

# Deep Learning Numerical Methods on High-Dimensional PDEs -- The Deep BSDE method

Presenter: Chenyang Fang  
10/31/19



---

## Papers of Interest

- Bing, Yu , Xiaojing Xing , Agus Sudjianto, Deep-learning based numerical BSDE method for barrier options, *arxiv*, 2017
- Weinan E , Jiequn Han , and Arnulf Jentzen, Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations -- *Communications in Mathematics and Statistics* 2017
- Côme Huré, Huyên Pham, Xavier Warin, Some machine learning schemes for high-dimensional nonlinear PDEs, *arxiv*, 2019

---

# Outline

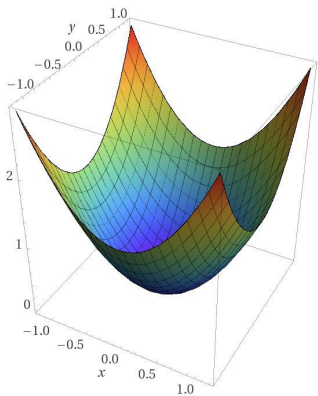
- Introduction
  - Problems of High Dimensional Nonlinear PDEs
  - Solution
  - Background
- Neural Network & Deep Learning Schemes
  - Neural Network as Approximation
  - Deep learning-based schemes for semi-linear PDEs
    - DBSP
    - DBDP1 & DBDP2
    - RDBDP
- Convergence Analysis
  - DBDP1
  - DBDP2
- Applications

# Problem of High Dimensional PDEs

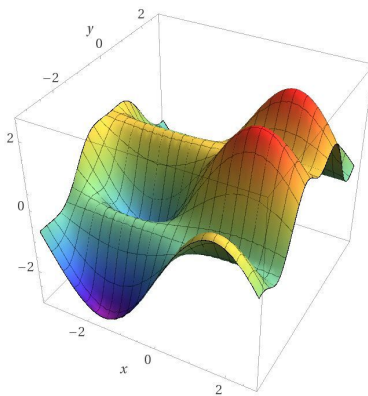
Part 1

# Problems of High Dimensional Nonlinear PDE

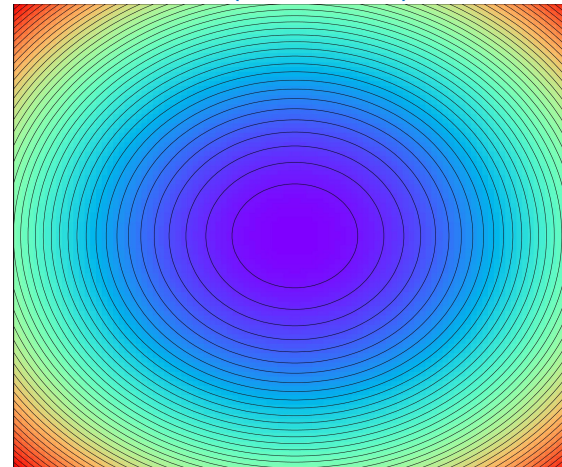
- “curse of dimensionality”
  - The BSDE (Backward Stochastic Differential Equation) approach has recently been developed to tackle this. Using the time discretization scheme, some effective algorithms based on regression manage to solve nonlinear PDEs in dimension above 4. However this approach is still not implementable in dimension above 6 or 7: the number of basis functions used for the regression still explodes with the dimension(Huré et al.)



Computed by WolframAlpha



Computed by WolframAlpha



High Dimensional Gradient  
Descent

High Dimensional Hamilton-Jacobi  
PDEs

Copyright to all pictures belongs to their respective owners. The pictures are used here for educational purposes only.

# An Elegant, Practical Solution

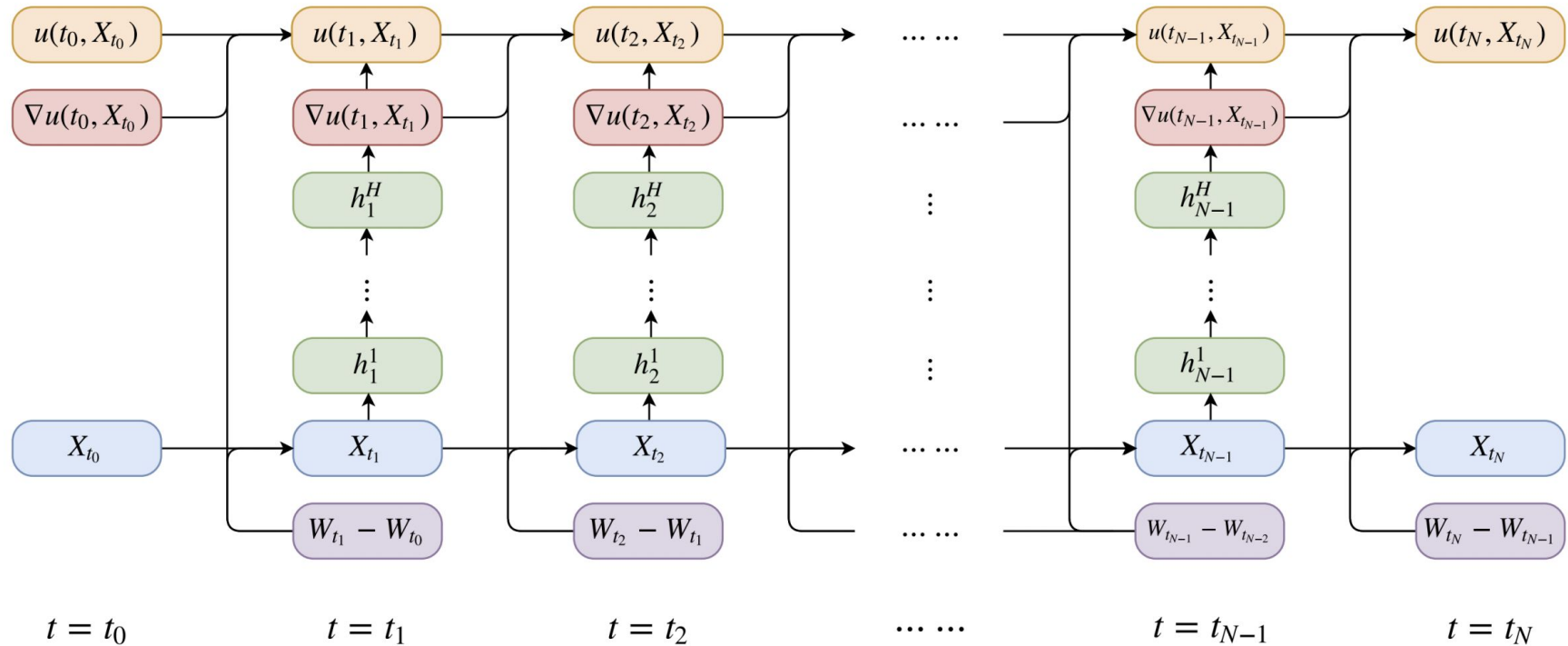


Fig. 1: Graphical illustration of the neural network diagram with  $H$  hidden layers and  $N$  time steps.

Details will be introduced

Figure adapted from *Solving high-dimensional partial differential equations using deep learning*, Han et al. PNAS 2018

# Neural Network & Deep Learning Based Schemes for Semi-linear PDEs

Part 2

---

# High Dimensional Semilinear PDE

$$\begin{aligned} \frac{\partial u}{\partial t}(t, x) + \frac{1}{2} \text{Tr} \left( \sigma \sigma^T(t, x) (\text{Hess}_x u)(t, x) \right) + \nabla u(t, x) \cdot \mu(t, x) \\ + f(t, x, u(t, x), \sigma^T(t, x) \nabla u(t, x)) = 0 \end{aligned}$$

- Terminal Condition:  $u(T, x) = g(x)$
- To fix ideas, we are interested in the solution at  $t = 0$ ,  $x = \xi$  for some vector  $\xi \in \mathbb{R}^d$



---

# Linear Quadratic Gaussian Control example for High Dimensional Semilinear PDE

- Consider a stochastic optimal control problem in  $\mathbb{R}^d$

$$dX_t = 2\sqrt{\lambda} m_t dt + \sqrt{2} dW_t$$

$$J(\{m_t\}_{0 \leq t \leq T}) = \mathbb{E} \left[ \int_0^T \|m_t\|_2^2 dt + g(X_T) \right]$$

- the optimal control is

$$m_t^* = \frac{\nabla u(t, x)}{\sqrt{2\lambda}}$$

$$u(t, x) = -\frac{1}{\lambda} \ln \left( \mathbb{E} \left[ \exp \left( -\lambda g(x + \sqrt{2} W_{T-t}) \right) \right] \right)$$

- via the Hamilton-Jacobi-Bellman equation,

$$\frac{\partial u}{\partial t}(t, x) + \Delta u(t, x) - \lambda \|\nabla u(t, x)\|_2^2 = 0$$

---

# Two Solutions

- (1) Estimate the solution and its gradient by a neural network.
- (2) Approximate the solution by a neural network while its gradient is estimated directly with some numerical differentiation techniques.

I will be focusing on 1 due to time

---

# Backward SDE

- Find a solution to this backward stochastic differential equation :

$$y(t) = (ry(t) + az(t)) * dt + z(t)dW_t$$

- with the terminal condition

$$y(T) = \xi$$

- with  $\xi$  is a square integrable,  $F$  measurable random variable in a filtered probability space and  $W$  is a one dimensional brownian motion. in addition  $r$  and  $a$  are constant

# Nonlinear Feynman-Kac Formula

- Markovian Backward SDEs give a nonlinear Feynman-Kac representation of some nonlinear parabolic PDEs  
(*Backward Stochastic Differential Equations in Finance*, N. et al. *Mathematical Finance* Vol 7. 1997)

- solution to the BSDE,

$$Y_t = g(\mathcal{X}_T) + \int_t^T f(s, \mathcal{X}_s, Y_s, Z_s) ds - \int_t^T Z_s^\top dW_s, \quad 0 \leq t \leq T,$$

- related to the solution  $u$  of the form,

$$\begin{cases} \partial_t u + \mathcal{L}u + f(., ., u, \sigma^\top D_x u) = 0, & \text{on } [0, T) \times \mathbb{R}^d, \\ u(T, .) = g, & \text{on } \mathbb{R}^d, \end{cases}$$

- via

$$Y_t = u(t, \mathcal{X}_t), \quad 0 \leq t \leq T$$

$$Z_t = \sigma^\top(t, \mathcal{X}_t) D_x u(t, \mathcal{X}_t), \quad 0 \leq t \leq T$$

# Connection Between PDEs and BSDEs

- Consider the following BSDE,

$$\begin{cases} X_t = \xi + \int_0^t \mu(s, X_s) ds + \int_0^t \sigma(s, X_s) dW_s, \\ Y_t = g(X_T) + \int_t^T f(s, X_s, Y_s, Z_s) ds - \int_t^T (Z_s)^\top dW_s, \end{cases}$$

- The solution is an adaptive process  $\{(X_t, Y_t, Z_t)\}$
- Under suitable regularity assumptions, the BSDE is well-posed and related to the PDE in the sense that for all  $t \in [0, T]$  it holds a.s. that

$$Y_t = u(t, X_t) \quad \text{and} \quad Z_t = \sigma^\top(t, X_t) \nabla u(t, X_t)$$

- the solution of PDE satisfies the following SDE

$$\begin{aligned} & u(t, X_t) - u(0, X_0) \\ &= - \int_0^t f(s, X_s, u(s, X_s), \sigma^\top(s, X_s) \nabla u(s, X_s)) ds \\ & \quad + \int_0^t [\nabla u(s, X_s)]^\top \sigma(s, X_s) dW_s. \end{aligned}$$

---

# Deep Learning Scheme

- Feedforward Neural Network
  - define linear feedforward neural network

$$x \in \mathbb{R}^d \mapsto A_L \circ \varrho \circ A_{L-1} \circ \dots \circ \varrho \circ A_1(x) \in \mathbb{R}^{d_1}$$

Here  $A_\ell$ ,  $\ell = 1, \dots, L$  are affine transformations:  $A_1$  maps from  $\mathbb{R}^d$  to  $\mathbb{R}^m$ ,  $A_2, \dots, A_{L-1}$  map from  $\mathbb{R}^m$  to  $\mathbb{R}^m$ , and  $A_L$  maps from  $\mathbb{R}^m$  to  $\mathbb{R}^{d_1}$ , represented by

$$A_\ell(x) = \mathcal{W}_\ell x + \beta_\ell,$$

# Time Discretization

- consider the simple Euler scheme of the BSDE, with a partition of the time interval  $[0, T]$ ,  $0 = t_0 < t_1 < \dots < t_N = T$

$$X_{t_{n+1}} - X_{t_n} \approx \mu(t_n, X_{t_n}) \Delta t_n + \sigma(t_n, X_{t_n}) \Delta W_n$$

$$\begin{aligned} & u(t_{n+1}, X_{t_{n+1}}) - u(t_n, X_{t_n}) \\ & \approx -f(t_n, X_{t_n}, u(t_n, X_{t_n}), \sigma^T(t_n, X_{t_n}) \nabla u(t_n, X_{t_n})) \Delta t_n \\ & \quad + [\nabla u(t_n, X_{t_n})]^T \sigma(t_n, X_{t_n}) \Delta W_n, \end{aligned}$$

$$\Delta t_n = t_{n+1} - t_n, \quad \Delta W_n = W_{t_{n+1}} - W_{t_n}$$

# Architecture

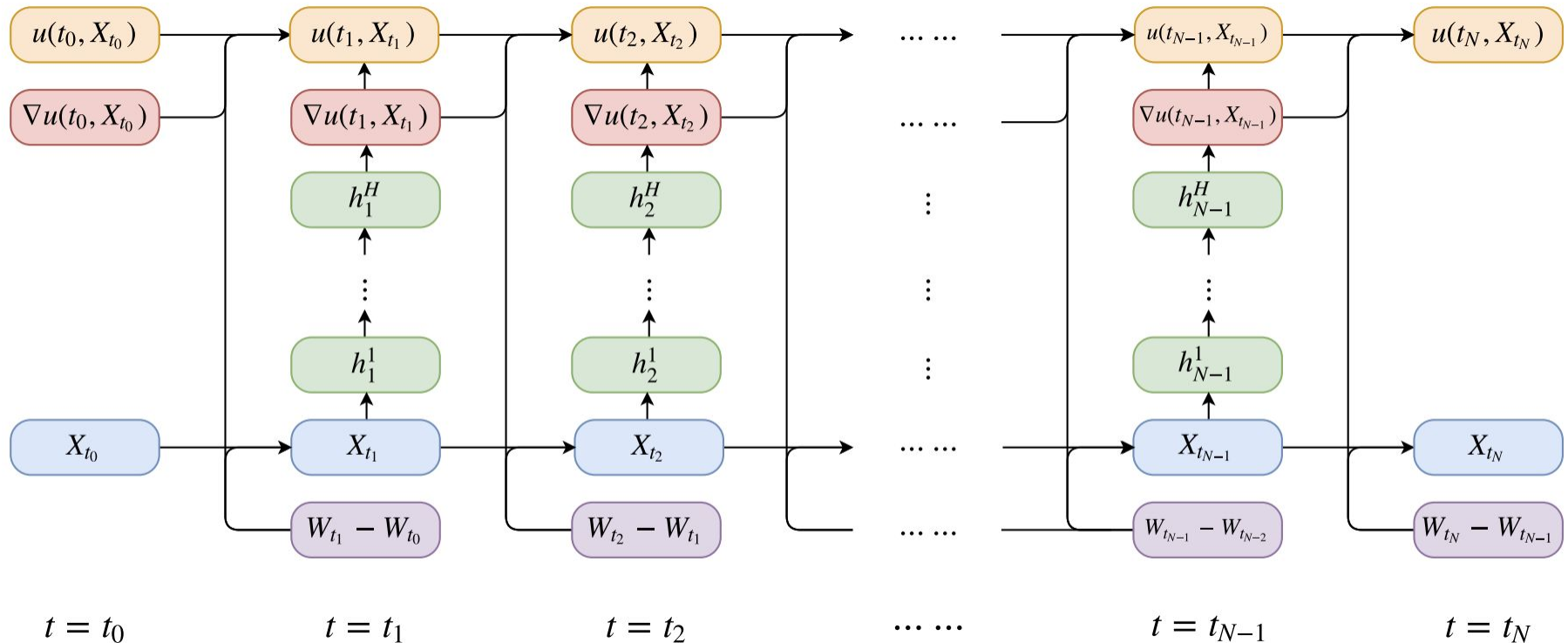


Figure adapted from *Solving high-dimensional partial differential equations using deep learning*, Han et al. PNAS 2018

- Each column corresponds to a subnetwork at time  $t = t_n$ . The whole network has  $(H + 1)(N - 1)$  layers in total that involve free parameters to be optimized simultaneously



---

# Approximation

- The most important step by Han et al. is to approximate

$$x \mapsto \sigma^T(t, x) \nabla u(t, x)$$

at each discretized time step  $t_n$  by a feedforward neural network

$$\begin{aligned} \sigma^T(t_n, X_{t_n}) \nabla u(t_n, X_{t_n}) &= (\sigma^T \nabla u)(t_n, X_{t_n}) \\ &\approx (\sigma^T \nabla u)(t_n, X_{t_n} | \theta_n), \end{aligned}$$

- “We can stack all the subnetworks together to form a deep neural network (DNN) as a whole, based on the time discretization”

## DBSP (this part is adapted from Han et al.'s presentation on the paper)

- According to Han et al. , the error in the matching of given terminal condition defines the expected loss function,
- This network takes the paths  $\{X_{t_n}\}_{0 \leq n \leq N}$  and  $\{W_{t_n}\}_{0 \leq n \leq N}$  as the input data and gives the final output, denoted by  $\hat{u}(\{X_{t_n}\}_{0 \leq n \leq N}, \{W_{t_n}\}_{0 \leq n \leq N})$ , as an approximation to  $u(t_N, X_{t_N})$

$$l(\theta) = \mathbb{E} \left[ \left| g(X_{t_N}) - \hat{u}(\{X_{t_n}\}_{0 \leq n \leq N}, \{W_{t_n}\}_{0 \leq n \leq N}) \right|^2 \right]$$

- The paths can be simulated easily. Therefore the commonly used SGD algorithm fits this problem well.
- Call the introduced methodology deep BSDE method since we use the BSDE and DNN as essential tool (Han, et al.)

---

# The Intuition

- Question:
  - Such deep neural networks can be trained?
  - Intuition: there are skip connections between different subnetworks (Han, et al.)

$$\begin{aligned} & u(t_{n+1}, X_{t_{n+1}}) - u(t_n, X_{t_n}) \\ & \approx -f(t_n, X_{t_n}, u(t_n, X_{t_n}), (\sigma^T \nabla u)(t_n, X_{t_n} | \theta_n)) \Delta t_n \\ & \quad + (\sigma^T \nabla u)(t_n, X_{t_n} | \theta_n) \Delta W_n, \end{aligned}$$

---

# Analogy to Deep Reinforcement Learning

- We want to draw insight from represent policy function (control) through neural networks in Deep Reinforcement Learning

Deep BSDE method		DRL	
BSDE	$\longleftrightarrow$	Markov decision model	
gradient of the solution	$\longleftrightarrow$	optimal policy function	

# DBDP1

- backward dynamic programming
  - use independent neural networks respectively for the approximation of  $u$  and for the approximation of  $\sigma$ .

- Initialize from an estimation  $\widehat{\mathcal{U}}_N^{(1)}$  of  $u(t_N, \cdot)$  with  $\widehat{\mathcal{U}}_N^{(1)} = g$
- For  $i = N-1, \dots, 0$ , given  $\widehat{\mathcal{U}}_{i+1}^{(1)}$ , use a pair of deep neural networks  $(\mathcal{U}_i(\cdot; \theta), \mathcal{Z}_i(\cdot; \theta)) \in \mathcal{NN}_{d,1,L,m}^{\mathcal{Q}}(\mathbb{R}^{N_m}) \times \mathcal{NN}_{d,d,L,m}^{\mathcal{Q}}(\mathbb{R}^{N_m})$  for the approximation of  $(u(t_i, \cdot), \sigma^\top(t_i, \cdot) D_x u(t_i, \cdot))$ , and compute (by SGD) the minimizer of the expected quadratic loss function

$$\begin{cases} \hat{L}_i^{(1)}(\theta) := \mathbb{E} \left| \widehat{\mathcal{U}}_{i+1}^{(1)}(X_{t_{i+1}}) - F(t_i, X_{t_i}, \mathcal{U}_i(X_{t_i}; \theta), \mathcal{Z}_i(X_{t_i}; \theta), \Delta t_i, \Delta W_{t_i}) \right|^2 & \text{Hure et al.} \\ \theta_i^* \in \arg \min_{\theta \in \mathbb{R}^{N_m}} \hat{L}_i^{(1)}(\theta). \end{cases}$$

Then, update:  $\widehat{\mathcal{U}}_i^{(1)} = \mathcal{U}_i(\cdot; \theta_i^*)$ , and set  $\widehat{\mathcal{Z}}_i^{(1)} = \mathcal{Z}_i(\cdot; \theta_i^*)$ .

---

## DBDP as to DBSDE

- The advantages of the two schemes, compared to the Deep BSDE algorithm, are the following:
  - We solve smaller problems that are less prone to be trapped in local minimizer. The memory needed in [HJE17] can be a problem when taking too many time steps.
  - at each time step, we initialize the weights and bias of the neural network to the weights and bias of the previous time step treated : this methodology always used in iterative solvers in PDE methods permits to have a starting point close to the solution, and then to avoid local minima too far away from the true solution. Besides the number of gradient iterations to achieve is rather small after the first resolution step (Hure et al.)

# Convergence Analysis

Part 3

# Convergence of DBDP1

- Proof is LOOOOOONG
  - Approximation error of  $u$  and  $z$  goes to zero as  $n$  goes to infinity according to the universal approximation theorem (HSW, 89).

$$\varepsilon_i^{\mathcal{N},v} := \inf_{\xi} \mathbb{E} |\hat{v}_i(X_{t_i}) - \mathcal{U}_i(X_{t_i}; \xi)|^2, \quad \varepsilon_i^{\mathcal{N},z} := \inf_{\eta} \mathbb{E} |\hat{z}_i(X_{t_i}) - \mathcal{Z}_i(X_{t_i}; \eta)|^2.$$

## – Consistency of DBDP1

**Theorem 4.1.** (Consistency of DBDP1) *Under (H1), there exists a constant  $C > 0$ , independent of  $\pi$ , such that*

$$\begin{aligned} \mathcal{E}[(\hat{\mathcal{U}}^{(1)}, \hat{\mathcal{Z}}^{(1)}), (Y, Z)] &\leq C \left( \mathbb{E} |g(\mathcal{X}_T) - g(X_T)|^2 + |\pi| + \varepsilon^Z(\pi) \right. \\ &\quad \left. + \sum_{i=0}^{N-1} (N \varepsilon_i^{\mathcal{N},v} + \varepsilon_i^{\mathcal{N},z}) \right). \end{aligned} \tag{4.8}$$



# Application

Part 4

# American Option

- Augment the classical Black Scholes model by some important factors in real markets, including defaultable securities, transactions costs, uncertainties in the model parameters, etc.
- “Ideally the pricing models should take into account the whole basket of financial derivative underlyings, resulting in high-dimensional nonlinear PDEs.”(Han et al.)
- To test the deep BSDE method, we study a special case of the recursive valuation model with default risk (Duffie et al. 1996, Bender et al. 2015)

# Formulation

- Risk neutral stock price,

$$dX_t^i = rX_t^i dt + \sigma_i X_t^i dW_t^i,$$

- The value at time  $t$  of an American option with payoff  $g$  and maturity  $T$  is given,

$$u(t, x) = \sup_{\tau \in \mathcal{T}_{t,T}} \mathbb{E}[e^{-r\tau} g(X_\tau)]$$

- $\tau, T$  is the solution to the variational inequality,

$$\begin{cases} \min [-\partial_t u - \hat{\mathcal{L}}u, u - g] = 0, & \text{on } [0, T) \times (0, \infty)^d \\ u(T, \cdot) = g, & \text{on } (0, \infty)^d, \end{cases}$$

- Solve for the BSDE,

$$\hat{\mathcal{L}}u(t, x) = \frac{1}{2} \sum_{i=1}^d \sigma_i^2 x_i^2 D_{x_i}^2 u(t, x) + r \sum_{i=1}^d x_i D_{x_i} u(t, x) - ru(t, x)$$

# Performance

	Nonlinear BS
ReLU	0.46% (0.0008)
Tanh	0.44% (0.0006)
Sigmoid	0.46% (0.0004)
Softplus	0.45% (0.0007)

Number of layers <sup>†</sup>	29	58	87	116	145
Mean of relative error	2.29%	0.90%	0.60%	0.56%	0.53%
Std. of relative error	0.0026	0.0016	0.0017	0.0017	0.0014

Han, et al.

Dimension	nb step	value	std	reference
1	10	0.06047	0.00023	0.060903
1	20	0.060789	0.00021	0.060903
1	40	0.061122	0.00015	0.060903
1	80	0.0613818	0.00019	0.060903
5	10	0.10537	0.00014	0.10738
5	20	0.10657	0.00011	0.10738
5	40	0.10725	0.00012	0.10738
5	80	0.107650	0.00016	0.10738
10	10	0.12637	0.00014	0.12996
10	20	0.128292	0.00011	0.12996
10	40	0.12937	0.00014	0.12996
10	80	0.129923	0.00016	0.12996
20	10	0.1443	0.00014	0.1510
20	20	0.147781	0.00012	0.1510
20	40	0.149560	0.00012	0.1510
20	80	0.15050	0.00010	0.1510
40	10	0.15512	0.00018	0.1680
40	20	0.16167	0.00015	0.1680
40	40	0.16487	0.00011	0.1680
40	80	0.16665	0.00013	0.1680
40	160	0.16758	0.00016	0.1680

Huré, et al.

As stated in the paper, “As it is done in python, the global graph creation takes much time as it is repeated for each time step and the global resolution is a little bit time consuming : as the dimension of the problem increases, the time difference decreases and it becomes hard to compare the computational time for a given accuracy when the dimension is above 5”, there is no time provided for computing the example

# Linear Quadratic Gaussian Control

- Consider a stochastic optimal control problem in  $\mathbb{R}^d$

$$dX_t = 2\sqrt{\lambda} m_t dt + \sqrt{2} dW_t$$

$$J(\{m_t\}_{0 \leq t \leq T}) = \mathbb{E} \left[ \int_0^T \|m_t\|_2^2 dt + g(X_T) \right]$$

- the optimal control is

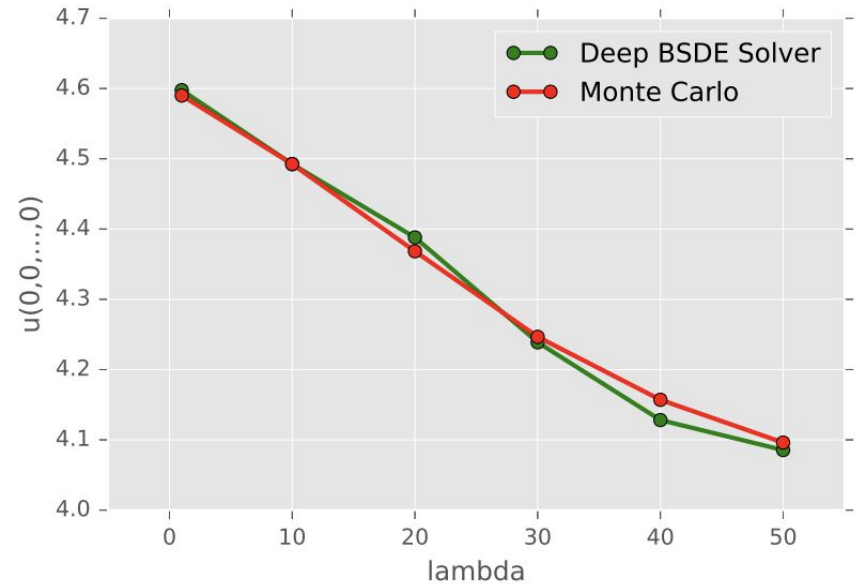
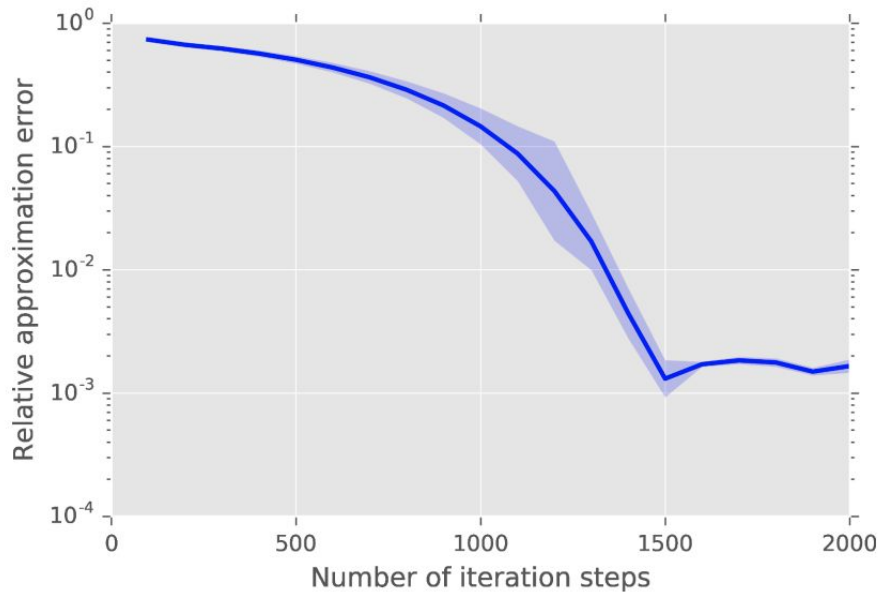
$$m_t^* = \frac{\nabla u(t, x)}{\sqrt{2\lambda}}$$

$$u(t, x) = -\frac{1}{\lambda} \ln \left( \mathbb{E} \left[ \exp \left( -\lambda g(x + \sqrt{2} W_{T-t}) \right) \right] \right)$$

- via the Hamilton-Jacobi-Bellman equation,

$$\frac{\partial u}{\partial t}(t, x) + \Delta u(t, x) - \lambda \|\nabla u(t, x)\|_2^2 = 0$$

# Result



THANK YOU !