

# Convergence of Time-Stepping Deep Gradient Flow Methods

Jasper Rou,<sup>†\*</sup> Chenguang Liu,<sup>†</sup> Antonis Papantoleon<sup>†</sup>

\*j.g.rou@tudelft.nl, <sup>†</sup>Delft Institute of Applied Mathematics, TU Delft

## Introduction

Neural networks are becoming increasingly popular in finance. However, good theoretical results are often lacking. In this research we prove convergence of the time deep gradient flow method, a neural network method which shows good numerical performance in solving partial differential equations (PDEs).

The price of an option can be written as the solution to the PDE:

$$\frac{\partial u}{\partial t} - \sum_{i,j=0}^n a^{ij} \frac{\partial^2 u}{\partial x_i \partial x_j} - \sum_{i=0}^n b^i \frac{\partial u}{\partial x_i} + ru = 0.$$

We split the PDE into a symmetric and an asymmetric part:

$$\begin{aligned} \frac{\partial u}{\partial t} &= \sum_{i,j=0}^n a^{ij} \frac{\partial^2 u}{\partial x_i \partial x_j} + \sum_{i=0}^n b^i \frac{\partial u}{\partial x_i} - ru \\ &= \sum_{i,j=0}^n \frac{\partial}{\partial x_j} \left( a^{ij} \frac{\partial u}{\partial x_i} \right) - \sum_{i,j=0}^n \frac{\partial a^{ij}}{\partial x_j} \frac{\partial u}{\partial x_i} + \sum_{i=0}^n b^i \frac{\partial u}{\partial x_i} - ru \\ &= \sum_{i,j=0}^n \frac{\partial}{\partial x_j} \left( a^{ij} \frac{\partial u}{\partial x_i} \right) - \sum_{i=0}^n \left( \sum_{j=0}^n \frac{\partial a^{ij}}{\partial x_j} - b^i \right) \frac{\partial u}{\partial x_i} - ru \\ &= \nabla \cdot (A \nabla u) - F(u) - ru. \end{aligned}$$

## Time Deep Gradient Flow

First, we discretize the PDE in time. Divide  $[0, T]$  in intervals  $(t_{k-1}, t_k]$  with  $h = t_k - t_{k-1}$  and seek approximations  $U^k$  such that

$$\frac{U^k - U^{k-1}}{h} - \nabla \cdot (A \nabla U^k) + rU^k + F(U^{k-1}) = 0.$$

### Theorem 1: Akrivis and Crouzeix 2004

There exists a constant  $C$  independent of  $h$  and  $k$  such that

$$\max_{0 \leq k \leq N} \|u(t_k) - U^k\| \leq Ch.$$

Second, we rewrite the solution of the discretized PDE as the minimizer of the functional

$$I^k(u) = \frac{1}{2} \|u - U^{k-1}\|^2 + h \int_{\mathbb{R}^d} \frac{1}{2} \left( (\nabla u)^T A \nabla u + ru^2 \right) + F(U^{k-1}) u dx.$$

### Theorem 2:

The minimizer  $w_* \in \mathcal{H}_0^1(\mathbb{R}^d)$  of  $I^k$  is the unique solution  $U^k$ .

### Proof

Define

$$i^k(\tau) = I^k(w_* + \tau v).$$

Since  $w_*$  minimizes  $I^k$ ,  $\tau = 0$  minimizes  $i^k$ . So using integration by parts,

$$\begin{aligned} 0 &= (i^k)'(0) \\ &= \int_{\mathbb{R}^d} \left( (w_* - U^{k-1}) + h \left( -\nabla \cdot (A \nabla w_*) + rw_* + F(U^{k-1}) \right) \right) v dx. \end{aligned}$$

### Theorem 3:

Let  $w_m$  be a sequence in  $\mathcal{H}_0^1(\mathbb{R}^d)$  and  $w_*$  the minimizer of  $I^k$ .

$$\lim_{m \rightarrow \infty} \|w_m - w_*\|_{\mathcal{H}_0^1} = 0 \iff \lim_{m \rightarrow \infty} I^k(w_m) = I^k(w_*).$$

Third, we approximate this minimizer by a neural network.

### Definition 1: Activation function

An activation function is a function  $\psi : \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\psi \in C_c^\infty(\mathbb{R}^d)$  and  $\int_{\mathbb{R}^d} \psi(x) dx \neq 0$ .

### Definition 2: Neural network

Denote the parameters of the neural network by  $\theta^N = (\beta^i, \alpha^i, c^i)_{i=1}^N$  and let  $\frac{1}{2} < \delta < 1$ . The space of neural networks with one layer and  $N$  neurons is

$$\mathcal{C}^N(\psi) = \left\{ V_t^N(\theta^N; x) = V^N(\theta_t^N; x) = N^{-\delta} \sum_{i=1}^N \beta^i \psi(\alpha^i x + c^i) \right\}.$$

### Theorem 4:

$\mathcal{C}(\psi) = \cup_{N \geq 1} \mathcal{C}^N(\psi)$  is dense in  $\mathcal{H}_0^1(\mathbb{R}^d)$ .

## Training

We train the network using gradient descent with learning rate  $\eta_N$ :

$$\begin{aligned} \frac{d\theta_t^N}{dt} &= -\eta_N \nabla_{\theta} I^k(V^N(\theta_t^N; x)), \quad \eta_N = N^{2\delta-1}, \\ \frac{dV_t^N(x)}{dt} &= \nabla_{\theta} V^N(\theta_t^N; x) \cdot \frac{d\theta_t^N}{dt} \\ &= -\eta_N \nabla_{\theta} V^N(\theta_t^N; x) \cdot \nabla_{\theta} I^k(V^N(\theta_t^N; x)) \\ &= -\langle \mathcal{D}I^k(V_t^N), Z_t^N(x, \cdot) \rangle_{\mathcal{H}_0^1}, \\ Z_t^N(x, y) &= N^{-1} \sum_{i=1}^N \nabla_{\beta, \alpha, c} \beta_t^i \psi(\alpha_t^i x + c_t^i) \cdot \nabla_{\beta, \alpha, c} \beta_t^i \psi(\alpha_t^i y + c_t^i). \end{aligned}$$

We can prove that as the number of neurons goes to infinity, the neural network converges to the following gradient flow:

$$\begin{aligned} \frac{dV_t(x)}{dt} &= -\langle \mathcal{D}I^k(V_t), Z(x, \cdot) \rangle_{\mathcal{H}_0^1}, \\ Z(x, y) &= \mathbb{E} \left[ \nabla_{\beta, \alpha, c} \beta_0^1 \psi(\alpha_0^1 x + c_0^1) \cdot \nabla_{\beta, \alpha, c} \beta_0^1 \psi(\alpha_0^1 y + c_0^1) \right]. \end{aligned}$$

### Theorem 5:

For any  $T > 0$ ,

$$\sup_{0 \leq t \leq T} \mathbb{E} \left[ \|V_t^N - V_t\|_{\mathcal{H}_0^1} \right] \xrightarrow{N \rightarrow \infty} 0.$$

Furthermore, we can prove that as the training time goes to infinity, this gradient flow converges to the minimizer of the functional.

### Theorem 6:

$$\lim_{t \rightarrow \infty} \|V_t - w_*\|_{\mathcal{H}_0^1} = 0.$$

### Proof

We define the operator  $\tilde{\mathcal{T}}$  by

$$\begin{aligned} \frac{d(V_t - w_*)}{dt}(x) &= -\langle \mathcal{D}I^k(V_t - w_* + w_*), Z(x, \cdot) \rangle_{\mathcal{H}_0^1} \\ &= -\tilde{\mathcal{T}}(V_t - w_*)(x). \end{aligned}$$

Then  $\tilde{\mathcal{T}}$  is a self-adjoint, positive definite trace class operator. So we can apply a spectral decomposition:

$$\tilde{\mathcal{T}}(\tilde{e}_i) = \lambda_i \tilde{e}_i,$$

with  $\lambda_1 \geq \lambda_2 \geq \dots > 0$  and an orthogonal basis  $\{\tilde{e}_i\}_{i=1}^{\infty}$ .

$$\begin{aligned} \frac{dh_t^i}{dt} &:= \frac{\langle d(V_t - w_*), \tilde{e}_i \rangle}{dt} = -\langle \tilde{\mathcal{T}}(V_t - w_*), \tilde{e}_i \rangle = -\langle V_t - w_*, \tilde{\mathcal{T}}(\tilde{e}_i) \rangle \\ &= -\lambda_i h_t^i. \end{aligned}$$

$h_t^i = e^{-\lambda_i t} h_0^i$ . By Parseval's identity we conclude

$$\|V_t - w_*\|^2 = \sum_{i=1}^{\infty} (h_t^i)^2 = \sum_{i=1}^{\infty} e^{-2\lambda_i t} (h_0^i)^2 \xrightarrow{t \rightarrow \infty} 0.$$