## Analysis of Neural Networks (2)

#### Jinchao Xu 许进超

Penn State University

xu@math.psu.edu http://www.math.psu.edu/xu/

CUHK, July 2, 2021

Acknowledgement: NSF: DMS-1819157

# Analysis of Neural Network

Jinchao Xu

This lecture series will provide some mathematical understanding of neural networks and machine learning by studying their close relationship with classic numerical methods such as finite element and multigrid methods. Applications will be given to image classification and numerical partial differential equations

#### Finite Element Connection and Approximation Theory

- ReLU neural networks = linear finite elements
- Largest function class that a stable neural network can approximate
- Optimal approximation rates for popular neural networks

#### 2 Multigrid and Image Classification

- Linear separable sets and logistic regression
- A model for feature extractions
- Image classification by multigrid method

#### Neural Network and Numerical PDEs

- Error analysis of neural network for numerical PDEs
- Numerical quadrature and Rademacher complexity analysis
- Training algorithm that achieves the best asymptotic convergence rate

## Introduction

- - Linear separable sets and logistic regression Nonlinear models

- - A simple example
  - Finite element method
  - Multigrid method
- - Multigrid as a CNN: MgNet

## A basic AI problem: classification

• Can a machine (function) tell the difference ?



### Supervised learning \leftrightarrow Function interpolation

- Each image = a big vector of pixel values
  - ▶  $n = 1280 \times 720 \times 3$  (width× height × RGB channel)  $\approx$  3M.



## Supervised learning \leftrightarrow Function interpolation

- Each image = a big vector of pixel values
  - $n = 1280 \times 720 \times 3$  (width× height × RGB channel)  $\approx 3$ M.
- Want to find a function  $f(\cdot, \Theta) : \mathbb{R}^n \to \mathbb{R}^3$  which separates the classes



## Supervised learning $\leftrightarrow$ Function interpolation

Each image = a big vector of pixel values

- $n = 1280 \times 720 \times 3$  (width× height × RGB channel)  $\approx 3$ M.
- Want to find a function  $f(\cdot, \Theta) : \mathbb{R}^n \to \mathbb{R}^3$  which separates the classes



• Mathematical formulation: Find  $f(\cdot; \Theta) : \mathbb{R}^n \to \mathbb{R}^3$  such that:

$$f(\underbrace{\bullet}) \approx \begin{pmatrix} 1\\0\\0 \end{pmatrix} \quad f(\underbrace{\bullet}) \approx \begin{pmatrix} 0\\1\\0 \end{pmatrix} \quad f(\underbrace{\bullet}) \approx \begin{pmatrix} 0\\0\\1 \end{pmatrix}$$

- Given data:  $\{x_i, y_i\}_{i=1}^N$ 
  - Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
  - Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
  - Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
  - ▶ ...

- Given data:  $\{x_i, y_i\}_{i=1}^N$ 
  - Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
  - Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
  - Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
  - ► ...
- How can we find  $f(\cdot, \Theta)$ ?.

- Given data:  $\{x_i, y_i\}_{i=1}^N$ 
  - Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
  - Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
  - Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
  - ► ...
- How can we find  $f(\cdot, \Theta)$ ?.
- Mathematically: solve the optimization problem

$$\min_{\Theta} \sum_{i=1}^{N} \ell(f(x_i, \Theta), y_i)$$

where  $\ell$  is an appropriate loss function.

(1)

• Given data:  $\{x_i, y_i\}_{i=1}^N$ 

- Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
- Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
- Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
- ► ...
- How can we find  $f(\cdot, \Theta)$ ?.
- Mathematically: solve the optimization problem

$$\min_{\Theta} \sum_{i=1}^{N} \ell(f(x_i, \Theta), y_i)$$

where  $\ell$  is an appropriate loss function.

• Application (generalization): image classification:

(1)

• Given data:  $\{x_i, y_i\}_{i=1}^N$ 

- Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
- Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
- Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
- ► ...
- How can we find  $f(\cdot, \Theta)$ ?.
- Mathematically: solve the optimization problem

$$\min_{\Theta} \sum_{i=1}^{N} \ell(f(x_i, \Theta), y_i)$$

where  $\ell$  is an appropriate loss function.

• Application (generalization): image classification:

$$f(\boxed{?};\Theta) = \begin{pmatrix} 0.7\\ 0.2\\ 0.1 \end{pmatrix}$$

(1)

• Given data:  $\{x_i, y_i\}_{i=1}^N$ 

- Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
- Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
- Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
- ► ...
- How can we find  $f(\cdot, \Theta)$ ?.
- Mathematically: solve the optimization problem

$$\min_{\Theta} \sum_{i=1}^{N} \ell(f(x_i, \Theta), y_i)$$

(1)

where  $\ell$  is an appropriate loss function.

• Application (generalization): image classification:

$$f(\fbox{?};\Theta) = \begin{pmatrix} 0.7\\ 0.2\\ 0.1 \end{pmatrix} \implies \fbox{?} =$$

• Given data:  $\{x_i, y_i\}_{i=1}^N$ 

- Image Classification: x<sub>i</sub> image, y<sub>i</sub> object class
- Speech Recognition: x<sub>i</sub> audio, y<sub>i</sub> text
- Face Recognition: x<sub>i</sub> graph with face, y<sub>i</sub> identity
- ► ...
- How can we find  $f(\cdot, \Theta)$ ?.
- Mathematically: solve the optimization problem

$$\min_{\Theta} \sum_{i=1}^{N} \ell(f(x_i, \Theta), y_i)$$

(1)

where  $\ell$  is an appropriate loss function.

• Application (generalization): image classification:

$$f(\fbox{?};\Theta) = \begin{pmatrix} 0.7\\ 0.2\\ 0.1 \end{pmatrix} \implies \fbox{?} = \operatorname{cat}$$

# Training and Generalization Error

- Training Error:
  - Fit the training dataset
  - Evaluate error or loss

# Training and Generalization Error

- Training Error:
  - Fit the training dataset
  - Evaluate error or loss
- Generalization Error:
  - Performance in the real world
  - Theoretical, cannot be measured exactly

# Overfitting

• Fit noisy function values



- Training error is 0
- Generalization Error is large

Jinchao Xu

## Regularization

- How can we prevent overfitting?
  - Overfitting occurs because the model class is to large relative to the data
  - Need to restrict the model class *F* in some way
- Typical way of restricting the model class:
  - Add a regularization term  $R(\theta)$  to the loss function
- Typical regularizers:
  - $\ell^2$ -loss:  $R(\theta) = ||\theta||_2^2$
  - $\ell^1$ -loss:  $R(\theta) = \|\theta\|_1$
- In deep learning SGD acts as a regularizer!
  - Deep learning model and training algorithm are tied together.

# Bounding the Generalization Error

Main Tool:

Uniform law of large numbers<sup>1</sup>:

$$\mathbb{P}\left(\max_{f\in\mathcal{F}}\left|\frac{1}{N}\sum_{i=1}^{N}L(f(x_i),y_i)-\mathbb{E}L(f(x),y)\right| > \epsilon(N)\right) < 1-\delta(N).$$
(2)

• If the model class  $\mathcal{F}$  is small,  $\epsilon(N), \delta(N) \rightarrow 0$  very quickly

- Main tools for bounding this
  - Rademacher Complexity
  - VC dimension
  - Metric Entropy (more advanced)

#### <sup>1</sup>vapnik2013nature.

Jinchao Xu

## Generalization Error (cont.)

The Rademacher complexity of a class of function *F* on Ω

$$R_N(F) = \mathbb{E}_{x_1, \dots, x_N} \mathbb{E}_{\xi_1, \dots, \xi_N} \left( \sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N \xi_i f(x_i) \right),$$
(3)

where  $x_i$  are drawn uniformly at random from  $\Omega$  and  $\xi_i$  are uniformly random signs.

- Can use the Rademacher complexity to obtain a uniform law of large numbers<sup>2</sup>
- Rademacher complexity of neural networks  $O(N^{-\frac{1}{2}})^{34}$

<sup>2</sup>vapnik2013nature.

<sup>3</sup>ma2019barron.

<sup>4</sup>hong2021rademacher.

Jinchao Xu

#### Introduction



#### Data classification

- Linear separable sets and logistic regressionNonlinear models
- 3 Convolutional neural networks (CNN)
- A model for feature extraction
- 5 Multigrid method: a brief introduction
  - A simple example
  - Finite element method
  - Multigrid method
  - Image as a finite element function
- 6 MgNet
  - Multigrid as a CNN: MgNet
- 7 Concluding remarks

## Linearly separable sets: 2-class

#### Definition

We say  $A_1$ ,  $A_2 \subset \mathbb{R}^n$  are linearly separable, if there exists a hyperplane

$$H_0 = \{x : wx + b = 0\},\$$

such that wx + b > 0 if  $x \in A_1$  and wx + b < 0 if  $x \in A_2$ .



(4)

# Linearly separable sets: 2-class

#### Definition

We say  $A_1$ ,  $A_2 \subset \mathbb{R}^n$  are linearly separable, if there exists a hyperplane

$$H_0 = \{x : wx + b = 0\},\$$

such that wx + b > 0 if  $x \in A_1$  and wx + b < 0 if  $x \in A_2$ .



Figure: Two linearly separable sets

(4)

## Linearly separable sets: k-class

#### Definition (Linearly Separable)

A collection of subsets  $A_1, ..., A_k \subset \mathbb{R}^n$  are linearly separable if there exist

$$W = \begin{pmatrix} W_1 \\ \vdots \\ W_k \end{pmatrix} \in \mathbb{R}^{k \times n}, b = \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} \in \mathbb{R}^{k \times n},$$
(5)

such that, for each  $1 \le i \le k$  and  $j \ne i$ , for any  $x \in A_i$ 

$$w_i x + b_i > w_j x + b_j$$
, or  $(w_i - w_j) x + b_i - b_j > 0$ .

(6)

### Linearly separable sets: k-class

#### Definition (Linearly Separable)

A collection of subsets  $A_1, ..., A_k \subset \mathbb{R}^n$  are linearly separable if there exist

$$W = \begin{pmatrix} w_1 \\ \vdots \\ w_k \end{pmatrix} \in \mathbb{R}^{k \times n}, b = \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} \in \mathbb{R}^{k \times n},$$
(5)

such that, for each  $1 \le i \le k$  and  $j \ne i$ , for any  $x \in A_i$ 

$$w_i x + b_i > w_j x + b_j$$
, or  $(w_i - w_j) x + b_i - b_j > 0.$  (6)

Special case: k = 2

$$(w, b) = (w_1 - w_2, b_1 - b_2).$$

### Linearly separable sets: k-class

#### Definition (Linearly Separable)

A collection of subsets  $A_1, ..., A_k \subset \mathbb{R}^n$  are linearly separable if there exist

$$W = \begin{pmatrix} W_1 \\ \vdots \\ W_k \end{pmatrix} \in \mathbb{R}^{k \times n}, b = \begin{pmatrix} b_1 \\ \vdots \\ b_k \end{pmatrix} \in \mathbb{R}^{k \times n},$$
(5)

such that, for each  $1 \le i \le k$  and  $j \ne i$ , for any  $x \in A_i$ 

$$w_i x + b_i > w_j x + b_j$$
, or  $(w_i - w_j) x + b_i - b_j > 0.$  (6)

Special case: k = 2

$$(w, b) = (w_1 - w_2, b_1 - b_2).$$

Examples of linear separable sets



Jinchao Xu

Soft-max:

softmax(x)

Soft-max:

$$softmax(x) = \frac{1}{\sum\limits_{i=1}^{k} e^{x_i}} \begin{pmatrix} e^{x_1} \\ e^{x_2} \\ \cdots \\ e^{x_k} \end{pmatrix},$$

Soft-max:

$$softmax(x) = rac{1}{\sum\limits_{i=1}^{k} e^{x_i}} \begin{pmatrix} e^{x_1} \\ e^{x_2} \\ \cdots \\ e^{x_k} \end{pmatrix},$$

Likelyhood function

$$P(\theta) = \prod_{i=1}^{k} \prod_{x \in A_i} p_i(\theta, x), \quad 0 < p_i(\theta, x) < 1, \quad \sum_{i=1}^{k} p_i(\theta, x) = 1.$$

$$(8)$$

where  $\boldsymbol{p}(\theta, x) = softmax(Wx + b)$  and  $\boldsymbol{p}(\theta, x) = (p_1(\theta, x), \cdots, p_k(\theta, x))$ .

(7)

Loss function

 $\theta_{\lambda} = (W_{\lambda}, b_{\lambda})$ 

Loss function

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\max_{ heta} [P( heta)]$$

Loss function

here

$$\theta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg \max_{\theta} [P(\theta)] = \arg \min_{\theta} \left( -\log P(\theta) \right), \tag{9}$$

$$-\log P(\theta) = -\sum_{j=1}^{N} y_j \cdot \log(\boldsymbol{p}(\theta, x_j)) \Leftrightarrow \text{cross-entropy.}$$
(10)

Classification

lf

 $p_i(\theta_\lambda, x_j) > p_m(\theta_\lambda, x_j), \quad \forall m \neq i,$  (11)

then

$$x_j \in A_i.$$
 (12)

### Existence

Consider:



#### Existence

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{ heta} \left( \sum_{j=1}^{N} \ell(y_j, p(\theta, x_j)) + \lambda R(\theta) \right)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

(13)
Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{ heta} \left( \sum_{j=1}^{N} \ell(y_j, \boldsymbol{p}(\theta, x_j)) + \lambda \boldsymbol{R}(\theta) \right)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{ heta} \left( \sum_{j=1}^{N} \ell(y_j, p( heta, x_j)) + \lambda R( heta) 
ight)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

If  $\lambda > 0$  is sufficiently small,  $W_{\lambda}x + b_{\lambda}$  classifies the sets correctly.

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{\theta} \left( \sum_{j=1}^{N} \ell(y_j, \boldsymbol{p}(\theta, x_j)) + \lambda R(\theta) \right)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

If  $\lambda > 0$  is sufficiently small,  $W_{\lambda}x + b_{\lambda}$  classifies the sets correctly.

2 If  $\lambda = 0$ , no global miminum

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{\theta} \left( \sum_{j=1}^{N} \ell(y_j, \boldsymbol{p}(\theta, x_j)) + \lambda \boldsymbol{R}(\theta) \right)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

If  $\lambda > 0$  is sufficiently small,  $W_{\lambda}x + b_{\lambda}$  classifies the sets correctly.

2 If  $\lambda = 0$ , no global miminum

Comments:

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{ heta} \left( \sum_{j=1}^{N} \ell(y_j, \boldsymbol{p}(\theta, x_j)) + \lambda R(\theta) \right)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

If  $\lambda > 0$  is sufficiently small,  $W_{\lambda}x + b_{\lambda}$  classifies the sets correctly.

2 If  $\lambda = 0$ , no global miminum

Comments:

•  $-\log P(\theta) \Leftrightarrow \text{cross-entropy.}$ 

Consider:

$$heta_{\lambda} = (W_{\lambda}, b_{\lambda}) = \arg\min_{ heta} \left( \sum_{j=1}^{N} \ell(y_j, oldsymbol{p}( heta, x_j)) + \lambda R( heta) 
ight)$$

with  $\ell(y, p(\theta, x)) = -y \cdot \log(p(\theta, x))$ 

Assume that  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$  is linearly separable,

If  $\lambda > 0$  is sufficiently small,  $W_{\lambda}x + b_{\lambda}$  classifies the sets correctly.

2 If  $\lambda = 0$ , no global miminum

#### Comments:

- $-\log P(\theta) \Leftrightarrow \text{cross-entropy.}$
- Regularization is necessary even for convex model for linearly separable sets!

### Definition

Given data  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$ , if there exist a continuous nonlinear (feature) map  $\phi(\cdot, \theta) : \mathbb{R}^d \mapsto \mathbb{R}^n$  so that

$$\{\phi(\mathbf{x}_i,\theta)\}_{i=1}^N = \bigcup_{i=1}^k \phi(\mathbf{A}_i,\theta)$$

is linearly separable, then  $\{x_j\}_{j=1}^N$  is called nonlinearly separable.

### Definition

Given data  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$ , if there exist a continuous nonlinear (feature) map  $\phi(\cdot, \theta) : \mathbb{R}^d \mapsto \mathbb{R}^n$  so that

$$\{\phi(\mathbf{x}_i,\theta)\}_{i=1}^N = \bigcup_{i=1}^{\kappa} \phi(\mathbf{A}_i,\theta)$$

is linearly separable, then  $\{x_j\}_{j=1}^N$  is called nonlinearly separable.

### Definition

Given data  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$ , if there exist a continuous nonlinear (feature) map  $\phi(\cdot, \theta) : \mathbb{R}^d \mapsto \mathbb{R}^n$  so that

$$\{\phi(\mathbf{x}_i,\theta)\}_{i=1}^N = \bigcup_{i=1}^{\kappa} \phi(\mathbf{A}_i,\theta)$$

is linearly separable, then  $\{x_j\}_{j=1}^N$  is called nonlinearly separable.



### Definition

Given data  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$ , if there exist a continuous nonlinear (feature) map  $\phi(\cdot, \theta) : \mathbb{R}^d \mapsto \mathbb{R}^n$  so that

$$\{\phi(\mathbf{x}_i,\theta)\}_{i=1}^N = \bigcup_{i=1}^{\kappa} \phi(\mathbf{A}_i,\theta)$$

is linearly separable, then  $\{x_j\}_{j=1}^N$  is called nonlinearly separable.



### Definition

Given data  $\{x_j\}_{j=1}^N = \bigcup_{i=1}^k A_i$ , if there exist a continuous nonlinear (feature) map  $\phi(\cdot, \theta) : \mathbb{R}^d \mapsto \mathbb{R}^n$  so that

$$\{\phi(\mathbf{x}_i,\theta)\}_{i=1}^N = \bigcup_{i=1}^{\kappa} \phi(\mathbf{A}_i,\theta)$$

is linearly separable, then  $\{x_j\}_{j=1}^N$  is called nonlinearly separable.



## Neural network classification models

#### Theorem

If data  $\{x_j\}_{j=1}^N$  is nonlinearly separable, then there exist a neural network function  $\phi$  such that  $\{\phi(x_j, \theta)\}_{i=1}^N$  is linearly separable.

J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.

Popular nonlinear models for  $\phi(\cdot, \theta)$ :

Popular nonlinear models for  $\phi(\cdot, \theta)$ :



Polynomials

Popular nonlinear models for  $\phi(\cdot, \theta)$ :



Polynomials

2 Kernel functions in SVM ( $k(x, y) = (\phi(x), \phi(y))$ )

#### Popular nonlinear models for $\phi(\cdot, \theta)$ :



- Polynomials
- 2 Kernel functions in SVM  $(k(x, y) = (\phi(x), \phi(y)))$ 
  - Neural Network functions

#### Popular nonlinear models for $\phi(\cdot, \theta)$ :



2

Kernel functions in SVM 
$$(k(x, y) = (\phi(x), \phi(y)))$$

Neural Network functions

Loss function using regularized logistic regression:

$$L = \sum_{j=1}^{N} \ell(y_j, p(\theta_L, \phi(x_j, \theta))) + \lambda \left( R(\theta_L) + R_N(\theta) \right)$$

# Example: MNIST data set

(Modified National Institute of Standards and Technology database)



#### Handwritten digits:

- Training set : 60,000
- Test set : 10,000
- Image size : 28\*28\*1
- Classes: 10

Logistic regression

$$P = LR(x) = softmax(Wx + b)$$

where,  $W \in \mathbb{R}^{10 \times 784}$ .



Logistic regression

$$P = LR(x) = softmax(Wx + b)$$

where,  $W \in \mathbb{R}^{10 \times 784}$ .

Training accuracy	Test accuracy
93.44	92.50

Table: Logistic regression on MNIST.

Logistic regression

$$P = LR(x) = softmax(Wx + b)$$

where,  $W \in \mathbb{R}^{10 \times 784}$ .

Training accuracy	Test accuracy
93.44	92.50

Table: Logistic regression on MNIST.

Logistic regression with DNN1

$$P = softmax(W_2\sigma(W_1x + b_1) + b_2)$$
(15)

where,  $W_1 \in \mathbb{R}^{1000 \times 784}$ ,  $W_2 \in \mathbb{R}^{10 \times 1000}$  and  $\sigma = \text{ReLU}$ .

Logistic regression

$$P = LR(x) = softmax(Wx + b)$$

where,  $W \in \mathbb{R}^{10 \times 784}$ .

Training accuracy	Test accuracy
93.44	92.50

Table: Logistic regression on MNIST.

Logistic regression with DNN<sub>1</sub>

$$P = softmax(W_2\sigma(W_1x + b_1) + b_2)$$
(15)

where,  $W_1 \in \mathbb{R}^{1000 \times 784}$ ,  $W_2 \in \mathbb{R}^{10 \times 1000}$  and  $\sigma = \text{ReLU}$ .

Training accuracy	Test accuracy
100.00	98.75

Table: LR with DNN1 on MNIST.

## Introduction

### 2 Data classification

Linear separable sets and logistic regressionNonlinear models

## 3 Convolutional neural networks (CNN)

A model for feature extraction

### 5 Multigrid method: a brief introduction

- A simple example
- Finite element method
- Multigrid method
- Image as a finite element function

### 6 MgNet

Multigrid as a CNN: MgNet

## 7 Concluding remarks

## Convolution: a special class of linear functions

$$(K * x)_{i,j} := \sum_{\ell=1}^{c} \sum_{s,t=-k}^{k} K_{s+k+1,t+k+1,\ell} x_{i+s,j+t,\ell},$$

• (i, j): pixel points,  $\ell$ : channel dimension





(16)



$$\begin{bmatrix}
-1 & -1 & -1 \\
2 & 2 & 2 \\
-1 & -1 & -1
\end{bmatrix}$$

$$\begin{bmatrix}
-1 & -1 & -1 \\
2 & 2 & 2 \\
-1 & -1 & -1
\end{bmatrix}$$



























# Sub-sampling (coarsening): stride and pooling









 $\mathcal{T}_1$ 

 $\mathcal{T}_2$ 

 $\mathcal{T}_3$ 

 $\mathcal{T}_4$ 


# Sub-sampling (coarsening): stride and pooling



Convolution with stride

$$(K *_{s} x)_{i,j} = \sum_{\ell=1}^{c} \sum_{p,q=-k}^{k} K_{p+k+1,q+k+1} x_{is+p,js+q}.$$

Pooling (max-pooling)

$$r(x)_{i,j} = \max\{x_{(i-1)s+p,(j-1)s+q} : p, q = 1 : s\}.$$

Jinchao Xu

**Classic CNN structure** 

**1** Initialization of inputs:  $g^{1,0}$ 

#### **Classic CNN structure**



1 Initialization of inputs:  $g^{1,0}$ 



#### **Classic CNN structure**



• For 
$$i = 1 : \nu_{\ell}$$

**Classic CNN structure** 



For 
$$i = 1 : \nu_{\ell}$$

$$g^{\ell,i} = \sigma\left(\theta^{\ell,i} * g^{\ell,i-1}\right) \tag{17}$$

#### **Classic CNN structure**



- 1 Initialization of inputs:  $g^{1,0}$
- 2 For ℓ = 1 : J For  $i = 1 : \nu_{\ell}$

$$g^{\ell,i} = \sigma\left(\theta^{\ell,i} * g^{\ell,i-1}\right)$$
(17)

- Restriction-Pooling
- 3 Output  $g^{J,\nu_J}$ .

#### **Neuron Activation Function**



$$g^{\ell+1,0} = R_{\ell}^{\ell+1} *_2 g^{\ell,\nu_{\ell}}$$
(18)

#### **Classic CNN structure**



- 1 Initialization of inputs:  $g^{1,0}$
- 2 For ℓ = 1 : J For  $i = 1 : \nu_{\ell}$

$$g^{\ell,i} = \sigma\left(\theta^{\ell,i} * g^{\ell,i-1}\right)$$
(17)

Restriction-Pooling

$$g^{\ell+1,0} = R_{\ell}^{\ell+1} *_2 g^{\ell,\nu_{\ell}}$$
(18)

3 Output  $g^{J,\nu_J}$ .

#### **Neuron Activation Function**





#### **Classic CNN structure**



- Initialization of inputs: g<sup>1,0</sup>
- **2** For  $\ell = 1 : J$ For  $i = 1 : \nu_{\ell}$

$$g^{\ell,i} = \sigma\left(\theta^{\ell,i} * g^{\ell,i-1}\right)$$
(17)

Restriction-Pooling

$$g^{\ell+1,0} = R_{\ell}^{\ell+1} *_2 g^{\ell,\nu_{\ell}}$$
(18)

3 Output  $g^{J,\nu_J}$ .

#### **Neuron Activation Function**







# Two classical CNN examples: LeNet-5 & AlexNet



Figure: First successful convolutional neural network: LeNet-5 (Y.LeCun, L.Bottou, Y.Bengio, and P.Haffner, 1998)

JII ON

# Two classical CNN examples: LeNet-5 & AlexNet



Figure: First successful convolutional neural network: LeNet-5 (Y.LeCun, L.Bottou, Y.Bengio, and P.Haffner, 1998)



Figure: The beginning of new revolution: AlexNet (A. Krizhevsky, I. Sutskever and G. Hinton, 2012)



Figure: From classical CNNs to ResNet and Pre-act ResNet (K. He, X. Zhang, S. Ren and J. Sun, 2015 and 2016)

ResNet

1 Initialization of inputs:  $r^{1,0}$ 

ResNet



Initialization of inputs:  $r^{1,0}$ 



ResNet



**2** For  $\ell = 1 : J$ 

For  $i = 1 : \nu_{\ell}$ 

Jinchao Xu

#### ResNet

2



Initialization of inputs: r<sup>1,0</sup>

For 
$$\ell = 1 : J$$
  
For  $i = 1 : \nu_{\ell}$ 

$$r^{\ell,i} = \sigma \left( r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * r^{\ell,i-1} \right).$$
<sup>(19)</sup>

#### **ResNet**



Initialization of inputs:  $r^{1,0}$ 

For 
$$\ell = 1 : J$$
  
For  $i = 1 : \nu_{\ell}$ 

$$\boldsymbol{r}^{\ell,i} = \sigma \left( \boldsymbol{r}^{\ell,i-1} + \boldsymbol{A}^{\ell,i} \ast \boldsymbol{\sigma} \circ \boldsymbol{B}^{\ell,i} \ast \boldsymbol{r}^{\ell,i-1} \right).$$
(19)

Restriction-Pooling

$$r^{\ell+1,0} = \sigma \left( R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 r^{\ell,\nu_{\ell}} \right).$$
(20)



#### Pre-act ResNet



Initialization of inputs: r<sup>1,0</sup>

#### ResNet



Initialization of inputs:  $r^{1,0}$ 

For 
$$\ell = 1 : J$$
  
For  $i = 1 : \nu_{\ell}$ 

$$r^{\ell,i} = \sigma \left( r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * r^{\ell,i-1} \right).$$
<sup>(19)</sup>

Restriction-Pooling

$$r^{\ell+1,0} = \sigma \left( R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 r^{\ell,\nu_{\ell}} \right).$$
(20)



#### Pre-act ResNet

Initialization of inputs: r<sup>1,0</sup>

2 For ℓ = 1 : J

#### ResNet



Initialization of inputs:  $r^{1,0}$ 

For 
$$\ell = 1 : J$$
  
For  $i = 1 : \nu_{\ell}$ 

$$r^{\ell,i} = \sigma \left( r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * r^{\ell,i-1} \right).$$
<sup>(19)</sup>

Restriction-Pooling

$$r^{\ell+1,0} = \sigma \left( R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 r^{\ell,\nu_{\ell}} \right).$$
(20)



#### Pre-act ResNet



Initialization of inputs:  $r^{1,0}$ 

For 
$$i = 1 : \nu_{\ell}$$

#### **ResNet**



Initialization of inputs:  $r^{1,0}$ 

For 
$$\ell = 1 : J$$
  
For  $i = 1 : \nu_{\ell}$ 

$$r^{\ell,i} = \sigma \left( r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * r^{\ell,i-1} \right).$$
(19)

Restriction-Pooling

$$r^{\ell+1,0} = \sigma \left( R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 r^{\ell,\nu_{\ell}} \right).$$
(20)



#### Pre-act ResNet



Initialization of inputs: r<sup>1,0</sup>

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}).$$
(21)

1 :  $\nu_{\ell}$ 

#### ResNet



Initialization of inputs:  $r^{1,0}$ 

For 
$$\ell = 1 : J$$
  
For  $i =$ 

$$r^{\ell,i} = \sigma \left( r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * r^{\ell,i-1} \right).$$
<sup>(19)</sup>

Restriction-Pooling

$$r^{\ell+1,0} = \sigma \left( R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 r^{\ell,\nu_{\ell}} \right).$$
(20)



#### Pre-act ResNet



Initialization of inputs:  $r^{1,0}$ 

🧿 For 
$$\ell = 1:$$
 J

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}).$$
(21)

Restriction-Pooling

$$r^{\ell+1,0} = R_{\ell}^{\ell+1} *_2 r^{\ell,\nu_{\ell}} + A^{\ell+1,0} \circ \sigma \circ B^{\ell+1,0} *_2 \sigma(r^{\ell,\nu_{\ell}}).$$
(22)

3 Output  $g^{J,\nu_J}$ .

# Applications of CNN

### CNN has been successfully used in:

- Computer vision
  - Classification, detection, segmentation...
  - Medical image processing,
  - Face recognition,
- Reinforcement learning
  - AlphaGo,
  - Automated driving,
- Natural language processing
  - Speech recognition,
  - Machine translation,



### Introduction

### 2 Data classification

Linear separable sets and logistic regression

Nonlinear models

### 3 Convolutional neural networks (CNN)

### A model for feature extraction

### 5 Multigrid method: a brief introduction

- A simple example
- Finite element method
- Multigrid method
- Image as a finite element function

### 6 MgNet

Multigrid as a CNN: MgNet

### 7 Concluding remarks

### Data space: g



### Data space: g



Feature extraction

### Data space: g



Feature extraction



### Data space: g





Feature extraction



### Data space: g



Feature extraction

### Feature space: u





Feature extraction,

# Feature space: u Data space: g Feature extraction $\gtrsim$ S Feature extraction,

A linear model: given an image g, find its feature u satisfying

A \* u = g

with a constraint

 $u \ge 0$ .

Jinchao Xu

(23)

A linear model: given an image g, find its feature u satisfying

A \* u = g

with a constraint

 $u \ge 0.$ 

Questions:



(23)

A linear model: given an image g, find its feature u satisfying

A \* u = g

with a constraint

 $u \ge 0.$ 

#### Questions:



What is A? (to be trained ...)

How to solve (23)?

(23)

A linear model: given an image g, find its feature u satisfying

A \* u = g

with a constraint

 $u \ge 0.$ 

#### Questions:



What is A? (to be trained ...)

How to solve (23)? (iterative methods)

(23)

### Iterative methods for Au = g: residual correction

Basic ideas:

 $u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$
$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:



**()** Form the residual:  $r = g - Au^{k-1}$ 

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:



Solve the residual eqn Ae = r approximately  $\hat{e} = Br$  with  $B \approx A^{-1}$ 

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

**1** Form the residual:  $r = g - Au^{k-1}$ 

2) Solve the residual eqn Ae = r approximately  $\hat{e} = Br$  with  $B \approx A^{-1}$ 

3 Update  $u^i = u^{i-1} + \hat{e}$ 

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

Form the residual: r = g - Au<sup>k-1</sup>
 Solve the residual eqn Ae = r approximately ê = Br with B ≈ A<sup>-1</sup>
 Update u<sup>i</sup> = u<sup>i-1</sup> + ê
 A basic iterative method:

$$u^{i} = u^{i-1} + B^{i}(g - Au^{i-1})$$
(25)

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

Form the residual: r = g - Au<sup>k-1</sup>
 Solve the residual eqn Ae = r approximately ê = Br with B ≈ A<sup>-1</sup>
 Update u<sup>i</sup> = u<sup>i-1</sup> + ê
 A basic iterative method:

$$u^{i} = u^{i-1} + B^{i}(g - Au^{i-1})$$
(25)

An example: A = D + L + U with D = DIAG(A)

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

Form the residual: r = g - Au<sup>k-1</sup>
 Solve the residual eqn Ae = r approximately ê = Br with B ≈ A<sup>-1</sup>
 Update u<sup>i</sup> = u<sup>i-1</sup> + ê
 A basic iterative method:

$$u^{i} = u^{i-1} + B^{i}(g - Au^{i-1})$$
(25)

```
An example: A = D + L + U with D = DIAG(A)

B = DIAG(A)^{-1}
```

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

● Form the residual: r = g - Au<sup>k-1</sup>
 ● Solve the residual eqn Ae = r approximately ê = Br with B ≈ A<sup>-1</sup>
 ● Update u<sup>i</sup> = u<sup>i-1</sup> + ê
 A basic iterative method:

$$J^{i} = u^{i-1} + B^{i}(g - Au^{i-1})$$
(25)

```
An example: A = D + L + U with D = DIAG(A)
```

•  $B = DIAG(A)^{-1}$  (Jacobi),

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

 Form the residual: r = g - Au<sup>k-1</sup>
 Solve the residual eqn Ae = r approximately ê = Br with B ≈ A<sup>-1</sup>
 Update u<sup>i</sup> = u<sup>i-1</sup> + ê
 A basic iterative method: u<sup>i</sup> = u<sup>i-1</sup> + B<sup>i</sup>(g - Au<sup>i-1</sup>)

$$u' = u'^{-1} + B'(g - Au'^{-1})$$
<sup>(25)</sup>

```
An example: A = D + L + U with D = DIAG(A)

1 B = DIAG(A)^{-1} (Jacobi),

2 B = TRIL(A)^{-1}
```

$$u^0, u^1, \ldots, u^{k-1} \longrightarrow u^k$$

Basic ideas:

• Form the residual:  $r = g - Au^{k-1}$ 

Solve the residual eqn Ae = r approximately  $\hat{e} = Br$  with  $B \approx A^{-1}$ 

3 Update 
$$u^i = u^{i-1} + \hat{e}$$

A basic iterative method:

$$u^{i} = u^{i-1} + B^{i}(g - Au^{i-1})$$
(25)

An example: A = D + L + U with D = DIAG(A)

1  $B = DIAG(A)^{-1}$  (Jacobi), 2  $B = TRIL(A)^{-1}$  (Gauss-Seidel)

# Iterative schemes for the constrained linear model

# Iterative schemes for the constrained linear model

Recall iterative methods without constraint

$$u^{i} = u^{i-1} + B^{i} * (g - A * u^{i-1})$$

## Iterative schemes for the constrained linear model

Recall iterative methods without constraint

$$u^{i} = u^{i-1} + B^{i} * (g - A * u^{i-1})$$

2 Image classification [ $\sigma = \text{ReLU}$ : dropping the negative values]

$$u^{i} = u^{i-1} + \sigma * B^{i} * \sigma (g - A * u^{i-1})$$

or, in terms of residual

$$r^{i} = (I - A * \sigma * B^{i} * \sigma)r^{i-1}.$$

# Numerical Examples for Constraint Linear Models

Model	CIFAR10	CIFAR100	# Parameters
ResNet18	93.45	74.45	11M
ResNet18-A <sup>ℓ</sup> -B <sup>ℓ,i</sup>	93.54	74.46	8.1M
pre-act ResNet18	93.75	74.33	11M
pre-act ResNet18-A <sup>ℓ</sup> -B <sup>ℓ,i</sup>	93.83	74.51	8.1M
ResNet34	94.71	77.20	21M
ResNet34-A <sup>ℓ</sup> -B <sup>ℓ,i</sup>	94.84	77.24	13M
pre-act ResNet34	94.76	77.25	21M
pre-act ResNet34- $A^{\ell}$ - $B^{\ell,i}$	94.84	77.40	13M

Table: The accuracy and number of parameters for ResNet, pre-act ResNet and their variants of modified versions on CIFAR10 and CIFAR100.



The simple constrained linear model seems adequate for image classifications

- The simple constrained linear model seems adequate for image classifications
- Peature extractions can be obtained using standard iterative methods

- The simple constrained linear model seems adequate for image classifications
- Peature extractions can be obtained using standard iterative methods
  - providing a validation of constrained linear model

- The simple constrained linear model seems adequate for image classifications
- 2 Feature extractions can be obtained using standard iterative methods
  - providing a validation of constrained linear model
- Otentially more efficient models can be developed.

# Introduction

- 2 Data classification
  - Linear separable sets and logistic regression
     Nonlinear models
  - Nonlinear models
- 3 Convolutional neural networks (CNN)
  - A model for feature extraction

#### 5 Multigrid method: a brief introduction

- A simple example
- Finite element method
- Multigrid method
- Image as a finite element function

#### 6 MgNet

Multigrid as a CNN: MgNet

### 7 Concluding remarks

Consider

$$Au = g \iff \arg\min \frac{1}{2}u^T Au - g^T u$$

Consider

$$Au = g \iff \arg\min \frac{1}{2}u^T Au - g^T u$$

• Gradient descent method:

$$u^{k+1} = u^k - \eta \nabla E(u^k)$$

Consider

$$Au = g \iff \arg\min \frac{1}{2}u^T Au - g^T u$$

• Gradient descent method:

$$u^{k+1} = u^k - \eta \nabla E(u^k) = u^k - \eta (Au^k - g)$$

Consider

$$Au = g \iff \arg\min \frac{1}{2}u^T Au - g^T u$$

• Gradient descent method:

$$u^{k+1} = u^k - \eta \nabla E(u^k) = u^k - \eta (Au^k - g)$$

Scaled gradient descent method:

$$u^{k+1} = u^k - \eta[\operatorname{diag}(A)]^{-1}(Au^k - g)$$

Consider

$$Au = g \iff \arg\min \frac{1}{2}u^T Au - g^T u$$

• Gradient descent method:

$$u^{k+1} = u^k - \eta \nabla E(u^k) = u^k - \eta (Au^k - g)$$

Scaled gradient descent method:

$$u^{k+1} = u^k - \eta[\operatorname{diag}(A)]^{-1}(Au^k - g)$$

• It converges for any symmetric, positive and definite A if  $\eta \ll 1$ 

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Note that  $\sigma(A_0) = \{3, 1, 0\}$ . Apply Gradient descent method with  $||Au^k - b|| \le 10^{-8}$ :

 $\epsilon$  # of iter = m

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Note that  $\sigma(A_0) = \{3, 1, 0\}$ . Apply Gradient descent method with  $||Au^k - b|| \le 10^{-8}$ :

 $\epsilon$  # of iter = m

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

$$\epsilon$$
 # of iter = m  
1 37

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

$$\begin{array}{c} \epsilon & \text{# of iter} = m \\ 1 & 37 \\ 10^{-1} \end{array}$$

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

$$\epsilon$$
 # of iter = m  
1 37  
 $10^{-1}$  236

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

ε	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

ε	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	1,918

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

ε	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	1,918
$10^{-3}$	

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

ε	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	1,918
$10^{-3}$	16, 115

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

$\epsilon$	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	1,918
10-3	16,115
$10^{-4}$	

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

$\epsilon$	# of iter $= m$
1	37
$10^{-1}$	236
$10^{-2}$	1,918
$10^{-3}$	16, 115
10 <sup>-4</sup>	130, 168
#### Gradient descent for a nearly singular system

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Note that  $\sigma(A_0) = \{3, 1, 0\}$ . Apply Gradient descent method with  $||Au^k - b|| \le 10^{-8}$ :

ε	# of iter $= m$
1	37
10 <sup>-1</sup>	236
10 <sup>-2</sup>	1,918
10 <sup>-3</sup>	16,115
10 <sup>-4</sup>	130, 168
0 [singular case]	

#### Gradient descent for a nearly singular system

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Note that  $\sigma(A_0) = \{3, 1, 0\}$ . Apply Gradient descent method with  $||Au^k - b|| \le 10^{-8}$ :

ε	# of iter $= m$
1	37
10 <sup>-1</sup>	236
10-2	1,918
10 <sup>-3</sup>	16,115
10 <sup>-4</sup>	130, 168
0 [singular case]	20

#### Gradient descent for a nearly singular system

Consider:  $(A_0 + \epsilon I)u = g$ 

$$A_0 = \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}, \quad g = \begin{pmatrix} -1 \\ -1 \\ 2 \end{pmatrix} \in R(A_0), \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in N(A_0).$$

Note that  $\sigma(A_0) = \{3, 1, 0\}$ . Apply Gradient descent method with  $||Au^k - b|| \le 10^{-8}$ :

ε	# of iter $= m$
1	37
10 <sup>-1</sup>	236
10-2	1,918
10 <sup>-3</sup>	16,115
10 <sup>-4</sup>	130, 168
0 [singular case]	20

Iterative method usually is OK for singular system, but subtle for nearly singular system! Ref for semi-definite case: Keller 1965

### Remedy of GD: Multilevel methods

$$V_1 = \mathbb{R}^3$$
,  $V_2 = \operatorname{span}\{p\} \subset V_1$ 

Initialization of inputs

 $A_1 = A_{\epsilon}, \quad g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

# Remedy of GD: Multilevel methods

$$V_1 = \mathbb{R}^3, \quad V_2 = \operatorname{span}\{p\} \subset V_1$$

Initialization of inputs

$$A_1 = A_{\epsilon}, \quad g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}$$

Iterate:

One step of GD method on V<sub>1</sub>

$$u_1 \leftarrow u_1 + \eta D_{\epsilon}^{-1}(g_1 - A_1 u_1).$$

2 Consider  $A_1 e_1 = r_1 \equiv g_1 - A_1 u_1$  and "pool" it to  $V_2$  and solve it:

 $A_2u_2=g_2,$ 

# Remedy of GD: Multilevel methods

$$V_1 = \mathbb{R}^3, \quad V_2 = \operatorname{span}\{p\} \subset V_1$$

Initialization of inputs

$$A_1 = A_{\epsilon}, \quad g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}$$

Iterate:

One step of GD method on V<sub>1</sub>

$$u_1 \leftarrow u_1 + \eta D_{\epsilon}^{-1}(g_1 - A_1 u_1).$$

2 Consider  $A_1e_1 = r_1 \equiv g_1 - A_1u_1$  and "pool" it to  $V_2$  and solve it:

$$A_2u_2 = g_2, \quad u_2 = A_2^{-1}g_2$$

with

$$A_2 = p^T A_1 p, g_2 = p^T r_1$$

3 update  $u_1 \leftarrow u_1 + pu_2$ .

# GD with $\eta = 0.7$		
ε	original	normalized expanded
1.	37	
10 <sup>-1</sup>	236	$\mathbf{D} \cdot \mathbf{O}$
10 <sup>-2</sup>	1,918	
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

# GD with $\eta = 0.7$		
ε	original	normalized expanded
1.	37	
10 <sup>-1</sup>	236	$\mathbf{D} \cdot \mathbf{O}$
10 <sup>-2</sup>	1,918	
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

# GD with $\eta = 0.7$		
ε	original	normalized expanded
1.	37	
10 <sup>-1</sup>	236	$\mathbf{D} \cdot \mathbf{O}$
10 <sup>-2</sup>	1,918	
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

# GD with $\eta = 0.7$		
ε	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	$\mathbf{D}$
10 <sup>-2</sup>	1,918	
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

$\#$ GD with $\eta = 0.7$		
$\epsilon$	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	14
10 <sup>-2</sup>	1,918	
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

# GD with $\eta = 0.7$		
$\epsilon$	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	14
10 <sup>-2</sup>	1,918	14
10 <sup>-3</sup>	16, 115	
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

# GD with $\eta = 0.7$		
$\epsilon$	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	14
10 <sup>-2</sup>	1,918	14
10 <sup>-3</sup>	16, 115	16
10 <sup>-4</sup>	130, 168	X
10 <sup>-5</sup>	> 1,000,000	
10 <sup>-9</sup>	> 1,000,000	
10 <sup>-10</sup>	21	
0.	20	

$\#$ GD with $\eta = 0.7$		
$\epsilon$	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	14
10 <sup>-2</sup>	1,918	14
10 <sup>-3</sup>	16, 115	16
10 <sup>-4</sup>	130, 168	16
10 <sup>-5</sup>	> 1,000,000	16
10 <sup>-9</sup>	> 1,000,000	15
10 <sup>-10</sup>	21	15
0.	20	

$\#$ GD with $\eta = 0.7$		
$\epsilon$	original	normalized expanded
1.	37	13
10 <sup>-1</sup>	236	14
10 <sup>-2</sup>	1,918	14
10 <sup>-3</sup>	16, 115	16
10 <sup>-4</sup>	130, 168	16
10 <sup>-5</sup>	> 1,000,000	16
10 <sup>-9</sup>	> 1,000,000	15
10 <sup>-10</sup>	21	15
0.	20	

• 1D Poisson equation on  $\Omega = (0, 1)$ 

-u'' = f, u(0) = u(1) = 0.

• 1D Poisson equation on  $\Omega = (0, 1)$ 

$$-u'' = f$$
,  $u(0) = u(1) = 0$ .

Variational principle:

$$-u''-f=0 \iff \int_0^1 (-u''-f)v \, dx=0$$

for all "good" test functions v.

• 1D Poisson equation on  $\Omega = (0, 1)$ 

-u'' = f, u(0) = u(1) = 0.

Variational principle:

$$-u''-f=0 \iff \int_0^1 (-u''-f)v \, dx=0$$

for all "good" test functions v.

U Function space:  $V = \{v \text{ is piecewise smooth and continuous }, v(0) = v(1) = 0\}.$ 

• 1D Poisson equation on  $\Omega = (0, 1)$ 

$$-u'' = f$$
,  $u(0) = u(1) = 0$ .

Variational principle:

$$-u''-f=0 \iff \int_0^1 (-u''-f)v \, dx=0$$

for all "good" test functions v.



$$\int_0^1 -u'' v \, dx = \int_0^1 u' v' \, dx - u' v |_0^1 = \int_0^1 u' v' \, dx$$

• 1D Poisson equation on  $\Omega = (0, 1)$ 

$$-u'' = f$$
,  $u(0) = u(1) = 0$ .

Variational principle:

$$-u'' - f = 0 \iff \int_0^1 (-u'' - f) v \, dx = 0$$

for all "good" test functions v.



$$\int_0^1 -u'' v \, dx = \int_0^1 u' v' \, dx - u' v |_0^1 = \int_0^1 u' v' \, dx.$$

Variational formulation: Find  $u \in V$ , such that

$$\int_0^1 u'v' \, dx = \int_0^1 fv \, dx \quad \forall v \in V.$$

# Galerkin method

### Galerkin method

• Galerkin method: Find  $u_h \in V_h$  such that

$$\int_0^1 u_h' v_h' \, dx = \int_0^1 f v_h \, dx \quad \forall v_h \in V_h.$$

#### Galerkin method

• Galerkin method: Find  $u_h \in V_h$  such that



### 1D linear system on uniform grid

#### 1D linear system on uniform grid

• Stiffness matrix 
$$A_{ij} = \int_0^1 \varphi'_j \varphi'_i \, dx, \, g_i = \int_0^1 f \varphi_i \, dx$$

$$A = \frac{1}{h} \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \quad g = h \begin{pmatrix} g(x_1) \\ g(x_2) \\ \vdots \\ g(x_{N-1}) \\ g(x_N) \end{pmatrix} + \mathcal{O}(h^3).$$

or

$$\frac{1}{h}(-u_{i-1}+2u_i-u_{i+1})=g_i, \quad 1 \le i \le n$$

Namely

A \* u = g

where

$$A = (-1, 2, -1)/h$$







$$A = 8 \times \text{tridiag}(-1, 2, -1) = \left( (\phi'_i, \phi'_j) \right) \in \mathbb{R}^{7 \times 7}$$
  
Note  
$$\phi_{3,1} = \frac{1}{2} \phi_{2,1} + \phi_{2,2} + \frac{1}{2} \phi_{2,3}$$



$$A = 8 imes ext{tridiag}(-1, 2, -1) = \left((\phi'_i, \phi'_j)\right) \in \mathbb{R}^{7 imes 7}$$

 $\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$ 

This can be written as

Note

 $\Phi^3 = \Phi^2 R$ 



$$A = 8 imes ext{tridiag}(-1, 2, -1) = \left( (\phi'_i, \phi'_j) 
ight) \in \mathbb{R}^{7 imes 7}$$

Note

$$\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$$

This can be written as

 $\Phi^3 = \Phi^2 R = R *_2 \Phi^2$ 



$$A = 8 imes ext{tridiag}(-1, 2, -1) = \left( (\phi'_i, \phi'_j) \right) \in \mathbb{R}^{7 imes 7}$$

 $\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$ 

This can be written as

 $\Phi^3 = \Phi^2 R = R *_2 \Phi^2$ 

Similarly

Note

 $\Phi^2 = \Phi^1 \mathbf{P}$ 



$$A = 8 imes ext{tridiag}(-1, 2, -1) = \left( (\phi'_i, \phi'_j) 
ight) \in \mathbb{R}^{7 imes 7}$$

 $\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$ 

This can be written as

 $\Phi^3 = \Phi^2 R = R *_2 \Phi^2$ 

Similarly

Note

 $\Phi^2 = \Phi^1 \mathbf{P} = \mathbf{R} *_2 \Phi^1$ 



$$A = 8 imes ext{tridiag}(-1, 2, -1) = \left( (\phi'_i, \phi'_j) 
ight) \in \mathbb{R}^{7 imes 7}$$

 $\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$ 

This can be written as

 $\Phi^3 = \Phi^2 R = R *_2 \Phi^2$ 

Similarly

Note

 $\Phi^2 = \Phi^1 \mathbf{P} = \mathbf{R} *_2 \Phi^1$ 

where

$$R = \begin{pmatrix} \frac{1}{2} \\ 1 \\ \frac{1}{2} \end{pmatrix}, P = \frac{1}{2} \begin{pmatrix} 1 & 0 & 0 \\ 2 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 2 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & 2 \end{pmatrix}$$

#### Two grid method for 1D case

 $V_1 = \operatorname{span}\{\phi_1, \phi_2, \phi_3, \phi_4, \phi_5, \phi_6, \phi_7\}, \quad V_2 = \operatorname{span}\{\phi_{2,1}, \phi_{2,2}, \phi_{2,3}\} \subset V_1$ 

Initialization of inputs

 $A_1 = A$ ,  $g_1 \leftarrow g$ ,  $u_1 \leftarrow$  random.
### Two grid method for 1D case

$$V_1 = \operatorname{span}\{\phi_1, \phi_2, \phi_3, \phi_4, \phi_5, \phi_6, \phi_7\}, \quad V_2 = \operatorname{span}\{\phi_{2,1}, \phi_{2,2}, \phi_{2,3}\} \subset V_1$$

Initialization of inputs

$$A_1 = A$$
,  $g_1 \leftarrow g$ ,  $u_1 \leftarrow$  random

Iterate:

A few steps of GD method on V<sub>1</sub>

$$u_1 \leftarrow u_1 + \eta(g_1 - A_1 u_1).$$

2 Consider  $A_1 e_1 = r_1 \equiv g_1 - A_1 u_1$  and "pool" it to  $V_2$  and solve it:

 $A_2u_2=g_2,$ 

### Two grid method for 1D case

$$V_1 = \operatorname{span}\{\phi_1, \phi_2, \phi_3, \phi_4, \phi_5, \phi_6, \phi_7\}, \quad V_2 = \operatorname{span}\{\phi_{2,1}, \phi_{2,2}, \phi_{2,3}\} \subset V_1$$

Initialization of inputs

$$A_1 = A$$
,  $g_1 \leftarrow g$ ,  $u_1 \leftarrow$  random

Iterate:

A few steps of GD method on V<sub>1</sub>

$$u_1 \leftarrow u_1 + \eta(g_1 - A_1 u_1).$$

2 Consider  $A_1e_1 = r_1 \equiv g_1 - A_1u_1$  and "pool" it to  $V_2$  and solve it:

$$A_2u_2 = g_2, \quad u_2 = A_2^{-1}g_2$$

with

$$A_2 = P^T A_1 P, \quad g_2 = P^T r_1$$

**3** update  $u_1 \leftarrow u_1 + Pu_2$ .

Images:

Images:

Piecewise (bi-)linear functions

Images:

- Piecewise (bi-)linear functions
- Initial grid T with size:

 $m = 2^{s} + 1, n = 2^{t} + 1$ 

Images:

- Piecewise (bi-)linear functions
- Initial grid T with size:

 $m = 2^{s} + 1, n = 2^{t} + 1$ 

#### Images:

- Piecewise (bi-)linear functions
- Initial grid T with size:

$$m = 2^s + 1, n = 2^t + 1$$

#### piecewise linear functions on multilevel grids



Model Problem:

$$\begin{aligned} -\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega, \quad \Omega = (0, 1)^d. \end{aligned}$$



Model Problem:

$$\begin{aligned} -\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial \Omega, \quad \Omega = (0, 1)^d. \end{aligned}$$



Discrete case:

Model Problem:

$$\begin{aligned} -\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial \Omega, \quad \Omega = (0, 1)^d. \end{aligned}$$



Discrete case:

Model Problem:

$$\begin{aligned} -\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial \Omega, \quad \Omega = (0, 1)^d. \end{aligned}$$



Discrete case:

▶ *d* = 2:

$$4u_{i,j} - u_{i+1,j} - u_{i-1,j} - u_{i,j+1} - u_{i,j-1} = g_{i,j},$$
(26)

Model Problem:

$$\begin{aligned} -\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial \Omega, \quad \Omega = (0, 1)^d. \end{aligned}$$



#### Discrete case:

▶ *d* = 2:

$$4u_{i,j} - u_{i+1,j} - u_{i-1,j} - u_{i,j+1} - u_{i,j-1} = g_{i,j},$$
(26)

with

$$A * u = f$$
, for  $A = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{pmatrix}$ , (27)

Smoother:

• Gradient descent, say  $\eta = 1/8$ 

$$u^{k+1} = u^k - \eta (A * u^k - g),$$

Smoother:

• Gradient descent, say  $\eta = 1/8$ 

$$u^{k+1} = u^k - \eta (A * u^k - g),$$

• Apply GD method twice, we get (from  $1 \times 1$  to  $3 \times 3$  kernels)

$$u^{k+1} = u^k - S * (A * u^k - g),$$

where

$$S = \frac{1}{64} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 12 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$
 (28)

Smoother:

• Gradient descent, say  $\eta = 1/8$ 

$$u^{k+1} = u^k - \eta (A * u^k - g),$$

• Apply GD method twice, we get (from  $1 \times 1$  to  $3 \times 3$  kernels)

$$u^{k+1} = u^k - S * (A * u^k - g),$$

where

$$S = \frac{1}{64} \begin{pmatrix} 0 & 1 & 0 \\ 1 & 12 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$
 (28)

Restriction (pooling)

$$\mathcal{R}_{\ell}^{\ell+1} : \mathbb{R}^{m_{\ell} \times n_{\ell}} \mapsto \mathbb{R}^{m_{\ell+1} \times n_{\ell+1}} \text{ is } \mathcal{R}_{\ell}^{\ell+1} = \mathcal{R}_{*2}$$

$$\mathcal{R} = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 1 & 1 \end{pmatrix} \quad \Phi^{\ell+1} = \mathcal{R} \oplus \Phi^{\ell}$$

where

 $\begin{pmatrix} \frac{2}{1} & \frac{1}{2} & \frac{2}{1} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$ 

(29)



Initialization of inputs



Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

2 Smoothing and restriction

Initialization of inputs



Initialization of inputs



Initialization of inputs



Initialization of inputs

**★** For 
$$i = 1 : \nu_{\ell}$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

2 Smoothing and restriction
 ▶ For ℓ = 1 : J

**★** For 
$$i = 1 : \nu_{\ell}$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

★ Form restricted residual and set initial guess:

$$u_{\ell+1,0} \leftarrow 0, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}), \quad A_{\ell+1} = R_{\ell} *_2 A_{\ell} * R_{\ell} *_2^{\top}.$$

Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

2 Smoothing and restriction
 ▶ For ℓ = 1 : J

**★** For 
$$i = 1 : \nu_{\ell}$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

★ Form restricted residual and set initial guess:

$$u_{\ell+1,0} \leftarrow 0, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}), \quad A_{\ell+1} = R_{\ell} *_2 A_{\ell} * R_{\ell} *_2^{\top}.$$

Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

# 2 Smoothing and restriction ▶ For ℓ = 1 : J

**★** For 
$$i = 1 : \nu_\ell$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

★ Form restricted residual and set initial guess:

$$u_{\ell+1,0} \leftarrow 0, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}), \quad A_{\ell+1} = R_{\ell} *_2 A_{\ell} * R_{\ell} *_2^{\top}.$$



Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

#### 2 Smoothing and restriction For $\ell = 1 : J$

**★** For 
$$i = 1 : \nu_\ell$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

★ Form restricted residual and set initial guess:

$$u_{\ell+1,0} \leftarrow 0, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}), \quad A_{\ell+1} = R_{\ell} *_2 A_{\ell} * R_{\ell} *_2^{\top}.$$



3 Update with prolongation: For  $\ell = J - 1$ : 1

$$u_{\ell} \leftarrow u_{\ell} + R_{\ell} *_2^{\top} (u_{\ell+1} - u_{\ell+1,0}).$$

Initialization of inputs

 $g_1 \leftarrow g, \quad u_1 \leftarrow \text{random}.$ 

#### 2 Smoothing and restriction For $\ell = 1 : J$

**★** For 
$$i = 1 : \nu_\ell$$

$$u_{\ell} \leftarrow u_{\ell} + S_{\ell} * (g_{\ell} - A_{\ell} * u_{\ell}). \tag{30}$$

★ Form restricted residual and set initial guess:

$$u_{\ell+1,0} \leftarrow 0, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}), \quad A_{\ell+1} = R_{\ell} *_2 A_{\ell} * R_{\ell} *_2^{\top}.$$



3 Update with prolongation: For  $\ell = J - 1$ : 1

$$u_{\ell} \leftarrow u_{\ell} + R_{\ell} *_2^{\top} (u_{\ell+1} - u_{\ell+1,0}).$$

	n = 10 <sup>6</sup>	$n = 10^{7}$	$n = 10^8$	$n = 10^9$	
G. E. on Fugaku					General dense matrix

	$n = 10^{6}$	$n = 10^{7}$	$n = 10^{8}$	$n = 10^9$	
G. E. on Fugaku	0.93s				General dense matrix

	n = 10 <sup>6</sup>	$n = 10^{7}$	$n = 10^8$	$n = 10^9$	
G. E. on Fugaku	0.93s	0.27 hr			General dense matrix

	n = 10 <sup>6</sup>	$n = 10^{7}$	$n = 10^{8}$	$n = 10^9$	
G. E. on Fugaku	0.93s	0.27 hr	11 days		General dense matrix

G.E. on the supercomputer: Fugaku(~400 petaFLOPs)

	<i>n</i> = 10 <sup>6</sup>	$n = 10^{7}$	$n = 10^{8}$	$n = 10^9$	
G. E. on Fugaku	0.93s	0.27 hr	11 days	30.1 Y	General d

eneral dense matrix

G. E. on Fugaku	$n = 10^6$ 0.93s	$n = 10^7$ 0.27 hr	$n = 10^8$ 11 days	$n = 10^9$ 30.1 Y	General dense matrix
Multigrid on 1 core	n = 10 <sup>6</sup>	$n = 10^7$	<i>n</i> = 10 <sup>8</sup>	<i>n</i> = 10 <sup>9</sup>	3D Poisson equation
G. E. on Fugaku	$n = 10^6$ 0.93s	$n = 10^7$ 0.27 hr	$n = 10^8$ 11 days	$n = 10^9$ 30.1 Y	General dense matrix
---------------------	--------------------------------------	----------------------------	----------------------------	----------------------------	----------------------
Multigrid on 1 core	<i>n</i> = 10 <sup>6</sup> 0.25 s	<i>n</i> = 10 <sup>7</sup>	<i>n</i> = 10 <sup>8</sup>	<i>n</i> = 10 <sup>9</sup>	3D Poisson equation

G. E. on Fugaku	n = 10 <sup>6</sup> 0.93s	$n = 10^7$ 0.27 hr	$n = 10^8$ 11 days	$n = 10^9$ 30.1 Y	General dense matrix
Multigrid on 1 core	<i>n</i> = 10 <sup>6</sup> 0.25 s	$n = 10^7$ 2.5 s	<i>n</i> = 10 <sup>8</sup>	<i>n</i> = 10 <sup>9</sup>	3D Poisson equation

G. E. on Fugaku	n = 10 <sup>6</sup> 0.93s	$n = 10^7$ 0.27 hr	$n = 10^8$ 11 days	$n = 10^9$ 30.1 Y	General dense matrix
Multigrid on 1 core	$n = 10^6$ 0.25 s	$n = 10^7$ 2.5 s	<i>n</i> = 10 <sup>8</sup> 25 s	<i>n</i> = 10 <sup>9</sup>	3D Poisson equation

G. E. on Fugaku	$n = 10^6$ 0.93s	$n = 10^7$ 0.27 hr	$n = 10^8$ 11 days	$n = 10^9$ 30.1 Y	General dense matrix
	n = 10 <sup>6</sup>	$n = 10^{7}$	<i>n</i> = 10 <sup>8</sup>	<i>n</i> = 10 <sup>9</sup>	
Multigrid on 1 core	0.25 s	2.5 s	25 s	250 s	3D Poisson equation

• Multigrid is powerful;

- Multigrid is powerful;
- Can the power of multigrid be transformed to CNN?

- Multigrid is powerful;
- Can the power of multigrid be transformed to CNN?
  better structure CNN with fewer weights?

- Multigrid is powerful;
- Can the power of multigrid be transformed to CNN?
  - better structure CNN with fewer weights?
  - faster training algorithms? [main original goal, making progress]

- - Linear separable sets and logistic regression Nonlinear models

- - A simple example
  - Finite element method
  - Multigrid method

#### MgNet

Multigrid as a CNN: MgNet

Partial differential equations

 $-\Delta u = g.$ 

(31)

Partial differential equations

$$-\Delta u = g. \tag{31}$$

Constrained linear model: Given an image g, find its feature u such that

$$A * u = g. \tag{32}$$

Partial differential equations

$$-\Delta u = g. \tag{31}$$

Constrained linear model: Given an image g, find its feature u such that

$$A * u = g. \tag{32}$$

#### Main idea:

 Use a geometric multigrid method for PDE (31) to solve the data-feature equation (32)!



Initialization of inputs:  $g_1 \leftarrow \theta * g$ ,  $u_1 \leftarrow$  random





Initialization of inputs:  $g_1 \leftarrow \theta * g$ ,  $u_1 \leftarrow$  random

• For 
$$i = 1 : \nu_{\ell}$$

$$u_{\ell} \leftarrow u_{\ell} + \sigma \circ S_{\ell} * \sigma(g_{\ell} - A_{\ell} * u_{\ell}).$$
(33)

1 Initialization of inputs:  $g_1 \leftarrow \theta * g$ ,  $u_1 \leftarrow random$ 

2 For 
$$\ell = 1 : J$$

For 
$$i = 1 : \nu_\ell$$

$$u_{\ell} \leftarrow u_{\ell} + \sigma \circ S_{\ell} * \sigma(g_{\ell} - A_{\ell} * u_{\ell}).$$
(33)

Restriction-Pooling

$$u_{\ell+1,0} \leftarrow \Pi_{\ell}^{\ell+1} u_{\ell}, \quad g_{\ell+1} \leftarrow R_{\ell} *_2 (g_{\ell} - A_{\ell} * u_{\ell}) + A_{\ell+1} * u_{\ell+1}^0,$$



 $\phi(g)=u_J.$ 

Multigrid:



Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori CNN:



Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

CNN:

Almost identically same structure as multigrid!

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

- Almost identically same structure as multigrid!
- $A_{\ell}^{i}, S_{\ell}^{i}, R_{\ell}^{i}$  are all trained!

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

- Almost identically same structure as multigrid!
- $A^i_{\ell}, S^i_{\ell}, R^i_{\ell}$  are all trained!
- Activation, ReLU, is introduced (to drop-off negative pixel values).

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

- Almost identically same structure as multigrid!
- $A^i_{\ell}, S^i_{\ell}, R^i_{\ell}$  are all trained!
- Activation, ReLU, is introduced (to drop-off negative pixel values).
- Extra channels are introduced.

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

- Almost identically same structure as multigrid!
- $A^i_{\ell}, S^i_{\ell}, R^i_{\ell}$  are all trained!
- Activation, ReLU, is introduced (to drop-off negative pixel values).
- Extra channels are introduced.

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

CNN:

- Almost identically same structure as multigrid!
- $A_{\ell}^{i}, S_{\ell}^{i}, R_{\ell}^{i}$  are all trained!
- Activation, ReLU, is introduced (to drop-off negative pixel values).
- Extra channels are introduced.

CNN versus multigrid: classic approaches versus MgNet

Multigrid:

•  $A_{\ell}, S_{\ell}, R_{\ell}$  are all given a priori

#### CNN:

- Almost identically same structure as multigrid!
- $A^i_{\ell}, S^i_{\ell}, R^i_{\ell}$  are all trained!
- Activation, ReLU, is introduced (to drop-off negative pixel values).
- Extra channels are introduced.

#### CNN versus multigrid: classic approaches versus MgNet

CNN:

Classic: Almost all the components are unrelated and need to be trained MgNet: Most of the components are related and can be given a priori

#### 2 Multigrid

Classic: Almost all the components are a priori given

MgNet: Some of components can be trained!

### Approximation properties of MgNet and deep CNNs

### Approximation properties of MgNet and deep CNNs

#### Theorem (He, Li and Xu 2021)

Assume  $f : \Omega \subset \mathbb{R}^{d \times d} \mapsto \mathbb{R}$  and  $\Omega$  is bounded, then there exist one version of MgNet  $\tilde{f} : \mathbb{R}^{d \times d} \mapsto \mathbb{R}$  with multichannel and  $3 \times 3$  kernels in smoothing, where

depth (number of smoothing):  $\leq d$ ,

width (number of channels for *i*-th smoothing):  $\leq 2(2i-1)^2$ , i = 1 : d - 1,

and number of channel for the last layer is N, such that

$$\left|f - \tilde{f}\right|_{L^{2}(\Omega)} \lesssim N^{-\frac{1}{2} - \frac{3}{2d^{2}}}.$$
(35)

(34)

### Approximation properties of MgNet and deep CNNs

#### Theorem (He, Li and Xu 2021)

Assume  $f : \Omega \subset \mathbb{R}^{d \times d} \mapsto \mathbb{R}$  and  $\Omega$  is bounded, then there exist one version of MgNet  $\tilde{f} : \mathbb{R}^{d \times d} \mapsto \mathbb{R}$  with multichannel and  $3 \times 3$  kernels in smoothing, where

depth (number of smoothing):  $\leq d$ ,

width (number of channels for i-th smoothing):  $\leq 2(2i-1)^2$ , i = 1 : d - 1,

and number of channel for the last layer is N, such that

$$\|f - \tilde{f}\|_{L^2(\Omega)} \lesssim N^{-\frac{1}{2} - \frac{3}{2d^2}}.$$
 (35)

Main properties:

- For 2D images with channels.
- Similar results hold for classical, ResNet and pre-act ResNet type networks.
- Convolution (the most commonly used type): 3 by 3 kernel with multichannel for both periodic and zero padding.
- Achieve the same asymptotic approximation rate to DNN with one-hidden layer.

(34)

#### Theorem (MgNet and pre-act ResNet, He and Xu 2019)

The MgNet model recovers the pre-act ResNet (K. He et al 2016) as follows

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : \nu_{\ell},$$
(36)

where

$$r^{\ell,i}=g^{\ell}-A^{\ell}*u^{\ell,i}.$$

#### Theorem (MgNet and pre-act ResNet, He and Xu 2019)

The MgNet model recovers the pre-act ResNet (K. He et al 2016) as follows

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : \nu_{\ell},$$
(36)

where

$$r^{\ell,i}=g^\ell-A^\ell*u^{\ell,i}.$$

Proof.

$$u^{\ell,i} = u^{\ell,i-1} + \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$

		-	M
 Inte	na	0	хп

#### Theorem (MgNet and pre-act ResNet, He and Xu 2019)

The MgNet model recovers the pre-act ResNet (K. He et al 2016) as follows

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : \nu_{\ell},$$
(36)

where

$$r^{\ell,i}=g^\ell-A^\ell*u^{\ell,i}.$$

Proof.

$$u^{\ell,i} = u^{\ell,i-1} + \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  
$$\Rightarrow A^{\ell} * u^{\ell,i} = A^{\ell} * u^{\ell,i-1} + A^{\ell} * \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
Pre-act ResNet: a "residual" version of MgNet

### Theorem (MgNet and pre-act ResNet, He and Xu 2019)

The MgNet model recovers the pre-act ResNet (K. He et al 2016) as follows

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : \nu_{\ell},$$
(36)

where

$$r^{\ell,i}=g^\ell-A^\ell*u^{\ell,i}.$$

Proof.

$$u^{\ell,i} = u^{\ell,i-1} + \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  

$$\Rightarrow A^{\ell} * u^{\ell,i} = A^{\ell} * u^{\ell,i-1} + A^{\ell} * \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  

$$\Rightarrow g^{\ell} - A^{\ell} * u^{\ell,i} = g^{\ell} - A^{\ell} * u^{\ell,i} - A^{\ell} * \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$

Pre-act ResNet: a "residual" version of MgNet

### Theorem (MgNet and pre-act ResNet, He and Xu 2019)

The MgNet model recovers the pre-act ResNet (K. He et al 2016) as follows

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : \nu_{\ell},$$
(36)

where

$$r^{\ell,i}=g^\ell-A^\ell*u^{\ell,i}.$$

Proof.

$$u^{\ell,i} = u^{\ell,i-1} + \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  

$$\Rightarrow A^{\ell} * u^{\ell,i} = A^{\ell} * u^{\ell,i-1} + A^{\ell} * \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  

$$\Rightarrow g^{\ell} - A^{\ell} * u^{\ell,i} = g^{\ell} - A^{\ell} * u^{\ell,i} - A^{\ell} * \sigma \circ B^{\ell,i} * \sigma(g^{\ell} - A^{\ell} * u^{\ell,i-1}),$$
  

$$\Rightarrow r_{\ell}^{i} = r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}).$$
(37)

### MgNet versus other CNNs

Model	Accuracy	# Parameters
ResNet18	95.28	11.2M
pre-act ResNet18	95.08	10.2M
MgNet[2,2,2,2],256	96.00	8.2M

Table: The comprison of MgNet and classical CNN on Cifar10

Model	Accuracy	# Parameters
ResNet18	77.54	11.2M
pre-act ResNet18	77.29	11.2M
MgNet[2,2,2,2],256	79.94	8.3M
MgNet[2,2,2,2],512	81.35	33.1M
MgNet[2,2,2,2],1024	82.46	132.2M

Table: The comprison of MgNet and classical CNN on Cifar100

Model	Accuracy	Parameters
ResNet18	72.12	11.2M
MgNet[2,2,2,2], [64,128,256,512]	73.36	13.0M
MgNet[3,4,6,3],[128,256,512,1024]	78.59	71.3M

Table: The comprison of MgNet and classical CNN on ImageNet.

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

### **Properties:**

A uniform framework for understanding and designing CNNs.

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- A uniform framework for understanding and designing CNNs.
- 2 Unify and develop many classical CNN structures.

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- A uniform framework for understanding and designing CNNs.
- 2 Unify and develop many classical CNN structures.
  - AlexNet, ResNet, U-Net, Audo-encoder, DenseNet ...

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- A uniform framework for understanding and designing CNNs.
- 2 Unify and develop many classical CNN structures.
  - AlexNet, ResNet, U-Net, Audo-encoder, DenseNet ...
- Construct the new

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

- A uniform framework for understanding and designing CNNs.
- 2 Unify and develop many classical CNN structures.
  - AlexNet, ResNet, U-Net, Audo-encoder, DenseNet ...
- Construct the new

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State), 2021.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354, https://doi.org/10.1007/s11425-019-9547-2
- J. He, Y. Chen, L. Zhang, J. Xu. Constrained Linear Data-feature Mapping for Image Classification. ArXiv:1911.10428, 2019.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. ArXiv:2010.14088, 2020
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. ArXiv:2109.00190, 2021.

### **Properties:**

- A uniform framework for understanding and designing CNNs.
- 2 Unify and develop many classical CNN structures.
  - AlexNet, ResNet, U-Net, Audo-encoder, DenseNet ....

### Construct the new

- reduce the free parameters 1%, .1%, .01%...?
- increase the generalization accuracy,
- accelerate the training speed,
- extend to general graph models,
- ► ...

- - Linear separable sets and logistic regression Nonlinear models

- - A simple example
  - Finite element method
  - Multigrid method
- - Multigrid as a CNN: MgNet

General data classification

General data classification

Linearly separable sets: logistic regressions

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods
- MgNet
  - A special CNN by making minor modification of multigrid
  - More mathematical structures lead to fewer weights

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods
- MgNet
  - A special CNN by making minor modification of multigrid
  - More mathematical structures lead to fewer weights
- Goals

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods
- MgNet
  - A special CNN by making minor modification of multigrid
  - More mathematical structures lead to fewer weights
- Goals
  - Slimmer CNN with higher generalization accuracy

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods
- MgNet
  - A special CNN by making minor modification of multigrid
  - More mathematical structures lead to fewer weights
- Goals
  - Slimmer CNN with higher generalization accuracy
  - Fast training algorithms (under development)

- General data classification
  - Linearly separable sets: logistic regressions
  - Nonlinear feature mapping: from nonlinear separable to linearly separable
  - Deep learning: feature mapping given by deep neural networks
- Image classification
  - Constraint linear model for data-feature
  - Feature extraction by iterative methods
- MgNet
  - A special CNN by making minor modification of multigrid
  - More mathematical structures lead to fewer weights
- Goals
  - Slimmer CNN with higher generalization accuracy
  - Fast training algorithms (under development)
  - More mathematical theory for deep learning