

基于深度学习的间断有限元求解方法

常芸凡¹

指导老师：杨顶辉教授¹、贺茜君教授²

¹清华大学数学科学系

²北京工商大学数学科学系

2021 年 12 月 2 日



- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

研究背景

对于物理空间中具有不连续解的方程，间断有限元(DG)方法精度高、数值频散小且容易处理复杂边界问题 [Cockburn and Shu, 1989]，但是由于CFL条件限制，我们的步长不能选得过大，较小的步长极大的增加了我们的计算负担，导致DG方法的计算速度不高。

目前有很多利用神经网络来快速求解高维偏微分方程的思想 [Weinan and Yu, 2017] [Sirignano and Spiliopoulos, 2018]等，这些方法有效的克服了高维问题的维数灾难。我们希望将深度学习引入到DG方法中来，同时结合神经网络和DG的优势，从而实现DG的加速。

- 1 研究背景
- 2 基于深度学习的间断有限元方法
 - 深度神经网络简介
 - 声波方程的间断有限元方法
 - 神经网络的近似方法
 - 边界条件的处理
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

- 1 研究背景
- 2 基于深度学习的间断有限元方法
深度神经网络简介
声波方程的间断有限元方法
神经网络的近似方法
边界条件的处理
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

DNN简介

Deep neural network(DNN)包含一系列的层，每一层都有若干个神经元与前层和后层神经元相连。神经元通过仿射变换和非线性激活函数连接。该universal approximation theorem(万能近似定理) [Cybenko, 1989]表明，DNN可以近似任何有限维空间的Borel可测函数。

一般而言DNN输入为 $\mathbf{z}^0 = (t, \mathbf{x})$ 输出为 $\mathcal{N}(t, \mathbf{x})$ 。每层神经网络之间的关系为：

$$\mathbf{z}^0 = (t, \mathbf{x}) \quad \text{input}$$

$$\mathbf{z}_k^{l+1} = \sigma_l \left(\mathbf{w}_k^{l+1} \cdot \mathbf{z}^l + b_k^l \right), \quad l = 0, 1, \dots, L-1, \quad 1 \leq k \leq m_{l+1}$$

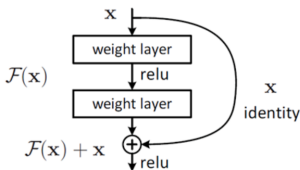
$$\mathcal{N}(t, \mathbf{x}) = \mathbf{w}^{L+1} \mathbf{z}^L \quad \text{output}$$

其中 m_l 是第 l 层的神经元个数， σ 是非线性激活函数。

ResNet

从经验来看，网络的深度对模型的性能至关重要，当增加网络层数后，网络可以进行更加复杂的特征模式的提取，所以当模型更深时理论上可以取得更好的结果(一些深度网络可表示的函数可能需要浅层网络指数级的隐藏单元才能表示 [Montúfar et al., 2014])。但是在实际应用中，深层网络存在着梯度消失或者爆炸的问题，这使得深度学习模型很难训练，效果也不好。

为此 [He et al., 2016]提出了残差学习(ResNet)来解决退化问题



这种网络类似于一种短路连接，保证了深层网络的效果至少不会比浅层网络效果差。

我们尝试的网络结构

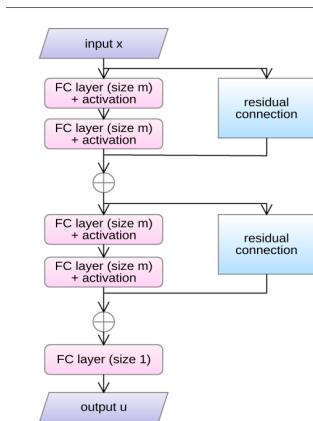


图 1: 残差网络

我们初步尝试的几种网络结构：3或4个残差块(也就是6个隐藏层或8个隐藏层)，每层网络的神经元为20或40。

- 1 研究背景
- 2 基于深度学习的间断有限元方法
 - 深度神经网络简介
 - 声波方程的间断有限元方法
 - 神经网络的近似方法
 - 边界条件的处理
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

我们以二维声波方程为例介绍基于通量函数的 DG 方法。二维声波方程的一般形式为：

$$\frac{\partial^2 u}{\partial t^2} = c^2 \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial z^2} \right) + f \tag{1}$$

其中 c 为声速， $f(t)$ 为震源函数。

声波方程的双曲守恒形式

引入变量 p, q , 使它们分别满足:

$$\frac{\partial p}{\partial t} = \frac{\partial u}{\partial x} \quad , \quad \frac{\partial q}{\partial t} = \frac{\partial u}{\partial z}$$

则方程(1)积分后可改写为:

$$\frac{\partial W}{\partial t} + \nabla \cdot F(W) = \hat{f}, \tag{2}$$

$$\text{其中 } W = (u, p, q)^T, F(W) = \begin{pmatrix} -c^2 p & -c^2 q \\ -u & 0 \\ 0 & -u \end{pmatrix}, \hat{f} = \begin{pmatrix} f \\ 0 \\ 0 \end{pmatrix}$$

空间离散

设区域为 Ω , 首先将 Ω 划分为互不相交的单元 $\{\Omega_i\}$ 。我们采用的试验函数空间为:

$$V_h = \left\{ v \in L^1(\Omega) : v|_{\Omega_i} \in P^k(\Omega_i) \right\}$$

其中 $P^k(\Omega_i)$ 表示单元 Ω_i 上至多 k 次多项式。

由变分法, 我们将问题转化为

$$\int_{\Omega_i} \left(v \frac{\partial \mathbf{W}}{\partial t} - \mathbf{F}(\mathbf{W}) \cdot \nabla v \right) dV + \int_{\partial \Omega_i} v \hat{\mathbf{F}}(\mathbf{W}_i, \mathbf{W}_e, \mathbf{n}) dS = \int_{\Omega_i} \hat{\mathbf{f}}_v dV \tag{3}$$

其中 $\hat{\mathbf{F}}(\mathbf{W}_i, \mathbf{W}_e, \mathbf{n})$ 表示依赖于边界内外两侧数值通量。

数值通量和基函数的选择

我们采用局部Lax-Friedrichs通量：

$$\hat{F}^{LLF}(a, b, n) = \frac{1}{2}(F(a) + F(b)) \cdot n - \frac{C}{2}(b - a)$$

其中 $C = \max_{\min(a,b) \leq s \leq \max(a,b)} |F'(s) \cdot n|$ 。

基函数方面，我们选取正交Legendre基函数：

$$\phi_0(x) = 1$$

$$\phi_1(x) = \frac{2}{\Delta x} \left(x - x_{i+\frac{1}{2}} \right)$$

$$\phi_2(x) = \frac{6}{\Delta x^2} \left(x - x_{i+\frac{1}{2}} \right)^2 - \frac{1}{2}$$

...

TVD时间离散格式

一般的TVD时间离散为：

$$\begin{aligned} \left(u^h\right)^{(i)} &= \sum_{\ell=0}^{i-1} \left[\alpha_{i\ell} \left(u^h\right)^{(\ell)} + \beta_{i\ell} \Delta t L_h \left(\left(u^h\right)^{(\ell)}, t^n + d_{\ell} \Delta t \right) \right], \\ i &= 1, \dots, r \end{aligned}$$

其中， $\left(u^h\right)^{(0)} = \left(u^h\right)^n$ ， $\left(u^h\right)^{(r)} = \left(u^h\right)^{n+1}$

我们常用的两种格式为：

- 2阶格式(r=2):

$$\alpha_{10} = \beta_{10} = 1, \alpha_{20} = \alpha_{21} = \beta_{21} = \frac{1}{2}, \beta_{20} = 0; d_0 = 0, d_1 = 1;$$

- 3阶格式(r=3):

$$\begin{aligned} \alpha_{10} = \beta_{10} = 1, \alpha_{20} = \frac{3}{4}, \beta_{20} = 0, \alpha_{21} = \beta_{21} = \frac{1}{4}, \alpha_{30} = \\ \frac{1}{3}, \beta_{30} = \alpha_{31} = \beta_{31} = 0, \alpha_{32} = \beta_{32} = \frac{2}{3}; d_0 = 0, d_1 = 1, d_2 = \\ \frac{1}{2}; \end{aligned}$$

具体而言，二阶精度的TVD时间离散格式为：

$$C^{(0)} = C^{(n)}$$

$$C^{(1)} = C^{(0)} + \Delta t L \left(C^{(0)} \right)$$

$$C^{(2)} = \frac{1}{2} C^{(0)} + \frac{1}{2} C^{(1)} + \frac{1}{2} \Delta t L \left(C^{(1)} \right)$$

$$C^{(n+1)} = C^{(2)}$$

三阶精度的TVD时间离散格式为：

$$C^{(0)} = C^{(n)}$$

$$C^{(1)} = C^{(0)} + \Delta t L \left(C^{(0)} \right)$$

$$C^{(2)} = \frac{3}{4} C^{(0)} + \frac{1}{4} C^{(1)} + \frac{1}{4} \Delta t L \left(C^{(1)} \right)$$

$$C^{(2)} = \frac{1}{3} C^{(0)} + \frac{2}{3} C^{(2)} + \frac{2}{3} \Delta t L \left(C^{(2)} \right)$$

$$C^{(n+1)} = C^{(3)}$$

① 研究背景

② 基于深度学习的间断有限元方法

深度神经网络简介

声波方程的间断有限元方法

神经网络的近似方法

边界条件的处理

③ 具体实现中的加速技巧和优势

④ 数值算例

⑤ 下一步工作安排

⑥ 参考文献

我们选取基函数 φ_i^j 为单元 I_i 上的 j 阶正交Legendre多项式(在其他单元上该函数值为0)，则DG方法最终的解可以写成

$$W_h(t, \mathbf{x}) = \sum_{j=0}^K \sum_{i=0}^M C_i^j(t) \varphi_i^j(x)$$

我们可以利用神经网络来近似系数 C_i^j ，具体而言，我们把解写成如下形式 [Chen et al., 2021]:

$$W_{h,\theta}(t, \mathbf{x}) = \sum_{j=0}^K \sum_{i=0}^M \mathcal{N}_{\theta}^j(t, x_{i+1/2}) \varphi_i^j(x) \quad (4)$$

这里我们一共用了 $K + 1$ 个神经网络 $\mathcal{N}_{\theta}^0, \mathcal{N}_{\theta}^1, \dots, \mathcal{N}_{\theta}^K$ ，每个神经网络的输入是时间 t 和空间坐标 x ，我们希望神经网络的输出 $\mathcal{N}_{\theta}^j(t, x_{i+1/2}) = C_i^j(t)$

声波方程推导

把神经网络的格式(4)带入声波方程的DG格式(3)，我们可以得到：

$$\begin{aligned} \frac{\partial \mathcal{N}_{\theta}^j(t, x_{i+1/2})}{\partial t} = & (\varphi_i^j(x), \varphi_i^j(x))^{-1} \left[\int_{\Omega_i} \mathbf{F}(\mathbf{W}) \cdot \nabla \varphi_i^j(x) \right. \\ & \left. - \int_{\partial \Omega_i} \varphi_i^j(x) \hat{\mathbf{F}}(\mathbf{W}_i, \mathbf{W}_e, \mathbf{n}) dS + \int_{\Omega_i} \hat{\mathbf{f}}_v dV \right] \end{aligned} \quad (5)$$

我们把方程简单的记作：

$$\frac{\partial \mathcal{N}_{\theta}^j(t, x_{i+1/2})}{\partial t} = L(\mathcal{N}_{\theta}(t, x)) \quad (6)$$

声波方程TVD时间离散推导

由方程(6)，我们利用TVD时间离散，如2阶格式：

$$C^{(0)} = \mathcal{N}_{\theta}^j(t_n, x_{i+1/2})$$

$$C^{(1)} = C^{(0)} + \Delta t L \left(C^{(0)} \right)$$

$$C^{(2)} = \frac{1}{2} C^{(0)} + \frac{1}{2} C^{(1)} + \frac{1}{2} \Delta t L \left(C^{(1)} \right)$$

$$C^{(n+1)} = C^{(2)}$$

按照传统方法，我们计算的 $C^{(n+1)}$ 应该是 $n+1$ 时间层的系数，也就是 $C^{(n+1)} = \mathcal{N}_{\theta}^j(t_{n+1}, x_{i+1/2})$ 应当成立，所以我们设

$$L_{i,j,n} = C^{(n+1)} - \mathcal{N}_{\theta}^j(t_{n+1}, x_{i+1/2})$$

相应的损失函数设为 $\mathcal{L}(\theta) = \left(\sum_{i,j,n} L_{i,j,n}^2 \right)^{1/2}$

- 1 研究背景
- 2 基于深度学习的间断有限元方法
 - 深度神经网络简介
 - 声波方程的间断有限元方法
 - 神经网络的近似方法
 - 边界条件的处理
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

常见的边界条件有:

- Dirichlet:

$$u(t, \mathbf{x}) = g(t, \mathbf{x}) \quad \mathbf{x} \in \partial D$$

- Neumann:

$$\frac{\partial u(t, \mathbf{x})}{\partial \nu} = g(t, \mathbf{x}) \quad \mathbf{x} \in \partial D$$

- 初值条件:

$$u(0, \mathbf{x}) = u_0(\mathbf{x})$$

两种处理方法

- 最直接的想法是把边界条件直接加入到损失函数里进行训练，直接让其收敛于0，比如Dirichlet边界条件我们可以在损失函数中加上 $\lambda \|u - g\|_{\partial D}^2$ 一项，其中 λ 是惩罚因子(超参数，可自行设置)
- 另一种思路是构造DNN网络使得其精确的满足边界条件。例如对于初值问题 $u(0, \mathbf{x}) = u_0(\mathbf{x})$ ，我们可以构造表达式

$$u_\theta(t, \mathbf{x}) = t\mathcal{N}_\theta(t, \mathbf{x}) + u_0(\mathbf{x})$$

- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

小批量训练

我们的神经网络 \mathcal{N}_θ^j 训练的是所有 j 阶Legendre多项式前的系数，并且相比于传统DG方法 $u_h(t, \mathbf{x}) = \sum_{j=0}^K \sum_{i=0}^M C_i^j(t) \varphi_i^j(\mathbf{x})$ ，当阶数 j 固定的时候， $C_i^j(t)$ 在每个网格上是相同的值，相当于是离散的点。

但是神经网络训练的 $\mathcal{N}_\theta^j(t, \mathbf{x})$ 实际上关于 t 和 \mathbf{x} 都是可连续计算的。相比于传统DG方法每个时间步都需要计算所有网格点前基函数的系数，神经网络的方法却不需要每次遍历所有网格，我们每次在 (t, \mathbf{x}) 组成的网格点中随机抽样(如10000个点)进行训练，这样既减少了我们每次训练的数量，又可知当空间维数变成3维时，由于我们每次依然随机抽样10000个点，每次训练的计算量其实并不会增加。

并行处理

我们每个神经网络的损失函数 $\mathcal{L}^j(\theta)$ 基本不相关，所以我们可以相对独立的同时训练每个网络，这意味着增加阶数并不会对我们的计算速度产生很多影响。

从数学推导中我们也可以发现，DG格式中，采用 n 阶正交多项式逼近和采用 $n+1$ 阶正交多项式逼近，其前 n 阶的系数是相同的。换言之，若 $K+1$ 阶DG格式的解是

$$u_h(t, \mathbf{x}) = \sum_{j=0}^{K+1} \sum_{i=0}^M c_i^j(t) \varphi_i^j(\mathbf{x})$$

则 K 阶DG格式的解就是

$$u_h(t, \mathbf{x}) = \sum_{j=0}^K \sum_{i=0}^M c_i^j(t) \varphi_i^j(\mathbf{x})$$

这意味着本身DG格式不同阶数的Legendre多项式前的系数是独立的，而我们的方法正是可以独立的同时近似这些系数。

优势

总得来说，相比于传统方法，神经网络近似有如下的优势

- 如果增加空间维度，由于我们每次训练都是抽样小批量训练，所以每次训练的计算量并不会增加(但是训练次数可能会增加);
- 如果增加近似阶数(使用更高阶的Legendre多项式)，我们可以采用并行处理，同时训练多个神经网络，也不会降低我们的计算速度。

但是我们的方法也只是近似DG格式中需要计算的 $C_i^j(t)$ ，所以肯定达不到DG方法的精度。但是神经网络的优势在于，它可以在很少的训练次数内达到较为满意的精度，从而实现对DG方法的加速。

- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

1维声波方程

我们考虑如下声波方程：

$$\begin{cases} \frac{\partial^2 u}{\partial t^2} - c^2 \frac{\partial^2 u}{\partial x^2} = 0 \\ u(x, 0) = \cos\left(-\frac{2\pi f_0}{c}x\right) \\ \frac{\partial u(x, 0)}{\partial t} = -2\pi f_0 \sin\left(-\frac{2\pi f_0}{c}x\right) \end{cases}$$

该方程有解析解：

$$u(t, x) = \cos\left(2\pi f_0 t - \frac{2\pi f_0}{c}x\right)$$

神经网络近似

0阶数值解我们如下构造：

$$u_{\theta}(t, x) = \begin{cases} \mathcal{N}_{\theta}\left(t, x_{i+\frac{1}{2}}\right) \varphi_i(x) & t > 0 \quad x \in [x_i, x_{i+1}) \\ \cos\left(-\frac{2\pi f_0}{c} x_{i+\frac{1}{2}}\right) & t = 0 \quad x \in [x_i, x_{i+1}) \end{cases}$$

$$q_{\theta}(t, x) = \begin{cases} \mathcal{M}_{\theta}\left(t, x_{i+\frac{1}{2}}\right) \varphi_i(x) & t > 0 \quad x \in [x_i, x_{i+1}) \\ -2\pi f_0 \sin\left(-\frac{2\pi f_0}{c} x_{i+\frac{1}{2}}\right) & t = 0 \quad x \in [x_i, x_{i+1}) \end{cases}$$

$$\text{其中 } \varphi_i(x) = \begin{cases} 1 & x \in [x_i, x_{i+1}) \\ 0 & \text{otherwise} \end{cases}$$

0阶格式的结果($N=320$)

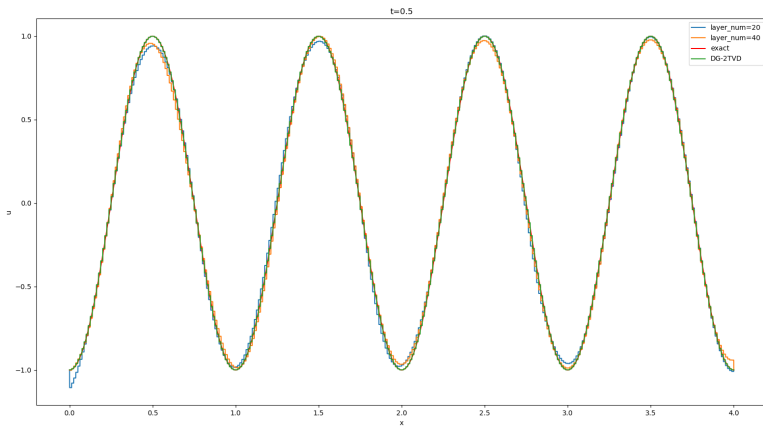


图 2: $N=320, t=0.5$

0阶格式的结果($N=320$)

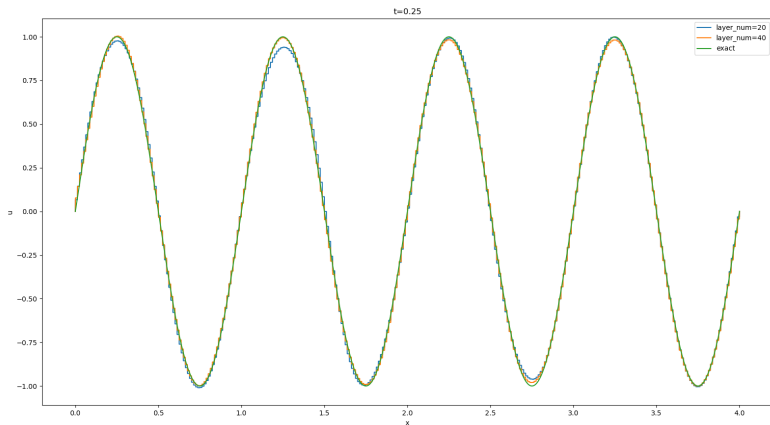


图 3: $N=320, t=0.25$

0阶格式的结果($N=320$)

$N = 320$	layer-num=20	layer-num=40	传统DG方法
训练步数	2000	3000	
$t = 0.25s$ 相对误差	3.758%	2.667%	
$t = 0.5s$ 相对误差	4.367%	4.049%	3.388%
运行时间	420s	1170s	120s

0阶格式的结果(N=640)

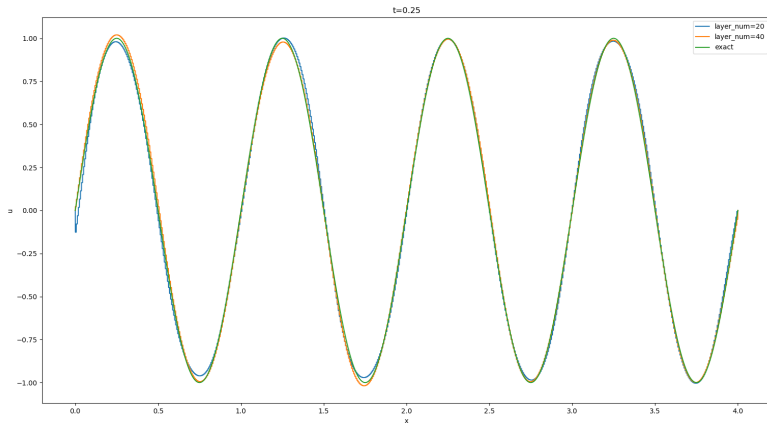


图 4: $N=640, t=0.25$

0阶格式的结果($N=640$)

$N = 640$	layer-num=20	layer-num=40	传统DG方法
训练步数	2500	3200	
$t = 0.25s$ 相对误差	4.310%	2.874%	
$t = 0.5s$ 相对误差	6.605%	3.693%	2.105%
运行时间	500s	1248s	462s

0阶格式的结果($N=1280$)

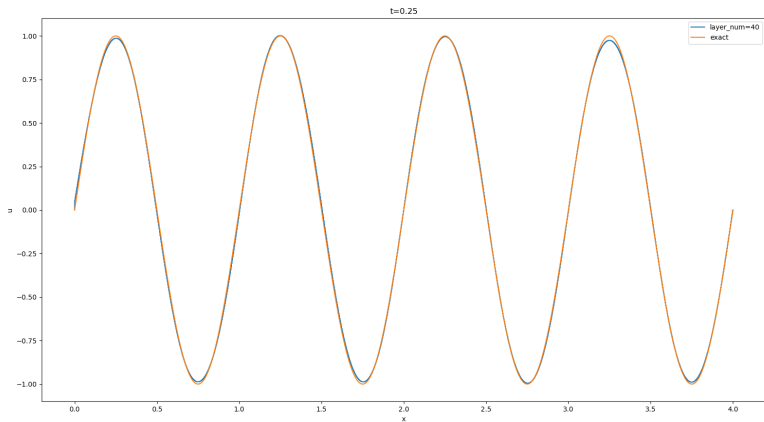


图 5: $N=1280, t=0.25$

0阶格式的结果($N=1280$)

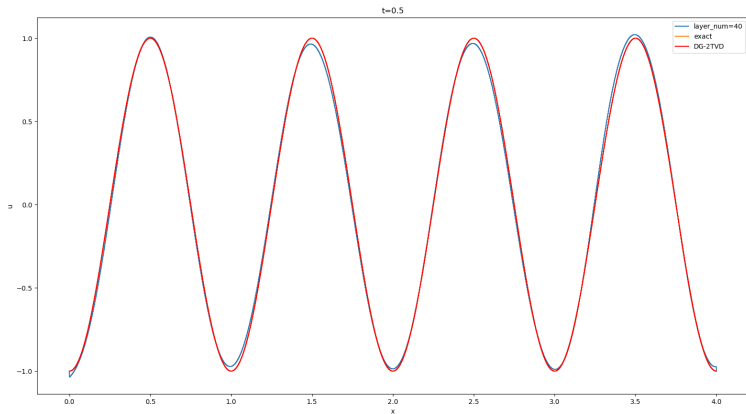


图 6: $N=1280, t=0.5$

0阶格式的结果($N=1280$)

$N = 1280$	layer-num=40	传统DG方法
训练步数	3000	
$t = 0.25s$ 相对误差	1.376%	
$t = 0.5s$ 相对误差	2.801%	1.374%
运行时间	1340s	2416s

- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

下一步工作安排

- ① 至少做到空间2阶格式;
- ② 向二维、三维空间尝试

- 1 研究背景
- 2 基于深度学习的间断有限元方法
- 3 具体实现中的加速技巧和优势
- 4 数值算例
- 5 下一步工作安排
- 6 参考文献

参考文献 I

- [Chen et al., 2021] Chen, J., Jin, S., and Lyu, L. (2021).
A deep learning based discontinuous galerkin method for hyperbolic equations with discontinuous solutions and random uncertainties.
arXiv preprint arXiv:2107.01127.
- [Cockburn and Shu, 1989] Cockburn, B. and Shu, C.-W. (1989).
Tvb runge-kutta local projection discontinuous galerkin finite element method for conservation laws ii: General framework.
Mathematics of Computation, 52(186):411–435.
- [Cybenko, 1989] Cybenko, G. (1989).
Approximation by superpositions of a sigmoidal function.
Mathematics of control, signals and systems, 2(4):303–314.
- [He et al., 2016] He, K., Zhang, X., Ren, S., and Sun, J. (2016).
Deep residual learning for image recognition.
In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778.
- [Montúfar et al., 2014] Montúfar, G., Pascanu, R., Cho, K., and Bengio, Y. (2014).
On the number of linear regions of deep neural networks.
arXiv preprint arXiv:1402.1869.

参考文献 II

- [Sirignano and Spiliopoulos, 2018] Sirignano, J. and Spiliopoulos, K. (2018).
Dgm: A deep learning algorithm for solving partial differential equations.
Journal of Computational Physics, 375:1339–1364.
- [Weinan and Yu, 2017] Weinan, E. and Yu, T. (2017).
The deep ritz method: A deep learning-based numerical algorithm for solving
variational problems.
Communications in Mathematics and Statistics, 6:1–12.

Thanks!