Neural network based policy iteration for values and controls of stochastic games

YUFEI ZHANG

based on work with Kazufumi Ito and Christoph Reisinger

ESSFM, 4 September 2019

Oxford Mathematics



Mathematical Institute





Oxford Mathematics

DNN for control

・ロト ・ 日 ト ・ モ ト ・ モ ト

Zermelo's Navigation Problem





Oxford Mathematics

DNN for control

ъ

(日)、

Stochastic games on domains and HJBI BVPs

Zermelo's Navigation Problem



A time-optimal control problem: find the optimal trajectories of a ship navigating a region Ω of strong winds.



- The domain Ω is an annulus.
- The wind "mainly" blows towards the positive x-axis with velocity v_c.
- The ship moves with a velocity v_s, and the captain can control its direction α.
- The captain prefers to exit from the inner circle.

Stochastic games on domains and HJBI BVPs

Stochastic exit time problem



Given a bounded open set $\Omega \subset \mathbb{R}^n$, strategies $\{\alpha_t\}_{t\geq 0}$ and $\{\beta_t\}_{t\geq 0}$, we consider $X^{x,\alpha,\beta}$ governed by:

$$dX_t = b(X_t, \alpha_t, \beta_t)dt + \sigma(X_t) dW_t, \quad t \in [0, \infty); \quad X_0 = x \in \Omega,$$

and the following value function:

$$u(x) = \inf_{\alpha \in \mathcal{A}} \sup_{\beta \in \mathcal{B}} \mathbb{E} \left[\int_0^{\tau_{\mathbf{x},\alpha,\beta}} f(X_t^{\mathbf{x},\alpha,\beta},\alpha_t,\beta_t) dt + g(X_{\tau_{\mathbf{x},\alpha,\beta}}^{\mathbf{x},\alpha,\beta}) \right],$$

where $\tau_{\mathbf{x},\alpha,\beta}$ is the first exit time of $X^{\mathbf{x},\alpha,\beta}$ from Ω , \mathcal{A} and \mathcal{B} contain all admissible controls valued in \mathbb{A} and \mathbb{B} , respectively.

Oxford Mathematics

DNN for control

イロト 不得 トイヨト イヨト

Stochastic games on domains and HJBI BVPs

HJBI Dirichlet boundary value problem



The value function u is a solution of the Hamilton–Jacobi– Bellman–Isaacs (HJBI) equation:

$$-a^{ij}(x)\partial_{ij}u + H(x, u, \nabla u) = 0, \quad \text{in } \Omega \subset \mathbb{R}^n,$$

 $u = g, \quad \text{on } \partial\Omega,$

where $a = \sigma \sigma^T / 2$, and the Hamiltonian is given by:

 $H(x, u, \nabla u) = \max_{\alpha \in \mathbb{A}} \min_{\beta \in \mathbb{B}} (b^{i}(x, \alpha, \beta)\partial_{i}u(x) + c(x, \alpha, \beta)u(x) - f(x, \alpha, \beta)).$

Oxford Mathematics

DNN for control

6

э.

・ロット (雪) (山) (山)

HJBI Dirichlet boundary value problem



Let *u* be a solution to the HJBI equation, the optimal (feedback) controls are given by: for all $x \in \Omega$,

```
(\alpha(x),\beta(x)) \\ \in \arg \max_{\alpha \in \mathbb{A}} \min_{\beta \in \mathbb{B}} (b^{i}(x,\alpha,\beta)\partial_{i}u(x) + c(x,\alpha,\beta)u(x) - f(x,\alpha,\beta)).
```

Artificial neural networks

Example one hidden layer







Definition

Let $\varrho:\mathbb{R}\to\mathbb{R}$ be a given function. For each $l=1,\ldots,L$, let

$$T_I: \mathbb{R}^{N_{I-1}} \to \mathbb{R}^{N_I}, \quad T_I(x) = W_I x + b_I$$

for some $W_l \in \mathbb{R}^{N_l \times N_{l-1}}$ and $b_l \in \mathbb{R}^{N_l}$. A function $F : \mathbb{R}^{N_0} \to \mathbb{R}^{N_L}$,

$$F(x) = (T_L \circ (\varrho \circ T_{L-1}) \circ \cdots (\varrho \circ T_1))(x), \quad x \in \mathbb{R}^{N_0},$$

is called (the realisation of) a feedforward neural network.

- L is the depth of F.
- N_1, \ldots, N_{L-1} are the dimensions of the hidden layers.
- The number of unknown $\{W_l, b_l\}_{l=1}^N$ is the complexity \mathscr{C} of F.

Existing numerical methods

Neural network-based methods



• Choose a family of neural networks \mathcal{F}_M parametrized by $\theta := \{W_l, b_l\}_{l=1}^L \in \mathbb{R}^M$, i.e., $\mathcal{F}_M = \{u^{\theta} \mid \theta \in \mathbb{R}^M\}$.

Formulate the semilinear Dirichlet problem

$$\begin{aligned} -a^{ij}(x)\partial_{ij}u + H(x,u,\nabla u) &= 0, \quad \text{in } \Omega, \\ u &= g, \quad \text{on } \partial\Omega \end{aligned}$$

into a nonlinear optimization problem over \mathcal{F}_M :

$$\underbrace{\min_{\theta \in \mathbb{R}^{M}}}_{\min_{u^{\theta} \in \mathcal{F}_{M}}} \| - a^{ij} \partial_{ij} u^{\theta} + H(\cdot, u^{\theta}, \nabla u^{\theta}) \|_{L^{2}(\Omega)}^{2} + \| u^{\theta} - g \|_{Y}^{2}.$$

ote that $Y = L^{2}(\partial \Omega)$ in most published works.

See, e.g., J. Sirignano and K. Spiliopoulos, DGM: A deep learning algorithm for solving partial differential equations, J. Comput. Phys., 375 (2018), pp. 1339–1364.

Oxford Mathematics

Ν

DNN for control

Existing numerical methods

Neural network-based methods



Discretize the optimization problem by using samples Ω_d ⊂ Ω and ∂Ω_d ⊂ ∂Ω: solve min_{θ∈ℝ^M} J(θ), where

$$\begin{split} J(\theta) &= \frac{1}{|\Omega_d|} \sum_{x_i \in \Omega_d} |-a^{ij}(x_i) \partial_{ij} u^{\theta}(x_i) + H(x_i, u^{\theta}(x_i), \nabla u^{\theta}(x_i))|^2 \\ &+ \frac{1}{|\partial \Omega_d|} \sum_{x_i \in \partial \Omega_d} |u^{\theta}(x_i) - g(x_i)|^2. \end{split}$$

イロト 不得 トイヨト イヨト

Existing numerical methods

Neural network-based methods



Discretize the optimization problem by using samples Ω_d ⊂ Ω and ∂Ω_d ⊂ ∂Ω: solve min_{θ∈ℝ^M} J(θ), where

$$\begin{split} J(\theta) &= \frac{1}{|\Omega_d|} \sum_{x_i \in \Omega_d} |-a^{ij}(x_i) \partial_{ij} u^{\theta}(x_i) + H(x_i, u^{\theta}(x_i), \nabla u^{\theta}(x_i))|^2 \\ &+ \frac{1}{|\partial \Omega_d|} \sum_{x_i \in \partial \Omega_d} |u^{\theta}(x_i) - g(x_i)|^2. \end{split}$$

Find θ ∈ arg min_{θ∈ℝ^M} J(θ) by minimizing J(θ) with stochastic gradient descent (SGD) method:

$$\theta^{k+1} \coloneqq \theta^k - \eta \nabla_{\theta} J(\theta^k), \quad \forall k \ge 1.$$

Neural network-based methods



Pros:

- It is a mesh-free method.
- Neural networks can approximate certain high-dimensional functions with reasonable complexity.

Cons:

- It is a non-convex non-smooth optimization problem due to the non-convexity of the Hamiltonian H.
- SGD requires evaluations of H and its gradient, which may not be possible if A and B are high-dimensional.
- The choice of Y = L²(∂Ω) does not ensure the convergence of ∇u, hence no convergence of optimal controls.



Iteratively linearize the HJBI equation:

$$-a^{ij}\partial_{ij}u + H(\cdot, u,
abla u) = 0,$$
 in Ω ,

and then solve the linear equations by neural networks.

Recall the optimal controls associated with u are given by:

$$(\alpha(\cdot),\beta(\cdot)) \\ \in \arg \max_{\alpha \in \mathbb{A}} \min_{\beta \in \mathbb{B}} (b^{i}(\cdot,\alpha,\beta)\partial_{i}u + c(\cdot,\alpha,\beta)u - f(\cdot,\alpha,\beta)).$$

Oxford Mathematics

ъ

・ロト ・ 母 ト ・ ヨ ト ・ ヨ ト



Algorithm (Policy iteration for HJBI equations)

Initialization: choose an initial guess u^0 and set k = 0. 1. Given the function u^k , update the control laws: $\forall x \in \Omega$,

$$(\alpha^{k}(x),\beta^{k}(x)) \\ \in \arg \max_{\alpha \in \mathbb{A}} \min_{\beta \in \mathbb{B}} \left(b^{i}(x,\alpha,\beta) \partial_{i} u^{k}(x) + c(x,\alpha,\beta) u^{k}(x) - f(x,\alpha,\beta) \right)$$

2. Solve the linear problem for u^{k+1} :

$$-a^{ij}\partial_{ij}u+b^i_k\partial_iu+c_ku-f_k=0, \quad ext{in } \Omega; \quad u=g, \quad ext{on } \partial\Omega, \ (1)$$

where $\phi_k(\cdot) \coloneqq \phi(\cdot, \alpha^k(\cdot), \beta^k(\cdot))$ for $\phi = b^i, c, f$.

Oxford Mathematics



Algorithm (Inexact policy iteration for HJBI equations)

Initialization: choose a family of trial functions $\mathcal{F} \subset H^2(\Omega)$, an initial guess $u^0 \in \mathcal{F}$, a sequence $\{\eta_k\}_{k \in \mathbb{N}}$ of non-negative scalars.

- 1. For each $k \ge 0$, given the function $u^k \in \mathcal{F}$, update the control laws (α^k, β^k) as before.
- 2. Find $u^{k+1} \in \mathcal{F}$ such that

$$\begin{split} \|(-a^{ij}\partial_{ij}+b^i_k\partial_i+c_k)u^{k+1}-f_k\|^2_{L^2(\Omega)}+\|u^{k+1}-g\|^2_{H^{3/2}(\partial\Omega)}\\ &\leq \eta_{k+1}\|u^{k+1}-u^k\|^2_{H^2(\Omega)}, \end{split}$$

where b_k , c_k and f_k depend on (α^k, β^k) as in (1).

Oxford Mathematics

DNN for control

Convergence analysis



Theorem

Suppose \mathcal{F} is dense in $H^2(\Omega)$, $\lim_{k\to\infty} \eta_k = 0$, and the coefficients are sufficiently regular.

Convergence analysis



Theorem

Suppose \mathcal{F} is dense in $H^2(\Omega)$, $\lim_{k\to\infty} \eta_k = 0$, and the coefficients are sufficiently regular. Let $u^* \in H^2(\Omega)$ be the solution to the HJBI equation.

Then for any initial guess $u^0 \in \mathcal{F}$, we have:

1. $\{u^k\}_{k\in\mathbb{N}}$ converge superlinearly to u^* in $H^2(\Omega)$, i.e., $\lim_{k\to\infty} \|u^{k+1} - u^*\|_{H^2(\Omega)}/\|u^k - u^*\|_{H^2(\Omega)} = 0.$

Convergence analysis



Theorem

Suppose \mathcal{F} is dense in $H^2(\Omega)$, $\lim_{k\to\infty} \eta_k = 0$, and the coefficients are sufficiently regular. Let $u^* \in H^2(\Omega)$ be the solution to the HJBI equation.

Then for any initial guess $u^0 \in \mathcal{F}$, we have:

- 1. $\{u^k\}_{k\in\mathbb{N}}$ converge superlinearly to u^* in $H^2(\Omega)$, i.e., $\lim_{k\to\infty} \|u^{k+1} - u^*\|_{H^2(\Omega)}/\|u^k - u^*\|_{H^2(\Omega)} = 0.$
- 2. $\lim_{k\to\infty} (u^k, \partial_i u^k, \partial_{ij} u^k)(x) = (u^*, \partial_i u^*, \partial_{ij} u^*)(x)$ for a.e. $x \in \Omega$, and for all i, j = 1, ..., n.

Convergence analysis



Theorem

Suppose \mathcal{F} is dense in $H^2(\Omega)$, $\lim_{k\to\infty} \eta_k = 0$, and the coefficients are sufficiently regular. Let $u^* \in H^2(\Omega)$ be the solution to the HJBI equation.

Then for any initial guess $u^0 \in \mathcal{F}$, we have:

- 1. $\{u^k\}_{k\in\mathbb{N}}$ converge superlinearly to u^* in $H^2(\Omega)$, i.e., $\lim_{k\to\infty} \|u^{k+1} - u^*\|_{H^2(\Omega)}/\|u^k - u^*\|_{H^2(\Omega)} = 0.$
- 2. $\lim_{k\to\infty} (u^k, \partial_i u^k, \partial_{ij} u^k)(x) = (u^*, \partial_i u^*, \partial_{ij} u^*)(x)$ for a.e. $x \in \Omega$, and for all i, j = 1, ..., n.
- 3. The control laws $\{(\alpha^k, \beta^k)\}_{k \in \mathbb{N}}$ converge to the optimal feedback control (α^*, β^*) .

Oxford Mathematics



The position of
$$\vec{X}_t^{x,\alpha} = (X_t^{x,\alpha}, Y_t^{x,\alpha})$$
 follows:

$$\begin{pmatrix} dX_t^{x,\alpha} \\ dY_t^{x,\alpha} \end{pmatrix} = \begin{pmatrix} v_c(X_t^{x,\alpha}, Y_t^{x,\alpha}) + v_s \cos(\alpha_t) \\ v_s \sin(\alpha_t) \end{pmatrix} dt + \begin{pmatrix} \sigma_x & 0 \\ 0 & \sigma_y \end{pmatrix} dW_t$$



► The velocity of the wind is $v_c = 1 - 0.2 \sin \left(\pi \frac{x^2 + y^2 - r^2}{R^2 - r^2} \right).$

イロト 不得 とうほう 不良 とう

•
$$\sigma_x = 0.5$$
 and $\sigma_y = 0.2$.

•
$$r = 0.5$$
 and $R = \sqrt{2}$.

Oxford Mathematics

DNN for control

Let $\alpha_t \in \mathbb{A} = [0, 2\pi]$ be the direction of the boat. We consider

$$u(x) = \inf_{\alpha \in \mathcal{A}} \sup_{\mathbb{Q} \in \mathcal{M}} \mathbb{E}_{\mathbb{Q}} [\tau_{x,\alpha} + g(\vec{X}^{x,\alpha}_{\tau_{x,\alpha}})].$$



- The first exit time is $\tau_{x,\alpha} \coloneqq \inf\{t \ge 0 \mid \vec{X}_t^{x,\alpha} \notin \Omega\}.$
- $g \equiv 0$ on $\partial B_r(0)$ and $g \equiv 1$ on $\partial B_R(0)$.
- *M* denotes the uncertainty from the unknown law of the random perturbation.

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

Oxford Mathematics

DNN for control

Mathematica Institute

Zermelo's Navigation Problem





Figure: The feedback control for the scenario where the ship moves slower than the wind ($v_s = 0.5$).

Oxford Mathematics

DNN for control

23

・ロト ・ 同ト ・ 日下 ・ 日

Zermelo's Navigation Problem





Figure: The value function for the scenario where the ship moves slower than the wind ($v_s = 0.5$).

Oxford Mathematics

DNN for control

э

(日)、

Numerical experiments: Zermelo's Navigation Problem



$\text{Residual} := \| - a^{ij} \partial_{ij} u^k + H(\cdot, u^k, \nabla u^k) \|_{2,\Omega,\text{val}}^2 + \| u^k - g \|_{3/2,\partial\Omega,\text{val}}^2.$



Figure: Residuals with respect to the number of policy iterations.

Oxford Mathematics

DNN for control

25

э

・ロト ・ 雪 ト ・ ヨ ト



Iterative fitting for stochastic games with controlled drifts, using

- Sobolev spaces (for feedback controls);
- Newton methods (a.k.a. policy iteration) in function space;
- stochastic gradient descent (weakest link).

Extensions to finite-horizon games and optimal stopping problems.

Reference:

1. K. Ito, C. Reisinger, Y. Zhang. A neural network based policy iteration algorithm with global H²-superlinear convergence for stochastic games on domains, arXiv:1906.02304.

Oxford Mathematics

イロト 不得 トイヨト イヨト