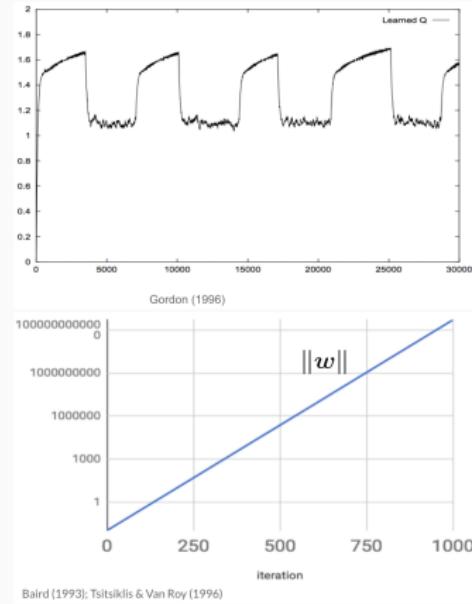
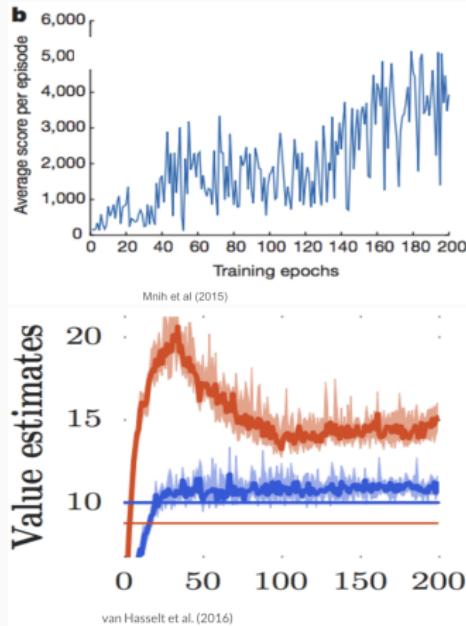
A major open problem for decades.

By Richard Bellman and Stuart Dreyfus

Math. Tables & Other Aids Comp. (1959)

Does It Matter in Practice?

Many empirical successes of (double, dueling) DQN, A3C, ...
in video games, Go, robotics, dialogue management, ...
but often with surprises:



The Deadly Triad (Controlled Case)

Existing RL algorithms risk divergence in the “**deadly triad**”:

- (nonlinear) function approximation
- bootstrapping
- off-policy learning

Outline

- Background
- **Saddle-point Reformulation of Bellman Equation**
- SBEED Learning
- Conclusions

Why Solving Bellman Equation is Hard?

A natural objective function for solving $V = \mathcal{T}V$:

$$\begin{aligned} & \min_V \underbrace{\|V - \mathcal{T}V\|^2}_{\text{"Bellman error/residual"}} \\ &= \min_V \mathbb{E}_s \left[(V(s) - \max_a (R(s, a) + \gamma \mathbb{E}_{s'|s,a}[V(s')])^2 \right] \end{aligned}$$

- **Difficulty #1:** breaks smoothness and continuity
- **Difficulty #2:** typical SGD gives **biased** gradient, known as “double sample” issue (Baird 95):

$$(\cdots + \gamma \mathbb{E}_{s'|s,a}[V_w(s')])^2 \neq \mathbb{E}_{s'|s,a} [(\cdots + \gamma V_w(s'))^2]$$

Addressing Difficulty #1: Nesterov Smoothing

$$\begin{aligned} V(s) &= \max_a (R(s, a) + \gamma \mathbb{E}_{s' | s, a}[V(s')]) \\ &\Downarrow \\ V(s) &= \underbrace{\max_{\pi(\cdot|s)} \sum_a \pi(a|s) (R(s, a) + \gamma \mathbb{E}_{s' | s, a}[V(s')])}_{:= \mathcal{T}_\lambda V(s)} \underbrace{+ \lambda H(\pi(\cdot|s))}_{\text{entropic regularization}} \end{aligned}$$

The smoothed Bellman operator \mathcal{T}_λ may be derived differently
Rawlik+ (12), Fox+ (16), Ne+ (17), Nachum+(17), Asadi & Littman (17), ...

Properties of \mathcal{T}_λ

$$\mathcal{T}_\lambda V(s) := \max_{\pi(\cdot|s)} \sum_a \pi(a|s) (R(s, a) + \gamma \mathbb{E}_{s'|s,a}[V(s')]) + \lambda H(\pi(\cdot|s))$$

- Still a γ -contraction
- Existence and uniqueness of fixed point V_λ^*
- Controlled bias: $\|V_\lambda^* - V^*\|_\infty = \mathcal{O}(\lambda/(1-\gamma))$
- **Temporal consistency** (as in PCL of Nachum+ (17))

$$\forall s, a : \quad V_\lambda^*(s) = R(s, a) + \gamma \mathbb{E}_{s'|s,a}[V_\lambda^*(s')] - \lambda \log \pi_\lambda^*(a|s)$$

Addressing Difficulty #2: Legendre-Fenchel Transformation

$$\min_V \mathbb{E}_s \left[(V(s) - \max_a (R(s, a) + \gamma \mathbb{E}_{s' | s, a} [\mathbf{V}(s')])^2 \right]$$

↓ (by Nesterov smoothing)

$$\min_{V, \pi} \mathbb{E}_{s, a} \left[\underbrace{\left(R(s, a) + \gamma \mathbb{E}_{s' | s, a} [\mathbf{V}(s')] - \lambda \log \pi(a | s) - V(s) \right)}_{\text{denoted } x_{sa}}^2 \right]$$

↓ (L-F transform: $x_{sa}^2 = \max_{y \in \mathbb{R}} (2x_{sa}y - y^2)$)

$$\min_{V, \pi} \max_{\nu \in \mathbb{R}^{S \times A}} \mathbb{E}_{s, a, s'} [(2\nu(s, a)x_{sa} - \nu(s, a)^2)]$$

The last step also applies the interchangeability principle
(Rockafellar & Wets 88; Shapiro & Dentcheva 14; Dai+ 17)

Reformulation of Bellman Equation

We have now turned a **fixed point** into a **saddle point**:

$$\min_{V, \pi} \max_{\nu} \mathbb{E}_{s, a, s'} \left[2\nu(s, a) \left(R(s, a) + \gamma V(s') - \lambda \log \pi(a|s) - V(s) \right) - \nu(s, a)^2 \right]$$

- Well-defined objective without requiring double samples
- May be optimized by gradient methods (SGD/BackProp, ...)
- See paper for a slightly more general version
- Special cases: GTD2, PCL, ...
- Inner maximum achieved when ν equals λ -smoothed Bellman error

Outline

- Background
- Saddle-point Reformulation of Bellman Equation
- **SBEED Learning**
- Conclusions

SBED: Smoothed Bellman Error Embedding

Algorithmic ideas

- Parameterize $(V, \pi; \nu)$ by $(w_V, w_\pi; w_\nu)$
- Stochastic-gradient updates on (w_V, w_π) and ascent on w_ν
 - Two-time-scale updates for primal and dual variables; or
 - Exact maximization if concave in w_ν
- Our implementation uses stochastic mirror descent

Algorithmic advantages

- Agnostic to whether data is on- or off-policy
- May be used in batch (e.g., experience replay) or online mode
- Extensible to multi-step and eligibility traces cases
- Efficiently implemented (only first-order updates needed)

SBED Convergence (Batch Case)

Define $\bar{\ell}(V, \pi) := \max_{\nu} L(V, \pi, \nu)$, and assume

- $\nabla \bar{\ell}$ is Lipschitz-continuous
- the stochastic gradient has finite variance
- stepsizes are properly set

Theorem. SBED solution satisfies $\mathbb{E}[\|\nabla \bar{\ell}(V_{\hat{w}}, \pi_{\hat{w}})\|] \rightarrow 0$

- Building on results of Ghadimi & Lan (13)
- See paper for variants of convergence results ...
- ... and statistical/generalization analyses

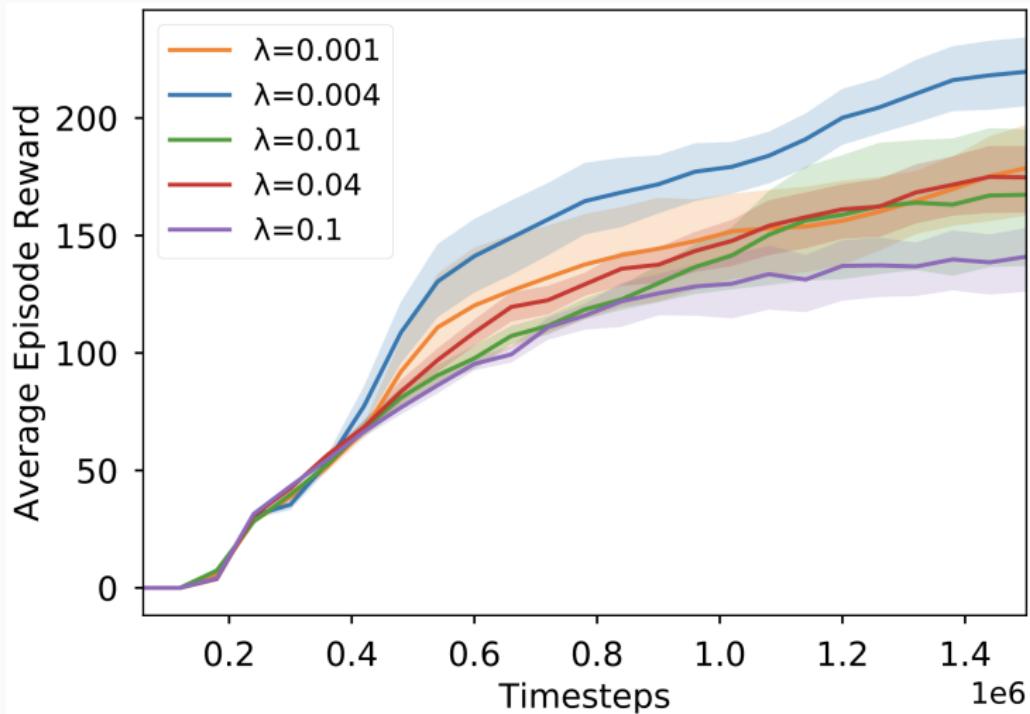
Experiments

- Use Mujoco on OpenAI as benchmark
- Compare to state-of-the-art baselines:
 - Dual-AC (Dai et al. 18)
 - TRPO (Schulman et al. 15)
 - DDPG (Lillicrap et al. 15)

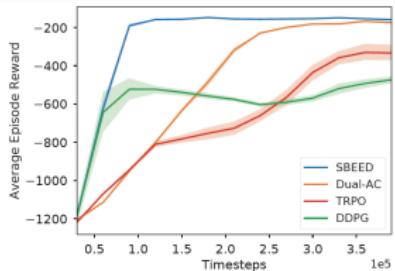


(from <http://www.mujoco.org>)

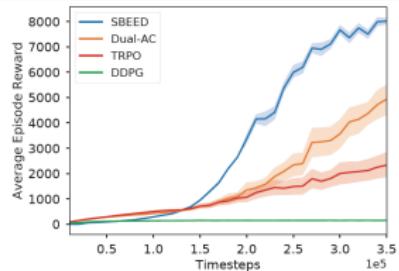
Role of Smoothing Parameter λ



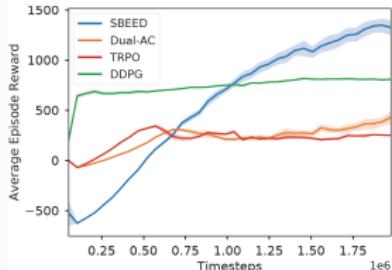
Comparison against Baselines



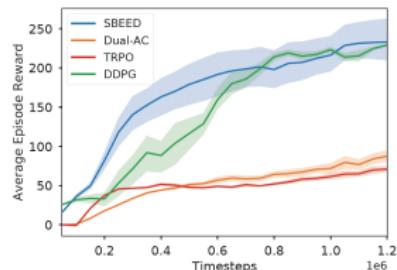
(a) Pendulum



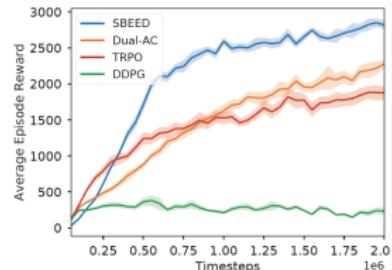
(b) InvertedDoublePendulum



(c) HalfCheetah



(d) Swimmer



(e) Hopper

Outline

- Background
- Saddle-point Reformulation of Bellman Equation
- SBEED Learning
- **Conclusions**

Conclusions

Contributions

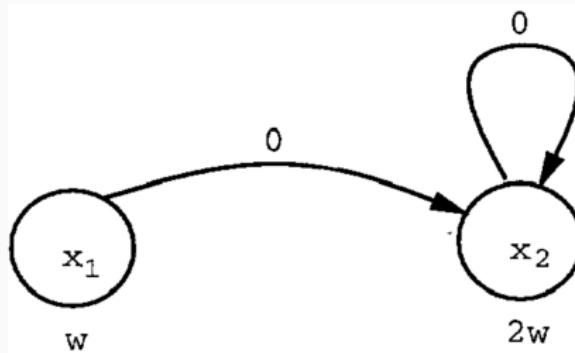
- A saddle-point reformulation of Bellman equation
- New algorithm SBEED with guaranteed convergence
- Promising empirical results on standard benchmark

Further remarks

- Many directions for future research
- Much efforts of finding true gradient RL algorithms
GTD (Sutton et al. 08), GTD2 (Sutton et al. 09), ...
- Deep connection to optimization
Mahadevan et al. (14), Macua et al. (15), ...
- Stronger algorithms based on new optimization techniques
Liu et al. (15, 16), Dai et al. (17, 18), Du et al. (17), Wang (17), **Chen et al. (18)**, ...

APPENDIX

Divergence Example of Tsitsiklis & Van Roy (96)



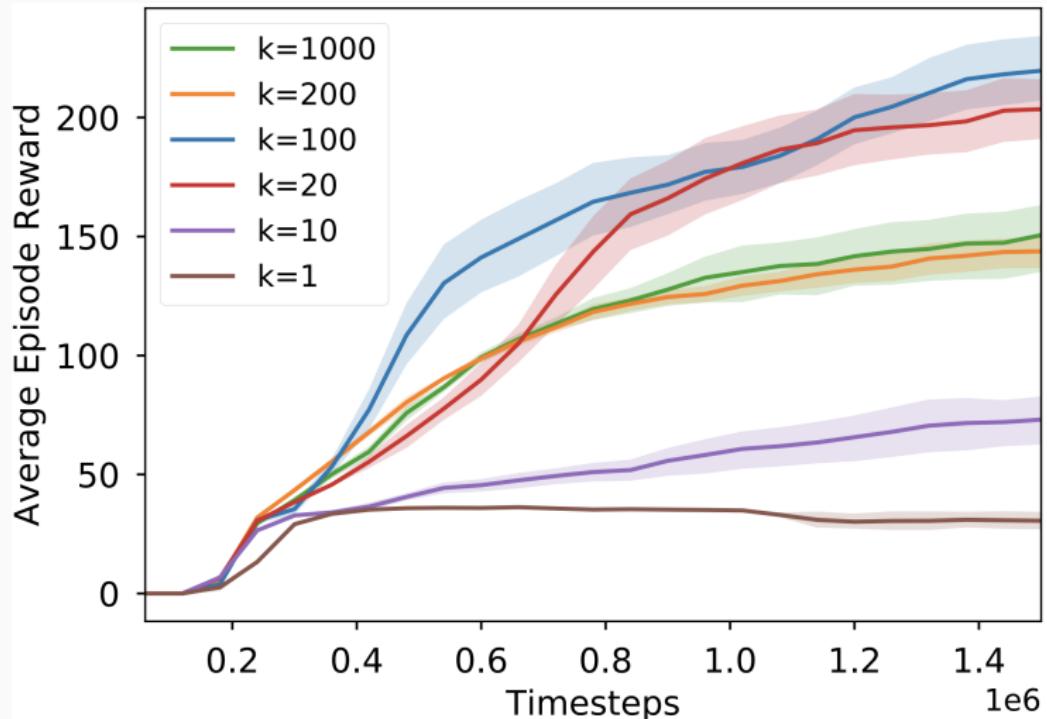
Starting with $w^{(0)} \neq 0$,
least-squares value iteration diverges when $\gamma > 5/6$,
although V^* may be exactly represented (with $w^* = 0$).

Online SBEED Learning with Experience Replay

Algorithm 1 Online SBEED learning with experience replay

- 1: Initialize $w = (w_V, w_\pi, w_\rho)$ and π_b randomly, set ϵ .
 - 2: **for** episode $i = 1, \dots, T$ **do**
 - 3: **for** size $k = 1, \dots, K$ **do**
 - 4: Add new transition (s, a, r, s') into \mathcal{D} by executing behavior policy π_b .
 - 5: **end for**
 - 6: **for** iteration $j = 1, \dots, N$ **do**
 - 7: Update w_ρ^j by solving
$$\min_{w_\rho} \hat{\mathbb{E}}_{\{s, a, s'\} \sim \mathcal{D}} [(\delta(s, a, s') - \rho(s, a))^2].$$
 - 8: Decay the stepsize ζ_j in rate $\mathcal{O}(1/j)$.
 - 9: Compute the stochastic gradients w.r.t. w_V and w_π as $\hat{\nabla}_{w_V} \bar{\ell}(V, \pi)$ and $\hat{\nabla}_{w_\pi} \bar{\ell}(V, \pi)$.
 - 10: Update the parameters of primal function by solving the prox-mappings, *i.e.*,
$$\begin{aligned} \text{update } V: \quad w_V^j &= P_{w_V^{j-1}}(\zeta_j \hat{\nabla}_{w_V} \bar{\ell}(V, \pi)) \\ \text{update } \pi: \quad w_\pi^j &= P_{w_\pi^{j-1}}(\zeta_j \hat{\nabla}_{w_\pi} \bar{\ell}(V, \pi)) \end{aligned}$$
 - 11: **end for**
 - 12: Update behavior policy $\pi_b = \pi^N$.
 - 13: **end for**
-

Role of Bootstrapping Distance k



Role of Dual Embedding η

