

Bellman Unbiasedness: Toward Provably Efficient Distributional Reinforcement Learning with General Value Function Approximation

Taehyun Cho¹ Seungyub Han¹ Seokhun Ju¹
Dohyeong Kim¹ Kyungjae Lee² Jungwoo Lee¹

¹Seoul National University, ²Korea University

ICML
June 14, 2025



Seoul National University
Cognitive Machine Learning Lab.

KOREA
UNIVERSITY

Motivation & Challenges

Why Distributional RL?

- **Distributional RL (DistRL)** models the entire distribution of returns, not just the expectation.

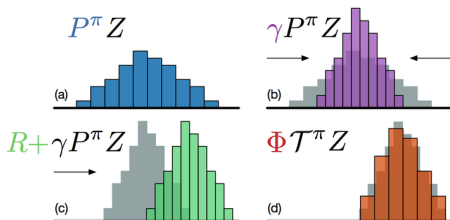


Figure: Distributional Bellman Update

Bellemare, Marc G., Will Dabney, and Rémi Munos. "A distributional perspective on reinforcement learning." International conference on machine learning. PMLR, 2017.

Why Distributional RL?

- **Distributional RL (DistRL)** models the entire distribution of returns, not just the expectation.
- Offers richer insight into **uncertainty**, such as variance, skewness, and quantiles.

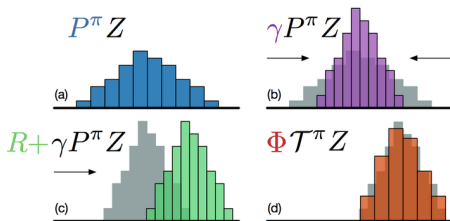


Figure: Distributional Bellman Update

Bellemare, Marc G., Will Dabney, and Rémi Munos. "A distributional perspective on reinforcement learning." International conference on machine learning. PMLR, 2017.

Why Distributional RL?

- **Distributional RL (DistRL)** models the entire distribution of returns, not just the expectation.
- Offers richer insight into **uncertainty**, such as variance, skewness, and quantiles.
- Facilitates safer and more effective decision-making by explicitly considering **risk**.

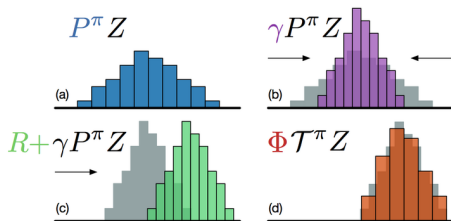


Figure: Distributional Bellman Update

Bellemare, Marc G., Will Dabney, and Rémi Munos. "A distributional perspective on reinforcement learning." International conference on machine learning. PMLR, 2017.

Two challenges in distRL

- 1 **Infinite-dimensionality**
- 2 **Online distributional update**

Two challenges in distRL

1 Infinite-dimensionality

- Return distributions contain an infinite amount of information.
- We must approximate it using a finite number of parameters or statistical functionals.
- However, not all statistical functionals can be **exactly learned** through the Bellman operator, as the meaning is not preserved.

2 Online distributional update

Two challenges in distRL

1 Infinite-dimensionality

- Return distributions contain an infinite amount of information.
- We must approximate it using a finite number of parameters or statistical functionals.
- However, not all statistical functionals can be **exactly learned** through the Bellman operator, as the meaning is not preserved.

2 Online distributional update

- Decoupling the policy update and the distribution estimation via additional rollouts is **sample-inefficient**.
- Limited rollouts inevitably introduce approximation errors into the estimated distribution.

Two challenges in distRL

1 Infinite-dimensionality

- Return distributions contain an infinite amount of information.
- We must approximate it using a finite number of parameters or statistical functionals.
- However, not all statistical functionals can be **exactly learned** through the Bellman operator, as the meaning is not preserved.

2 Online distributional update

- Decoupling the policy update and the distribution estimation via additional rollouts is **sample-inefficient**.
- Limited rollouts inevitably introduce approximation errors into the estimated distribution.

“Can we design a representation that is both exactly learned and provably sample-efficient?”

Backgrounds

- ① Bellman Closedness
- ② Bellman Unbiasedness
- ③ Statistical Functional Bellman Completeness (SFBC)
- ④ Statistical Functional Least Square Value Iteration (SF-LSVI)

Statistical Functional, Sketch; [Bellemare, 2023]

A **statistical functional** is a mapping from a probability distribution to a real value $\psi : \mathcal{P}(\mathbb{R}) \rightarrow \mathbb{R}$. A **sketch** is a vector-valued function $\psi_{1:N} : \mathcal{P}(\mathbb{R}) \rightarrow \mathbb{R}^N$ specified by an N -tuple where each component is a statistical functional,

$$\psi_{1:N}(\cdot) = (\psi_1(\cdot), \dots, \psi_N(\cdot)).$$

Bellman Closedness (BC)

Bellman Closedness; [Rowland, 2019]

A sketch $\psi_{1:N}$ is **Bellman closed** if there exists an operator $\mathcal{T}_{\psi_{1:N}} : I_{\psi_{1:N}}^{\mathcal{S}} \rightarrow I_{\psi_{1:N}}^{\mathcal{S}}$ such that

$$\psi_{1:N}(\mathcal{T}\bar{\eta}) = \mathcal{T}_{\psi_{1:N}}\psi_{1:N}(\bar{\eta}) \quad \text{for all } \bar{\eta} \in \mathcal{P}(\mathbb{R})^{\mathcal{S}}$$

which is closed under a distributional Bellman operator $\mathcal{T} : \mathcal{P}(\mathbb{R})^{\mathcal{S}} \rightarrow \mathcal{P}(\mathbb{R})^{\mathcal{S}}$.

Bellman Closedness (BC)

Theorem ([Rowland, 2019])

*The only finite sets of statistics **of the form** $\psi(\bar{\eta}) = \mathbb{E}_{Z \sim \bar{\eta}}[h(Z)]$ that are Bellman closed are given by the collections of ψ_1, \dots, ψ_N where its linear span $\{\sum_{n=0}^N \alpha_n \psi_n \mid \alpha_n \in \mathbb{R}, \forall N\}$ is equal to the set of exponential polynomial functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l \exp(\lambda Z)] \mid l = 0, 1, \dots, L, \lambda \in \mathbb{R}\}$, where ψ_0 is the constant functional equal to 1.*

In discount setting, it is equal to the linear span of the set of moment functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l] \mid l = 0, 1, \dots, L\}$ for some $L \leq N$.

Bellman Closedness (BC)

Theorem ([Rowland, 2019])

*The only finite sets of statistics **of the form** $\psi(\bar{\eta}) = \mathbb{E}_{Z \sim \bar{\eta}}[h(Z)]$ that are Bellman closed are given by the collections of ψ_1, \dots, ψ_N where its linear span $\{\sum_{n=0}^N \alpha_n \psi_n \mid \alpha_n \in \mathbb{R}, \forall N\}$ is equal to the set of exponential polynomial functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l \exp(\lambda Z)] \mid l = 0, 1, \dots, L, \lambda \in \mathbb{R}\}$, where ψ_0 is the constant functional equal to 1.*

In discount setting, it is equal to the linear span of the set of moment functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l] \mid l = 0, 1, \dots, L\}$ for some $L \leq N$.

Although both the first and second moments are Bellman closed, the variance is **nonlinear**.

As a result, its Bellman closedness cannot be determined by the existing theory, which only applies to **linear** statistical functionals.

Theorem

Quantile functional cannot be Bellman closed under any additional sketch.

Along with a technical clarification of the proof in [Rowland, 2019], we provide an improved version of the proof.

Bellman Closedness (BC)

Theorem

Quantile functional cannot be Bellman closed under any additional sketch.

Along with a technical clarification of the proof in [Rowland, 2019], we provide an improved version of the proof.

Theorem

Maximum and minimum functional are nonlinear and Bellman closed.

- $\mathcal{T}_{\psi_{\max}}\left(\psi_{\max}(\bar{\eta}(s))\right) = \max_{s' \sim \mathbb{P}(\cdot|s,a)} \left(r + \psi_{\max}(\bar{\eta}(s'))\right).$
- $\mathcal{T}_{\psi_{\min}}\left(\psi_{\min}(\bar{\eta}(s))\right) = \min_{s' \sim \mathbb{P}(\cdot|s,a)} \left(r + \psi_{\min}(\bar{\eta}(s'))\right).$

Key Concepts

- ① Bellman Closedness
- ② **Bellman Unbiasedness**
- ③ Statistical Functional Bellman Completeness (SFBC)
- ④ Statistical Functional Least Square Value Iteration (SF-LSVI)

Bellman Unbiasedness (BU)

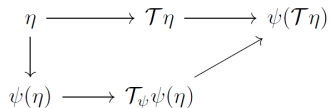


Figure 3. Bellman Closedness

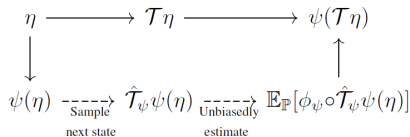


Figure 4. Bellman Unbiasedness

Bellman Unbiasedness (BU)

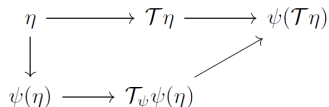


Figure 3. Bellman Closedness

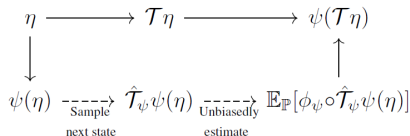


Figure 4. Bellman Unbiasedness

Bellman Closedness

- Exact learnability
- Exact update in finite dimensional space.

Bellman Unbiasedness (BU)

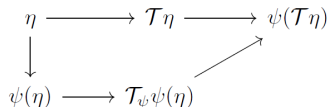


Figure 3. Bellman Closedness

Bellman Closedness

- Exact learnability
- Exact update in finite dimensional space.

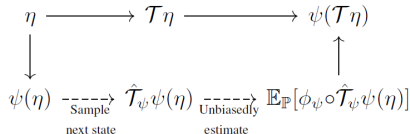


Figure 4. Bellman Unbiasedness

Bellman Unbiasedness

- Provable efficiency
- Unbiased update using sampled sketches

Bellman Unbiasedness (BU)

Bellman Unbiasedness

A sketch $\psi(= \psi_{1:N})$ is **Bellman unbiased** if a vector-valued estimator $\phi_\psi = \phi_\psi(\psi(\cdot), \dots, \psi(\cdot)) : (I_\psi^S)^k \rightarrow I_\psi^S$ exists where the sketch of expected distribution can be unbiasedly estimated by ϕ_ψ using the k sampled sketches from the sample distribution, i.e.,

$$\mathbb{E}_{s'_i \sim \mathbb{P}} \left[\phi_\psi \left(\underbrace{\psi((B_r)_\# \bar{\eta}(s'_1)), \dots, \psi((B_r)_\# \bar{\eta}(s'_k))}_{k \text{ sampled sketches from sample distribution } \hat{\mathcal{T}}_\psi, \psi(\bar{\eta}(s))} \right) \right] = \psi((B_r)_\# \mathbb{E}_{s' \sim \mathbb{P}(\cdot|s,a)}[\bar{\eta}(s')]).$$

Example) Mean-variance sketch

$$\begin{aligned} (\bar{\mu}, \bar{\sigma}^2) &= \phi_{(\mu, \sigma^2)}((\hat{\mu}_1, \hat{\sigma}_1^2), \dots, (\hat{\mu}_k, \hat{\sigma}_k^2)) \\ &= \left(\frac{1}{k} \sum_{i=1}^k \hat{\mu}_i, \frac{1}{k} \sum_{i=1}^k (\hat{\mu}_i - \frac{1}{k} \sum_{j=1}^k \hat{\mu}_j)^2 + \hat{\sigma}_i^2 \right) \end{aligned}$$

Bellman Unbiasedness (BU)

Lemma

Let $F_{\bar{\eta}}$ be a CDF of the probability distribution $\bar{\eta} \in \mathcal{P}_{\psi}(\mathbb{R})^S$. Then a sketch is Bellman unbiased if and only if the sketch is **homogeneous** over $\mathcal{P}_{\psi}(\mathbb{R})^S$ of degree k , i.e., there exists some vector-valued function $h = h(x_1, \dots, x_k) : \mathcal{X}^k \rightarrow \mathbb{R}^N$ such that

$$\psi(\bar{\eta}) = \int \cdots \int h(x_1, \dots, x_k) dF_{\bar{\eta}}(x_1) \cdots dF_{\bar{\eta}}(x_k).$$

Example) Variance is nonlinear but homogeneous of degree 2.

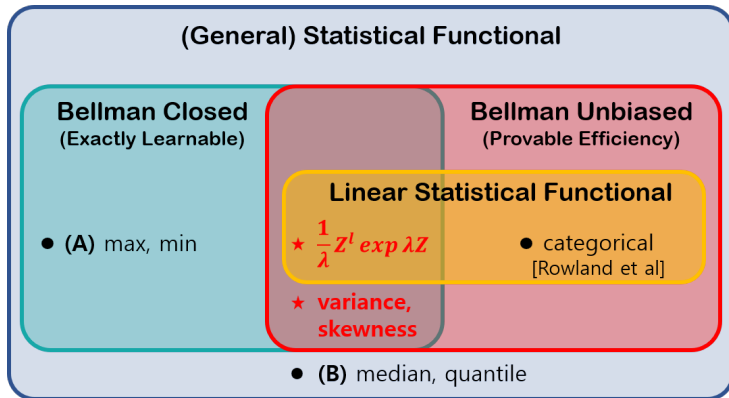
$$\begin{aligned} \text{Var}(\bar{\eta}) &= \mathbb{E}_{Z_1 \sim \bar{\eta}}[(Z_1 - \mathbb{E}_{Z_2 \sim \bar{\eta}}[Z_2])^2] \\ &= \mathbb{E}_{Z_1, Z_2 \sim \bar{\eta}}[Z_1^2 - 2Z_1Z_2 + Z_2^2] = \mathbb{E}_{Z_1, Z_2 \sim \bar{\eta}}[h(Z_1, Z_2)] \end{aligned}$$

Theorem

*The only finite statistical functionals that are **both Bellman unbiased and closed** are given by the collections of ψ_1, \dots, ψ_N where its linear span $\{\sum_{n=0}^N \alpha_n \psi_n \mid \alpha_n \in \mathbb{R}, \forall N\}$ is equal to the set of exponential polynomial functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l \exp(\lambda Z)] \mid l = 0, 1, \dots, L, \lambda \in \mathbb{R}\}$, where ψ_0 is the constant functional equal to 1.*

In discount setting, it is equal to the linear span of the set of moment functionals $\{\eta \rightarrow \mathbb{E}_{Z \sim \eta}[Z^l] \mid l = 0, 1, \dots, L\}$ for some $L \leq N$.

Venn-Diagram of Statistical Functional Classes



Assumption & Algorithm

- ① Bellman Closedness
- ② Bellman Unbiasedness
- ③ Statistical Functional Bellman Completeness (SFBC)
- ④ Statistical Functional Least Square Value Iteration (SF-LSVI)

Statistical Functional Bellman Completeness

Distributional Bellman Completeness (DistBC)

For any distribution $\bar{\eta} : \mathcal{S} \rightarrow \mathcal{P}([0, H])$ and $h \in [H]$, there exists $f_{\bar{\eta}} \in \mathcal{H} (\subseteq \mathcal{F}^\infty)$ which satisfies

$$f_{\bar{\eta}}(s, a) = (\mathcal{B}_{r_h})_{\#}[\mathbb{P}\bar{\eta}](s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

Statistical Functional Bellman Completeness

Distributional Bellman Completeness (DistBC)

For any distribution $\bar{\eta} : \mathcal{S} \rightarrow \mathcal{P}([0, H])$ and $h \in [H]$, there exists $f_{\bar{\eta}} \in \mathcal{H} (\subseteq \mathcal{F}^\infty)$ which satisfies

$$f_{\bar{\eta}}(s, a) = (\mathcal{B}_{r_h})_{\#}[\mathbb{P}\bar{\eta}](s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$



(Relax the representation space to statistical functionals)

Statistical Functional Bellman Completeness (SFBC)

For any distribution $\bar{\eta} : \mathcal{S} \rightarrow \mathcal{P}([0, H])$ and $h \in [H]$, there exists $f_{\bar{\eta}} \in \mathcal{F}^N$ which satisfies

$$f_{\bar{\eta}}(s, a) = \psi_{1:N}\left((\mathcal{B}_{r_h})_{\#}[\mathbb{P}\bar{\eta}](s, a)\right) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

Algorithm 1 Statistical Functional Least Square Value Iteration (**SF-LSVI**(δ))

Input: failure probability $\delta \in (0, 1)$ and the number of episodes K

```

1: for episode  $k = 1, 2, \dots, K$  do
2:   Receive initial state  $s_1^k$ 
3:   Initialize  $\psi_{1:N}(\bar{\eta}_{H+1}^k(\cdot)) \leftarrow \mathbf{0}^N$ 
4:   for step  $h = H, H-1, \dots, 1$  do
5:      $\mathcal{D}_h^k \leftarrow \left\{ s_{h'}^\tau, a_{h'}^\tau, \psi_{1:N} \left( (\mathcal{B}_{r_{h'}^\tau})_{\#} \bar{\eta}_{h+1}^k(s_{h'+1}^\tau) \right) \right\}_{(\tau, h') \in [k-1] \times [H]}$  // Data collection
6:      $\tilde{f}_{h,\bar{\eta}}^k \leftarrow \arg \min_{f \in \mathcal{F}^N} \|f\|_{\mathcal{D}_h^k}$  // Distribution Estimation
7:      $b_h^k(\cdot, \cdot) \leftarrow w^{(1)}((\mathcal{F}^N)_h^k, \cdot, \cdot)$ 
8:      $Q_h^k(\cdot, \cdot) \leftarrow \min\{(\tilde{f}_{h,\bar{\eta}}^k)^{(1)}(\cdot, \cdot) + b_h^k(\cdot, \cdot), H\}$ 
9:      $\pi_h^k(\cdot) = \arg \max_{a \in \mathcal{A}} Q_h^k(\cdot, a), V_h^k(\cdot) = Q_h^k(\cdot, \pi_h^k(\cdot))$  // Optimistic planning
10:     $\psi_1(\bar{\eta}_h^k(\cdot, \cdot)) \leftarrow Q_h^k(\cdot, \cdot), \psi_{2:N}(\bar{\eta}_h^k(\cdot, \cdot)) \leftarrow \left( \min\{(\tilde{f}_{h,\bar{\eta}}^k)^{(n)}(\cdot, \cdot), H\} \right)_{n \in [2:N]}$ 
11:     $\psi_1(\bar{\eta}_h^k(\cdot)) \leftarrow V_h^k(\cdot), \psi_{2:N}(\bar{\eta}_h^k(\cdot)) \leftarrow \psi_{1:N}(\bar{\eta}_h^k(\cdot, \pi_h^k(\cdot)))_{n \in [2:N]}$ 
12:  for  $h = 1, 2, \dots, H$  do
13:    Take action  $a_h^k \leftarrow \pi_h^k(s_h^k)$ 
14:    Observe reward  $r_h^k(s_h^k, a_h^k)$  and get next state  $s_{h+1}^k$ .
```

Moment least square regression

$$\tilde{f}_{h,\bar{\eta}}^k \leftarrow \arg \min_{f \in \mathcal{F}} \sum_{\tau=1}^{k-1} \sum_{h'=1}^H \left(\sum_{n=1}^N f^{(n)}(s_{h'}^\tau, a_{h'}^\tau) - \psi_n \left((\mathcal{B}_{r_{h'}^\tau})_{\#} \bar{\eta}_{h+1}^k(s_{h'+1}^\tau) \right) \right)^2$$

Theorem

*Under SFBC assumption, with probability at least $1 - \delta$, **SF-LSVI** achieves a regret bound of*

$$\text{Reg}(K) \leq 2H \dim_E(\mathcal{F}^N, 1/T) + 4H \sqrt{KH \log(1/\delta)}.$$

Theorem

Under SFBC assumption, with probability at least $1 - \delta$, **SF-LSVI** achieves a regret bound of

$$\text{Reg}(K) \leq 2H \dim_E(\mathcal{F}^N, 1/T) + 4H\sqrt{KH \log(1/\delta)}.$$

Table 1. Comparison for different methods under distributional RL framework. \mathcal{H} represents a subspace of infinite-dimensional space \mathcal{F}^∞ . To bound the eluder dimension d_E , Wang et al. (2023) and Chen et al. (2024) assumed the discretized reward MDP.

Algorithm	Regret	Eluder dimension d_E	Bellman Completeness	MDP assumption	Finite Representation	Exactly Learnable
O-DISCO (Wang et al., 2023)	$\tilde{O}(\text{poly}(d_E H) \sqrt{K})$	$\dim_E(\mathcal{H}, \epsilon)$	distributional BC	discretized reward, small-loss bound	✗	✗
V-EST-LSR (Chen et al., 2024)	$\tilde{O}(d_E H^2 \sqrt{K})^2$	$\dim_E(\mathcal{H}, \epsilon)$	distributional BC	discretized reward, lipschitz continuity	✗	✗
SF-LSVI [Ours]	$\tilde{O}(d_E H^{\frac{3}{2}} \sqrt{K})$	$\dim_E(\mathcal{F}^N, \epsilon)$	statistical functional BC	none	✓	✓

→ Compared to previous distRL methods, **SF-LSVI** achieves a **tighter** regret bound under a **weaker** structural assumption.

To sum up,

- **Bellman Unbiasedness** provides a foundation for designing *exactly learnable* and *provably efficient* distRL algorithm.
- We show that only **moment-based functionals** can be exactly learned—even among *nonlinear* statistical functionals.
- **SF-LSVI** achieves a *tighter* regret bound under a *weaker* assumption, **SFBC**.

Thank you!

ArXiv Link



References



[Mark Rowland \(2019\)](#)

Statistics and samples in distributional reinforcement learning



[Marc G. Bellemare \(2023\)](#)

Distributional reinforcement learning



[Kaiwen Wang \(2023\)](#)

The Benefits of Being Distributional: Small-Loss Bounds for Reinforcement Learning



[Yu Chen \(2024\)](#)

Provable Risk-Sensitive Distributional Reinforcement Learning with General Function Approximation