

# Distributional Hamilton-Jacobi-Bellman Equations for Continuous-Time Reinforcement Learning

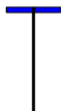
Harley Wiltzer    David Meger    Marc G. Bellemare



**CIM** CENTRE FOR  
INTELLIGENT  
MACHINES

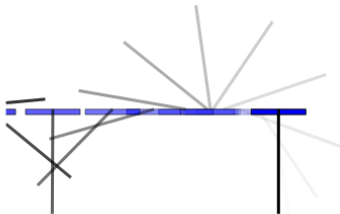


# Continuous-Time Reinforcement Learning



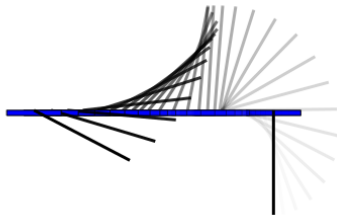
# Continuous-Time Reinforcement Learning

Control frequency:  $\omega$



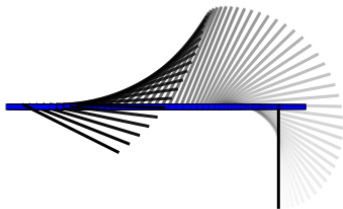
# Continuous-Time Reinforcement Learning

Control frequency:  $2\omega$



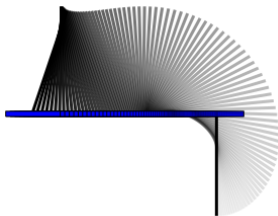
# Continuous-Time Reinforcement Learning

Control frequency:  $5\omega$



# Continuous-Time Reinforcement Learning

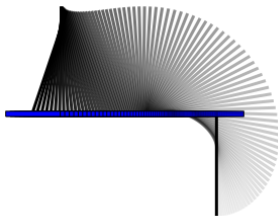
Control frequency:  $10\omega$



- ▶ Challenge: learn a value function that converges as  $\omega \uparrow \infty$ .

# Continuous-Time Reinforcement Learning

Control frequency:  $10\omega$



- ▶ Challenge: learn a value function that converges as  $\omega \uparrow \infty$ .
- ▶ Characterized by the HJB Equation:

$$V^\pi(x) \log \gamma + r(x) + \langle \nabla V^\pi(x), \mu_\pi(x) \rangle + \frac{1}{2} \text{Tr}(\sigma_\pi(x)^\top \mathbf{H} V^\pi(x) \sigma_\pi(x)) = 0$$

## Contribution 1: The *Distributional* HJB Equation

- ▶ Goal: rather than learn the expected return, learn the return *distribution*.



## Contribution 1: The *Distributional* HJB Equation

- ▶ Goal: rather than learn the expected return, learn the return *distribution*.

$$V^\pi(x) \log \gamma + r(x) + \langle \nabla V^\pi(x), \mu_\pi(x) \rangle + \frac{1}{2} \text{Tr}(\sigma_\pi(x)^\top H V^\pi(x) \sigma_\pi(x)) = 0$$

(HJB)

## Contribution 1: The *Distributional* HJB Equation

- ▶ Goal: rather than learn the expected return, learn the return *distribution*.

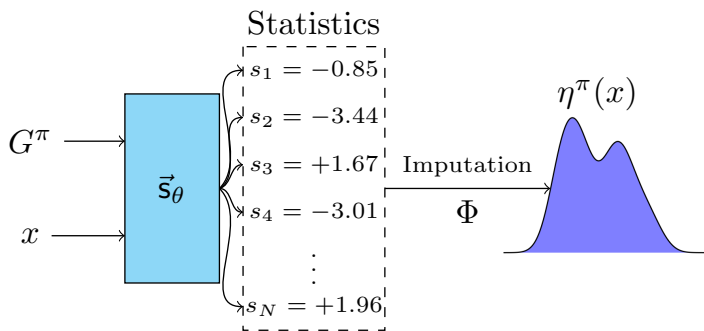
$$V^\pi(x) \log \gamma + r(x) + \langle \nabla V^\pi(x), \mu_\pi(x) \rangle + \frac{1}{2} \text{Tr}(\sigma_\pi(x)^\top \mathbf{H} V^\pi(x) \sigma_\pi(x)) = 0 \quad (\text{HJB})$$

⇓

$$\langle \nabla_x F_{\eta^\pi}(x, z), \mu_\pi(x) \rangle - (r(x) + z \log \gamma) \frac{\partial}{\partial z} F_{\eta^\pi}(x, z) + \frac{1}{2} \text{Tr}(\sigma_\pi(x)^\top \mathbf{H}_x F_{\eta^\pi}(x, z) \sigma_\pi(x)) = 0 \quad (\text{DHJB})$$

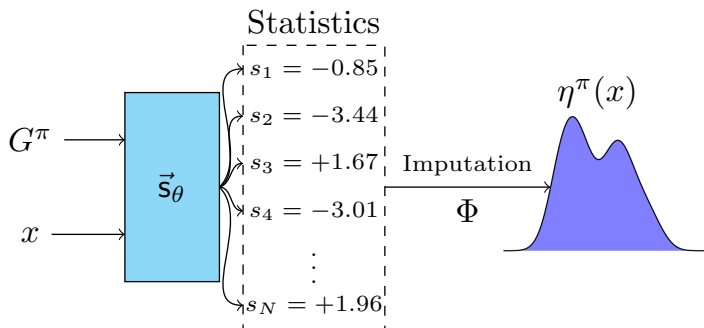
## Contribution 2: The Statistical HJB Loss

- ▶ Goal: Represent the DHJB equation in finite space.



## Contribution 2: The Statistical HJB Loss

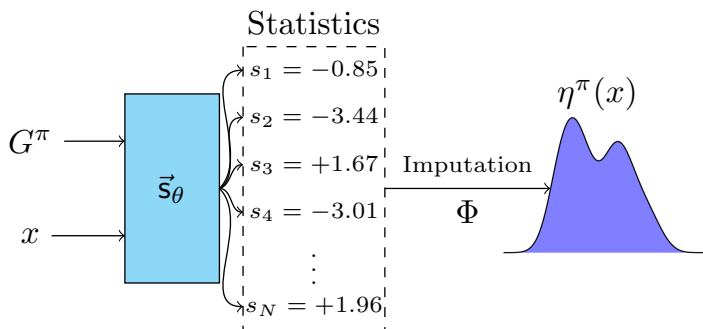
- ▶ Goal: Represent the DHJB equation in finite space.



$$\begin{aligned} \nabla_{\vec{s}(x)} \Phi(\vec{s}(x), z)^\top \mathbf{J}\vec{s}(x) \mu_\pi(x) - (r(x) + z \log \gamma) \frac{\partial}{\partial z} \Phi(\vec{s}(x), z) & \quad (\text{SHJB}) \\ + \frac{1}{2} \text{Tr} \left[ \sigma_\pi(x)^\top (\mathbf{K}_\Phi^x(x, z) + \mathbf{K}_\Phi^s(x, z)) \sigma_\pi(x) \right] \xrightarrow{N \uparrow \infty} 0 \end{aligned}$$

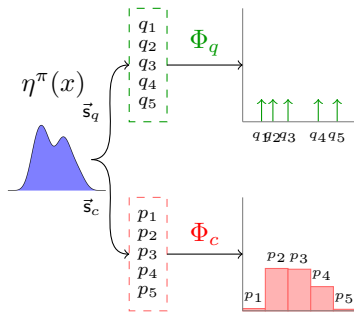
## Contribution 2: The Statistical HJB Loss

- Goal: Represent the DHJB equation in finite space.



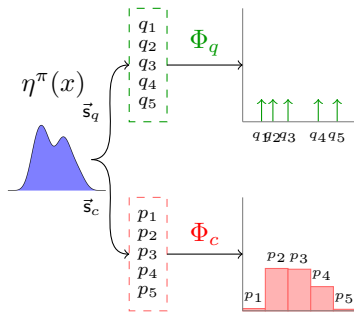
$$\begin{aligned} & \nabla_{\vec{s}(x)} \Phi(\vec{s}(x), z)^\top \mathbf{J}\vec{s}(x) \mu_\pi(x) - (r(x) + z \log \gamma) \frac{\partial}{\partial z} \Phi(\vec{s}(x), z) \quad (\text{SHJB}) \\ & + \frac{1}{2} \text{Tr} \left[ \sigma_\pi(x)^\top \underbrace{(\mathbf{K}_\Phi^x(x, z) + \mathbf{K}_\Phi^s(x, z))}_{\text{"Spatial \& Statistical Diffusivity"}} \sigma_\pi(x) \right] \xrightarrow{N \uparrow \infty} 0 \end{aligned}$$

# Diffusivities of some common imputation strategies



- ▶ When  $\Phi_q$  is the *quantile imputation strategy*, the statistical diffusivity vanishes.

# Diffusivities of some common imputation strategies



- ▶ When  $\Phi_q$  is the *quantile imputation strategy*, the statistical diffusivity vanishes.
- ▶ And when  $\Phi_c$  is the *categorical imputation strategy*, the SHJB is very complex.

# The Quantile Case

We show that, when return distributions are represented as empirical distributions,

$$\begin{cases} \langle \nabla_x \vec{s}_k(x), \mu_\pi(x) \rangle + r(x) + \vec{s}_k(x) \log \gamma + \frac{1}{2} \text{Tr} (\sigma_\pi(x)^\top \mathbf{H}_x \vec{s}_k(x) \sigma_\pi(x)) = 0 \\ \vec{s}_k(x) = F_{\eta^\pi}^{-1}(\hat{\tau}_k) \\ \hat{\tau}_k = \frac{k - \frac{1}{2}}{N} \\ k \in [N] \end{cases}$$



# The Quantile Case

We show that, when return distributions are represented as empirical distributions,

$$\left\{ \begin{array}{l} \overbrace{\langle \nabla_x \vec{s}_k(x), \mu_\pi(x) \rangle + r(x) + \vec{s}_k(x) \log \gamma + \frac{1}{2} \text{Tr} (\sigma_\pi(x)^\top \mathbf{H}_x \vec{s}_k(x) \sigma_\pi(x))}^{\text{The HJB Equation!}} = 0 \\ \vec{s}_k(x) = F_{\eta^\pi}^{-1}(\hat{\tau}_k) \\ \hat{\tau}_k = \frac{k - \frac{1}{2}}{N} \\ k \in \{1, 2, \dots, N\} \end{array} \right.$$

# The Quantile Case

We show that, when return distributions are represented as empirical distributions,

$$\left\{ \begin{array}{l} \overbrace{\langle \nabla_x \vec{s}_k(x), \mu_\pi(x) \rangle + r(x) + \vec{s}_k(x) \log \gamma + \frac{1}{2} \text{Tr} (\sigma_\pi(x)^\top \mathbf{H}_x \vec{s}_k(x) \sigma_\pi(x))}^{\text{The HJB Equation!}} = 0 \\ \vec{s}_k(x) = F_{\eta^\pi}^{-1}(\hat{\tau}_k) \\ \hat{\tau}_k = \frac{k - \frac{1}{2}}{N} \\ k \in \{1, 2, \dots, N\} \end{array} \right.$$

- ▶ Notably, distributional dynamic programming reduces to dynamic programming.

# Thanks!

Check out our poster, #4711, to learn more.