

Finite Neuron Method
Iterative Methods and Frequency Principle

Jinchao Xu

KAUST and Penn State

xu@multigrid.org

CIRM, July 18th, 2023

CEMRACS 2023 Summer School

KAUST Baseline Research Fund

Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

Iterative methods for $Au = f$

$$u^0, u^1, \dots, u^{m-1} \longrightarrow u^m$$

Basic ideas:

- 1 Form the residual: $r = f - Au^{m-1}$
- 2 Solve the residual equation $Ae = r$ approximately $\hat{e} = Br$ with $B \approx A^{-1}$
- 3 Update $u^m = u^{m-1} + \hat{e}$

Linear iterative method:

$$u^m = u^{m-1} + B(f - Au^{m-1}) \quad (1)$$

Let $A = L + D + U$. Thus,

- Jacobi iteration: $B = D^{-1}$,
- Gauss-Seidel iteration: $B = (L + D)^{-1}$.

Examples: basic iterative methods

- Richardson iteration:

$$u^m = u^{m-1} + \omega(f - Au^{m-1}), \quad m = 1, 2, \dots, \quad (2)$$

- Modified Jacobi:

$$u^m = u^{m-1} + \omega D^{-1}(f - Au^{m-1}), \quad m = 1, 2, \dots, \quad (3)$$

- Modified Gauss-Seidel:

$$u^m = u^{m-1} + (\omega^{-1}D + L)^{-1}(f - Au^{m-1}), \quad m = 1, 2, \dots, \quad (4)$$

Thus, the iterative method converges if the following operator is SPD:

$$(B')^{-1} + B^{-1} - A = \begin{cases} 2\omega^{-1} - A > 0 & \text{if } 0 < \omega < \frac{2}{\rho(A)} \\ 2\omega^{-1}D - A > 0 & \text{if } 0 < \omega < \frac{2}{\rho(D^{-1}A)} \\ (2 - \omega)\omega^{-1}D > 0 & \text{if } 0 < \omega < 2 \end{cases}$$

Richardson;
Modified Jacobi;
Modified G.-S.

Iterative methods as gradient descent (GD)

If A is SPD, then we have the following equivalence:

$$Au = f \iff \min \underbrace{\frac{1}{2} u^T Au - f^T u}_{J(u)}$$

- Richardson for $Au = f \Leftrightarrow$ Gradient descent for $f(u)$

$$u^m = u^{m-1} + \eta(f - Au^{m-1}) = u^{m-1} - \eta \nabla J(u^{m-1})$$

- Jacobi for $Au = f \Leftrightarrow$ Scaled gradient descent for $f(u)$

$$u^m = u^{m-1} + \eta D^{-1}(f - Au^{m-1}) = u^{m-1} - \eta [\text{diag}(A)]^{-1} \nabla J(u^{m-1})$$

- Gauss–Seidel for $Au = f \Leftrightarrow$ Preconditioned gradient descent for $f(u)$

$$u^m = u^{m-1} + (\eta D + L)^{-1}(f - Au^{m-1}) = u^{m-1} - P \nabla J(u^{m-1}), \quad P = (\eta D + L)^{-1}$$

GD for a nearly singular system

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

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

Note that $\sigma(A_0) = \{3, 1, 0\}$. Apply scaled gradient descent method with $\|A_\epsilon u^k - g\| \leq 10^{-8}$:

ϵ	# of iter = m
1.	37
10^{-1}	236
10^{-2}	1, 918
10^{-3}	16, 115
10^{-4}	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; Lee, Wu, Xu and Zikatanov 2007

Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

Model problem and frequencies

$$\begin{cases} -u''(x) = f, & x \in (0, 1), \\ u(0) = 0, \quad u'(1) = 0. \end{cases}$$

Consider eigenvalue problem

$$\begin{cases} -u_k''(x) = \lambda_k u_k(x), & x \in (0, 1), \\ u_k(0) = 0, \quad u'_k(1) = 0, \end{cases}$$

We have

$$\lambda_k = (k - \frac{1}{2})^2 \pi^2, \quad u_k(x) = \sin\left((k - \frac{1}{2})(\pi x)\right), \quad k = 1, 2, 3, \dots.$$

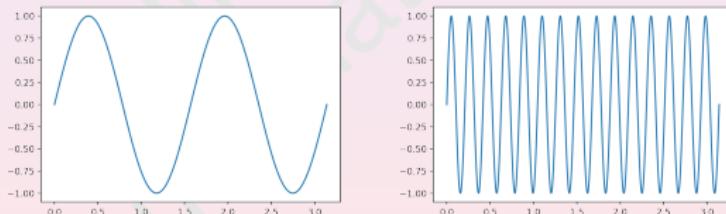


Figure: Frequencies with smaller k and larger k

Frequency bias of GD

For any SPD matrix $A \in \mathbb{R}^{n \times n}$ and vector $b \in \mathbb{R}^n$, the gradient descent method solving

$$\min_{v \in \mathbb{R}^n} I(v) \quad \text{with} \quad I(v) = \frac{1}{2} v^T A v - v^T b$$

reads as

$$v^{\ell+1} = v^\ell - \eta \nabla_v I(v^\ell), \quad \ell = 0, 1, \dots,$$

with initial guess v^0 .

Since that $\nabla_v I(v) = Av - b$, we have

$$v^{\ell+1} = v^\ell - \eta(Av^\ell - b), \quad \ell = 0, 1, \dots.$$

Convergence of GD with $\eta = \frac{1}{\lambda_{n,A}}$

$$v - v^\ell = \sum_{k=1}^n \alpha_k \left(1 - \frac{\lambda_{k,A}}{\lambda_{n,A}}\right)^\ell \xi_A^k$$

where $\xi_A^k, k = 1, 2, \dots, n$ are the eigenvector of A .

- Fast on algebraic frequencies corresponding to large eigenvalues.
- Slow on algebraic frequencies corresponding to small eigenvalues.

H^1 fitting

Given $f \in L^2(\Omega)$

$$J(v) = \frac{1}{2}a(v, v) - (f, v)$$

Consider to fit a target function $u(x) \in V$ by a function $u_h(x) \in V_h$.

$$a(u, v) = (u', v')_{L^2}, \quad H^1 \text{ fitting.}$$

Finite element: Piecewise linear functions

- Uniform grid \mathcal{T}_h

$$0 = x_0 < x_1 < \cdots < x_{N+1} = 1, \quad x_j = \frac{j}{N+1} \quad (j = 0 : N+1).$$

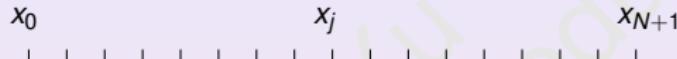


Figure: 1D uniform grid

- Linear finite element space

$$V_h = \{v_h : v \text{ is continuous and piecewise linear w.r.t. } \mathcal{T}_h, v_h(0) = 0\}.$$

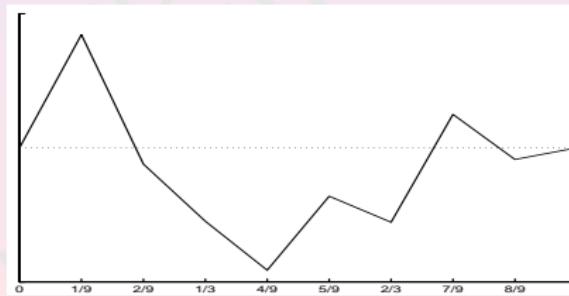


Figure: Typical finite element functions.

Two basis of the finite element space V_h

- Hat basis:

$$\varphi(x) = \begin{cases} x & x \in [0, 1] \\ 2 - x & x \in [1, 2] \\ 0, & \text{others} \end{cases}.$$

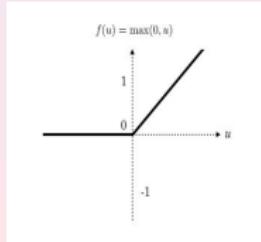
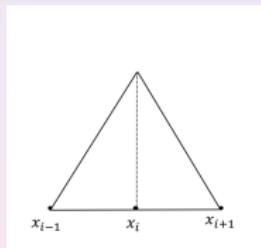
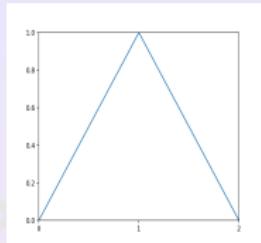
$$\varphi_i(x) = \varphi\left(\frac{x - x_{i-1}}{h}\right) = \varphi(w_h x + b_i).$$

with $w_h = \frac{1}{h}$, $b_i = \frac{-x_{i-1}}{h}$.

- ReLU basis: $\text{ReLU}(x) = \max(0, x)$ and

$$r_i(x) = \text{ReLU}\left(\frac{x - x_{i-1}}{h}\right) = \text{ReLU}(w_h x + b_i)$$

- $V_h = \text{span} \{\text{ReLU}(w_h x + b_i)\} = \text{span} \{\varphi(w_h x + b_i)\}$



Hat and ReLU bases on a uniform grid

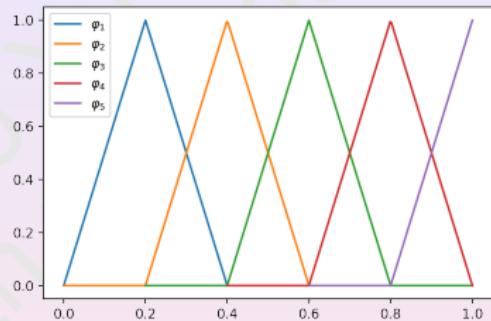
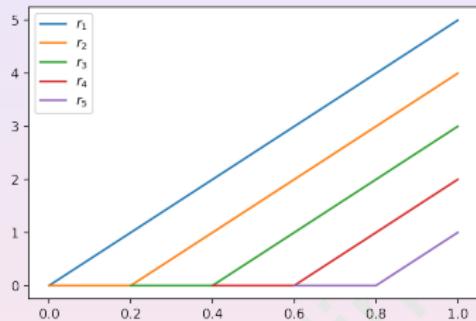


Figure: Left: ReLU bases. Right: Hat bases.

H^1 -fitting

Stiffness matrix for Hat basis A_{Hat} is given by

$$A_{Hat} = \left(\int_0^1 \varphi_j'(x) \varphi_i'(x) dx \right) = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 1 \end{pmatrix} \in \mathbb{R}^{n \times n}. \quad (5)$$

Lemma

The eigenvalues $\lambda_{k,A_{hat}}$, $1 \leq k \leq n$ and corresponding eigenvectors

$\xi_{A_{hat}}^k = (\xi_{A_{hat},j}^k)_{j=1}^n$, $1 \leq k \leq n$ of A_{hat} are

$$\lambda_{k,A_{hat}} = 4(n+1)^2 \sin^2 \frac{(k-\frac{1}{2})\pi}{2n+1} \approx \lambda_k,$$

$$\xi_{A_{hat},j}^k = \sin \left((k - \frac{1}{2})\pi x_j \right) \text{ with } x_j = \frac{2j}{2n+1}, 1 \leq j \leq n.$$

Frequency bias for hat basis

① GD for stiffness matrix of Hat bases:

- ▶ $\|\alpha - \alpha_\ell\| = \mathcal{O}((1 - cn^{-2})^\ell)$.
- ▶ Low frequency converges slowly: $\mathcal{O}((1 - cn^{-2})^\ell)$.
- ▶ High frequency converges fast: $\mathcal{O}(1 - \delta)^\ell$ for $0 < \delta < 1$.

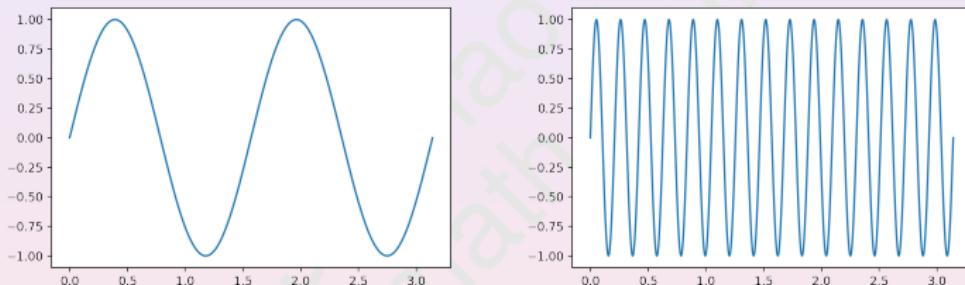


Figure: Low and high frequencies

Ref: Q. Hong, Q. Tan, J.W. Siegel, and J. Xu. On the activation function dependence of the spectral bias of neural networks. arXiv:2208:04924 (2022).

Relationship between ReLU basis and hat basis

- We have

$$\varphi(x) = \mathbf{1} \cdot \text{ReLU}(x) - \mathbf{2} \cdot \text{ReLU}(x-1) + \mathbf{1} \cdot \text{ReLU}(x-2). \quad (6)$$

- Let $\Psi(x) = (r_1(x), r_2(x), \dots, r_n(x))^T$ and $\Phi(x) = (\varphi_1(x), \varphi_2(x), \dots, \varphi_n(x))^T$. Then

$$\Phi = C\Psi, \quad (7)$$

where

$$C = \frac{1}{h^2} \begin{pmatrix} 1 & -2 & 1 & & & \\ & 1 & -2 & 1 & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1 & -2 & 1 \\ & & & & 1 & -2 \\ & & & & & 1 \end{pmatrix}. \quad (8)$$

Spectral analysis of H^1 -fitting

Stiffness matrix A_{ReLU} is given by

$$A_{ReLU} = \left(\int_0^1 r_j'(x) r_i'(x) dx \right) = h^2 \begin{pmatrix} n & n-1 & n-2 & \cdots & 1 \\ n-1 & n-1 & n-2 & \cdots & 1 \\ n-2 & n-2 & n-2 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{n \times n}. \quad (9)$$

Theorem

$$A_{ReLU} = E A_{\text{Hat}}^{-1} E^{-1} \quad \text{with} \quad E = \begin{pmatrix} & & & 1 \\ & \dots & & 1 \\ 1 & & & \end{pmatrix}. \quad (10)$$

The eigenvalues $\lambda_{k,A_{ReLU}}$, $1 \leq k \leq n$ and the corresponding eigenvectors $\xi_{A_{ReLU}}^k$, $1 \leq k \leq n$ of A_{ReLU} are as follows:

$$\lambda_{k,A_{ReLU}} = \lambda_{n+1-k,A_{\text{Hat}}}^{-1}, \quad \xi_{A_{ReLU}}^k = E \xi_{A_{\text{Hat}}}^{n+1-k}. \quad (11)$$

Spectral analysis of H^1 -fitting

Proof:

By direct computation, we have

$$A_{ReLU} = h^2 A_1, \quad \text{with} \quad A_1 = \begin{pmatrix} n & n-1 & n-2 & \cdots & 1 \\ n-1 & n-1 & n-2 & \cdots & 1 \\ n-2 & n-2 & n-2 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 1 \end{pmatrix} \quad (12)$$

and

$$A_1^{-1} = \begin{pmatrix} 1 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}. \quad (13)$$

By inspection, we have

$$\frac{1}{h^2} A_1^{-1} = \begin{pmatrix} & & & 1 \\ & \dots & 1 & \\ 1 & & & \end{pmatrix} A_{hat} \begin{pmatrix} & & & 1 \\ 1 & \dots & 1 & \end{pmatrix}. \quad (14)$$

Spectral analysis of H^1 -fitting: eigenvectors

Theorem

Let $e_k(x) = \xi_{A_{ReLU}}^k \cdot \Psi(x) = \sum_{i=1}^n \xi_{A_{ReLU},i}^k r_i(x)$, then we have

$$e_k(x_j) = \sin \frac{\pi t_k}{2} + \sin \left((n - k + \frac{1}{2})\pi t_j - \frac{\pi t_k}{2} \right) \text{ and } t_j = \frac{2j}{2n+1}.$$

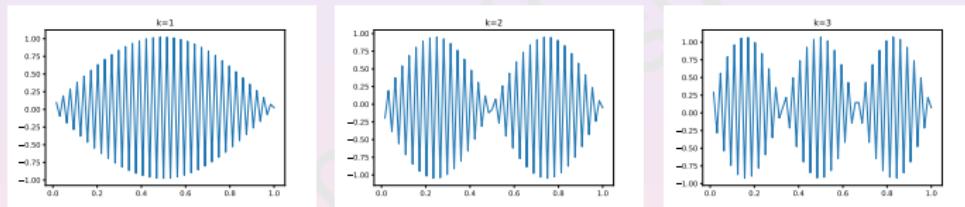


Figure: Functions: $e_1(x)$, $e_2(x)$ and $e_3(x)$.

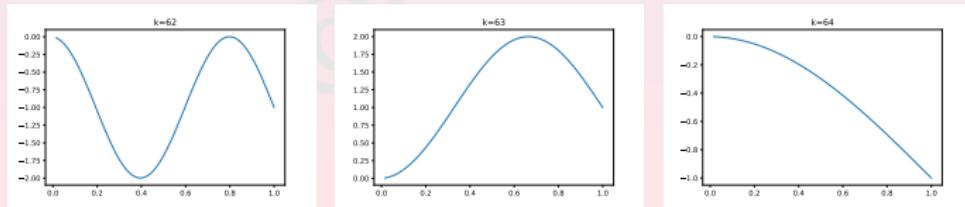


Figure: Functions: $e_{62}(x)$, $e_{63}(x)$ and $e_{64}(x)$.

Frequency bias for ReLU basis

1 GD for the stiffness matrix of ReLU basis:

- ▶ $\|\alpha - \alpha_\ell\| = \mathcal{O}((1 - cn^{-2})^\ell)$.
- ▶ Low frequency converges fast: $\mathcal{O}(1 - \delta)^\ell$ for $0 < \delta < 1..$
- ▶ High frequency converges slowly: $\mathcal{O}((1 - cn^{-2})^\ell)$.

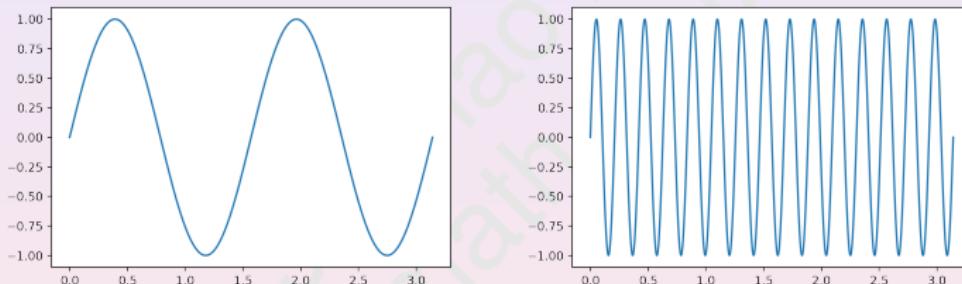


Figure: Low and high frequencies

Ref: Q. Hong, Q. Tan, J.W. Siegel, and J. Xu. On the activation function dependence of the spectral bias of neural networks. arXiv:2208:04924 (2022).

GD for H^1 -fitting

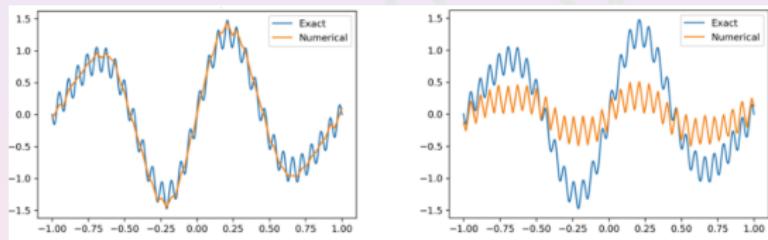


Figure: Results of Hat basis.

Figure: Results of ReLU basis.

Frequency bias for training neural network

- A special case of neural network functions: linear problems
- The frequency principle is still true for nonlinear problems with neural network functions.



Poisson equation. Left: ReLU activation. Right: Hat activation.

Activation dependence of training neural network

ReLU neural networks

- Prioritize learning low frequency modes in H^1 fitting
- Prioritize learning low frequency modes in L^2 fitting
- Training loss decreases slowly in L^2 fitting due to the frequency bias

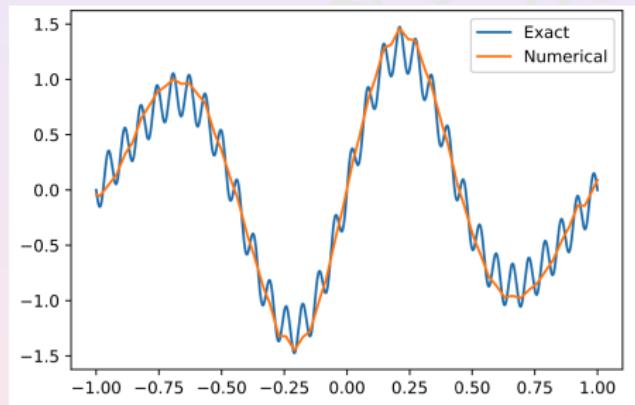
Hat neural networks

- Prioritize learning the high frequency modes in H^1 fitting
- Learn both the low frequency and high frequency modes in L^2 fitting
- Training loss decreases very fast in L^2 fitting since there is no frequency bias

● Rahaman, N., Baratin, A., Arpit, D., Draxler, F., Lin, M., Hamprecht, F. A., Bengio, Y. & Courville, A. (2019), Xu, Z. (2018), Cai, W. & Xu, Z. (2019), Xu, Z., Zhang, Y., Luo, T., Xiao, Y. & Ma, Z (2019), Hong, Q., Seigel, J., Tan, Q., & Xu, J. (2022).

"Convergence" of SGD or Adam Algorithms for NN-based PDE Solver

- SGD and Adam converge rather quickly for low frequency, and hence capture the "profile" of physical solutions reasonably well.



- This provides a theoretical explanation of the success of methods such as PINN.

"Non-convergence" of SGD or Adam Algorithms for NN-based PDE Solver

$$u_n = \arg \min_{v_n \in \Sigma_n^{ReLU}} J(v_n). \quad (15)$$

We have proved that one can NOT use SGD or Adam to numerically find $\tilde{u}_n \approx u_n$ such that

$$\|u - \tilde{u}_n\| \leq cn^{-\alpha} \quad (16)$$

for any $\alpha > 0$ for large n .

- H^1 -fitting by ReLU NN:

$$1 - cn^{-2}$$

Taking $n = 10^6$: how many iterations do we need such that

$$(1 - 10^{-12})^k \leq 10^{-7} \quad (17)$$

- ▶ $k \geq 1.61 \times 10^{25}$
- ▶ 32 years for the fastest computer in the world (Frontier, 1.1 EFLOPS)

New training algorithms are required to achieve sufficiently good accuracy!

Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

GD for a nearly singular system

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

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

Note that $\sigma(A_0) = \{3, 1, 0\}$. Apply scaled gradient descent method with $\|A_\epsilon u^k - g\| \leq 10^{-8}$:

ϵ	# of iter = m
1.	37
10^{-1}	236
10^{-2}	1,918
10^{-3}	16,115
10^{-4}	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; Lee, Wu, Xu and Zikatanov 2007

Remedy of GD: Expanded system (Over-parametrization)

Write $u \in \mathbb{R}^3 = u_1 e_1 + u_2 e_2 + u_3 e_3$ as

$$u = \underline{u}_1 e_1 + \underline{u}_2 e_2 + \underline{u}_3 e_3 + \underline{u}_4 p = P\underline{u}$$

where

$$P = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}, \quad p = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \in \ker(\mathbf{A}_0).$$

Namely, we consider the coarse level with “lowest” frequency $p \in \ker(\mathbf{A}_0)$.

The equation $A_\epsilon u = g$ becomes

$$A_\epsilon P\underline{u} = g \iff (\mathbf{P}^T A_\epsilon \mathbf{P})\underline{u} = \mathbf{P}^T g,$$

leading to a semi-definite system:

$$\begin{pmatrix} 1+\epsilon & -1 & 0 & \epsilon \\ -1 & 2+\epsilon & -1 & \epsilon \\ 0 & -1 & 1+\epsilon & \epsilon \\ \epsilon & \epsilon & \epsilon & 3\epsilon \end{pmatrix} \underline{u} = \begin{pmatrix} -1 \\ -1 \\ 2 \\ 0 \end{pmatrix}.$$

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

Over-parametrization \iff Two-level methods



$$V_1 = \mathbb{R}^3 \xrightarrow[p^T]{\text{Pooling}} V_2 = \mathbb{R}$$

- 1 Initialization of inputs

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

- 2 Iterate:

- 1 One step of GD method on V_1

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

- 2 Consider $A_1 \theta_1 = r_1 \equiv g_1 - A_1 u_1$ and "pool" it to V_2 and solve it:

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

with

$$A_2 = p^T A_1 p = 3\epsilon, \quad g_2 = p^T r_1$$

- 3 update $u_1 \leftarrow u_1 + p u_2$.

Multilevel method: over-parameterization using multilevel frame

Multilevel frame over-parameterization \iff Multigrid

$$V_J \subset V_{J-1} \subset V_{J-2} \subset \dots \subset V_1 \equiv V.$$

Frame:

$$\{\phi_{k,i} : i = 1 : n_k, k = 1 : J\}$$

Frame expansion (not unique):

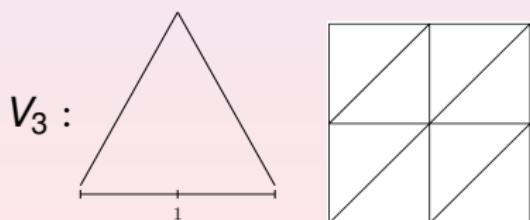
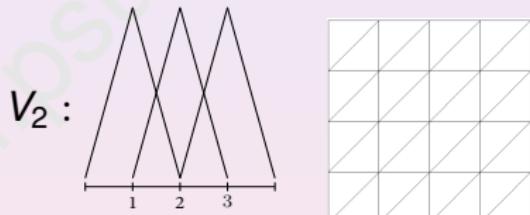
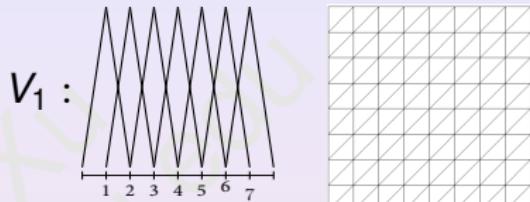
$$u_h = \sum_{k=1}^J \sum_{x_{k,i} \in N_k} \mu_{i,k} \phi_{k,i}.$$

Expanded system

$$\underline{A}\underline{\mu} = \underline{b}$$

where \underline{A} is the frame stiffness matrix

$$\underline{A} = \left((\phi_{k,i}, \phi_{l,j})_A \right) \in R^{N \times N}, \quad N = \sum_{k=1}^J n_k$$



An equivalent formulation of multigrid

Smoothing and restriction

- For $k = 1 : J$
 - For $i = 1 : n_k$

$$u_k \leftarrow u_k + S_k * (g_k - A_k * u_k).$$

- Form restricted residual and set initial guess:

$$u_{k+1,0} \leftarrow \Pi_k^{k+1} u_k,$$

$$g_{k+1} \leftarrow R_k *_2 (g_k - A_k * u_k) + A_{k+1} * u_{k+1}^0.$$

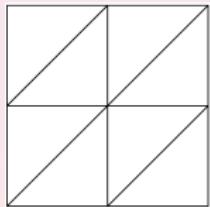
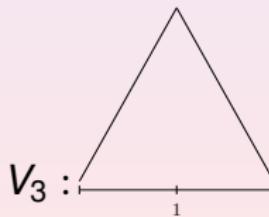
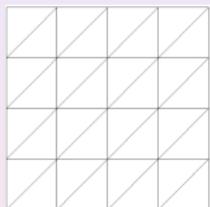
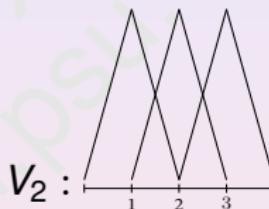
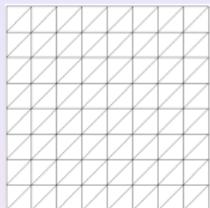
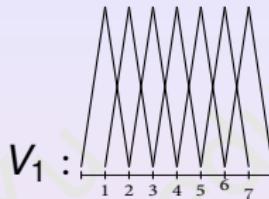
Prolongation with post-smoothing

- For $k = 1 : J - 1$

$$u_k \leftarrow u_k + R_k *_2^\top (u_{k+1} - u_{k+1}^0).$$

- For $i = 1 : n'_k$

$$u_k \leftarrow u_k + S'_k * (g_k - A_k * u_k)$$

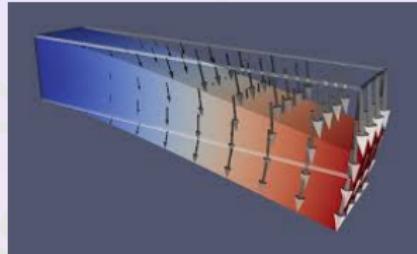


$$\phi_{3,1} = \frac{1}{2}\phi_{2,1} + \phi_{2,2} + \frac{1}{2}\phi_{2,3}$$

2D linear system on a uniform grid

- Model Problem:

$$\begin{aligned}-\Delta u &:= -(u_{xx} + u_{yy}) = g, \text{ in } \Omega, \\ u &= 0 \text{ on } \partial\Omega, \quad \Omega = (0, 1)^2.\end{aligned}$$



- Discrete case:

$$4u_{i,j} - u_{i+1,j} - u_{i-1,j} - u_{i,j+1} - u_{i,j-1} = g_{i,j}, \quad (18)$$

with

$$A * u = g, \quad \text{for } A = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{pmatrix}, \quad (19)$$

GD for the over-parameterized multilevel system

Original system in terms of a basis

$$Au = g, \quad A = ((\nabla \phi_j, \nabla \phi_i)) \in \mathbb{R}^{n_1 \times n_1}$$

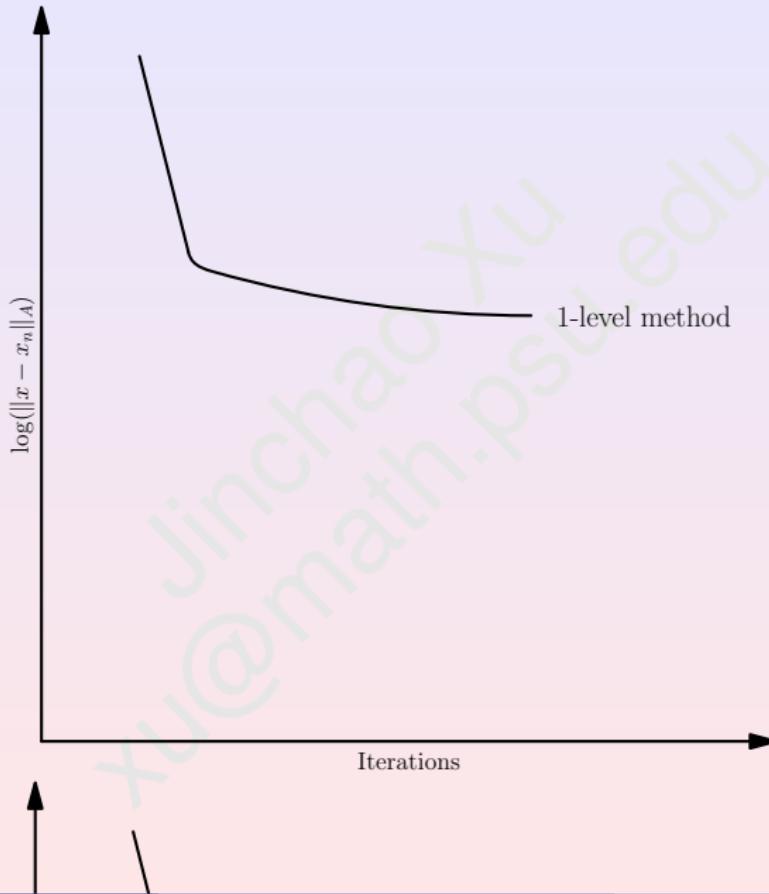
Expanded system in terms of a multilevel frame (over-parameterization):

$$\underline{A}\underline{u} = \underline{g}, \quad \underline{A} = ((\nabla \phi_{\ell,j}, \nabla \phi_{k,i})) \in \mathbb{R}^{N \times N}, \quad N = \sum_{k=1}^J n_k$$

Solve \underline{A} by Gradient Descent:

Size	GD for A	GD for \underline{A}
4^2	56	16
16^2	954	21
64^2	14,758	26
256^2	223,630	26
1024^2	>1,000,000	26

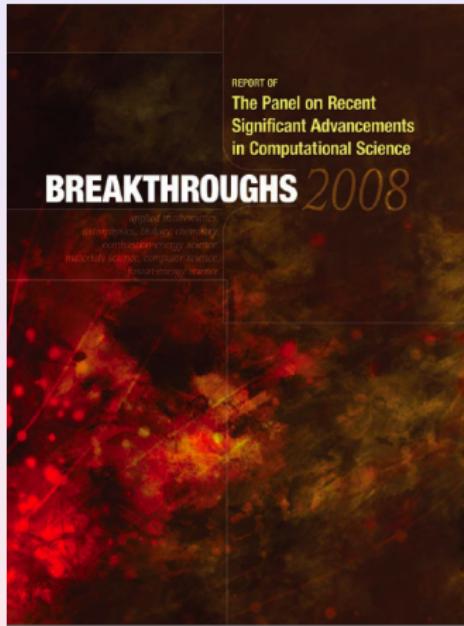
Performance of multigrid:



A success story: HX preconditioner

(Hiptmair and Xu 2005, 2007, Xu 2014)

A DOE report to the U. S. Congress



Novel Solver Enables Scalable Electromagnetic Simulations

ABSTRACT: A team based at Lawrence Livermore National Laboratory has developed the first provably scalable solver code for linear systems of equations of electromagnetic equations that are fundamental to numerous areas of physics and engineering. This new software technology enables researchers to solve larger computational problems with greater accuracy.

able to handle complex geometries and problems with large jumps in the material coefficients. In contrast, some of the old codes can only handle simple geometries and accurate results when faced with systems of equations that have relatively smooth, homogeneous electromagnetic properties, which are common in engineering.

AMG works by reducing the original system of equations into smaller subproblems that can be individually handled using classical multigrid or other regular arrangement of points. The overall system of equations is tackled by a solid theoretical framework. The theory behind the AMG algorithm shows the fundamental mathematical research can lead to important software advances in scientific computing. In addition, James and Tasso Ruge of LLNL, co-authors of the paper, are currently working on a follow-up paper with Timothy and Rolf Kippel of CFD, Zürich.

Electromagnetic simulations have a wide range of physical and engineering applications, from the design and development of semiconductor chips, medical devices, and electronic generators. As the complexity of problems grows, researchers need to be able to take advantage of the power of modern computing power.

The AMS solver's scalability is the result of its availability. Specifically, AMS exhibits "near" parallel efficiency, meaning that the execution time is proportional to the problem size and processor workload simultaneously increase. The new algorithm is

AMS computation of a higher-fidelity problem used in pulsed-power experiments. Image courtesy of LLNL, UC.

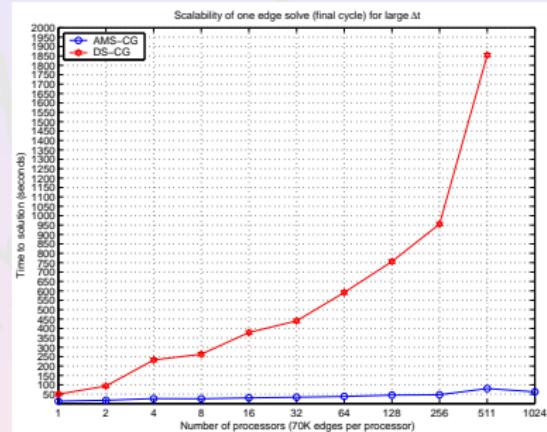
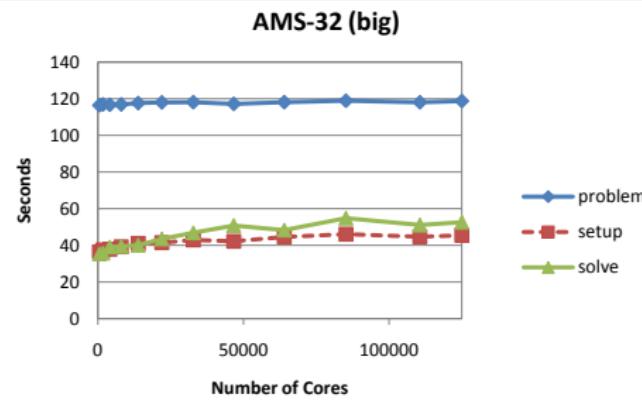
Large-scale electromagnetic simulations of complex structures are often bottlenecked by the number of unknowns. A team of researchers at LLNL, a newly developed algorithm, known as the auxiliary-space Multigrid (AMG) method, has developed faster solution techniques by as much as 25 times, giving researchers a significant advantage in solving large-scale problems. They seek to answer inevitably larger engineering questions.

The solver's scalability is the result of its availability. Specifically, AMS exhibits "near" parallel efficiency, meaning that the execution time is proportional to the problem size and processor workload simultaneously increase. The new algorithm is

"AMS is a perfect example of how fundamental mathematical research can lead to important software advances in high-performance computing."

One application: LLNL

Scalability of HX preconditioner to 125,000 cores



Left: Auxiliary-space Maxwell Solver. Total problem size is 12 billion.

Right: Scalability (70K edge unknowns per processor)

Ref: A. Baker, R. Falgout, T. Kolev, and U. Yang 2012

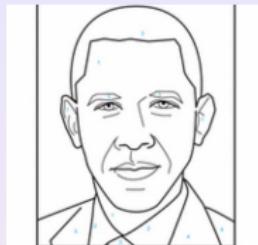
Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

A mathematical model of feature extraction



$$g \xleftarrow{A*u=g} u$$



Model: given an image g , find its feature u satisfying

$$A * u = g \quad (20)$$

with a constraint

$$u \geq 0. \quad (21)$$

Questions:

- 1 What is A ? (to be trained ...) \leftarrow data
- 2 How to solve (20)? (iterative methods) \leftarrow scientific computing

Ref: J. He and J. Xu 2019, J. He, J. Xu, L. Zhang and J. Zhu 2021.

Data and feature spaces

Partial differential equations

$$-\Delta u = g. \quad (22)$$

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

$$A * u = g. \quad (23)$$

Main idea:

- Use a geometric multigrid method for PDE (22) to solve the data-feature equation (23)!

Ref: J. He and J. Xu 2019, J. He, J. Xu, L. Zhang and J. Zhu 2021.

Iterative methods for $Au = g$: residual correction

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

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

A basic iterative method:

$$u^i = u^{i-1} + B^i(g - Au^{i-1}) \quad (24)$$

An example: $A = D + L + U$ with $D = \text{DIAG}(A)$

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

Iterative schemes for the constrained linear model

- 1 Recall iterative methods without constraint

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

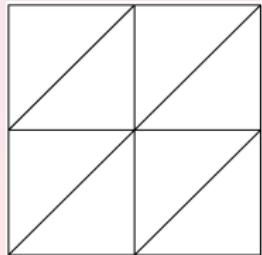
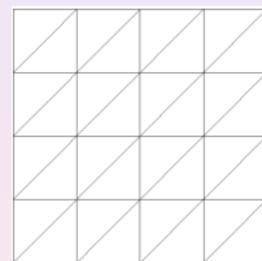
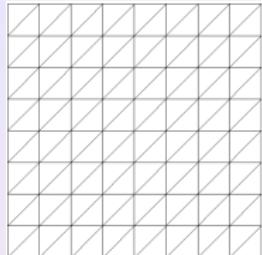
- 2 Image classification [$\sigma = \text{ReLU}$: dropping the negative values]

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

or, in terms of residual

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

MgNet: a “trained” multigrid method



Initialization of inputs:

$$g_1 \leftarrow \theta * g, \quad u_1 \leftarrow 0$$

Smoothing and restriction

- For $\ell = 1 : J$
 - ▶ For $i = 1 : v_\ell$

$$u^\ell \leftarrow u^\ell + \sigma \circ B^\ell * \sigma(g^\ell - A^\ell * u^\ell).$$

- ▶ Form restricted residual and set initial guess:

$$u^{\ell+1,0} \leftarrow \sigma \circ \Pi_\ell^{\ell+1} *_2 u^\ell,$$

$$g^{\ell+1} \leftarrow \sigma \circ R_\ell^{\ell+1} *_2 (g^\ell - A^\ell * u^\ell) + A^{\ell+1} * u^{\ell+1,0}.$$

Output: ...

$$\phi(g) \leftarrow u^J$$

Ref: He, J. & Xu, J. (2019)

Batch normalization

DNN as example:

- Original model

$$\begin{cases} f^1(x^i) &= W^1 x^i, \\ f^\ell &= W^\ell \sigma(f^{\ell-1}), \quad \ell = 2, \dots, L. \end{cases} \quad (25)$$

- Batch normalization:

- For t -th step of SGD training on mini-batch \mathcal{B}_t
- For the ℓ -th layer

$$\mu_{\mathcal{B}_t}^\ell \leftarrow \frac{1}{m} \sum_{i \in \mathcal{B}_t} f^\ell(x^i) \quad \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}_t}^\ell \leftarrow \frac{1}{m} \sum_{i \in \mathcal{B}_t} (f^\ell(x^i) - \mu_{\mathcal{B}_t}^\ell)^2 \quad \text{mini-batch variance}$$

$$\tilde{f}^\ell(x) \leftarrow \frac{f^\ell(x) - \mu_{\mathcal{B}_t}^\ell}{\sqrt{\sigma_{\mathcal{B}_t}^\ell + \epsilon}} \quad \text{normalize}$$

$$\text{BN}_{\mathcal{B}_t}(f^\ell) \leftarrow \gamma^\ell \tilde{f}^\ell(x) + \beta^\ell \quad \text{scale and shift}$$

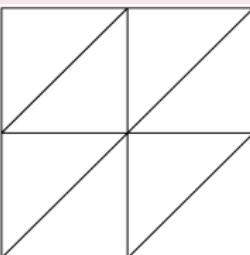
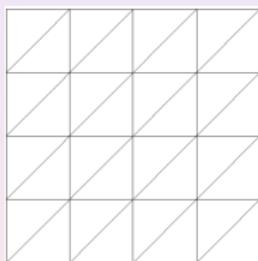
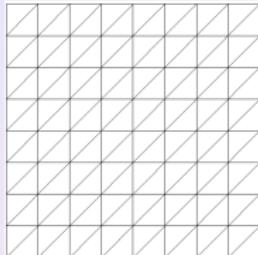
- Model with batch normalization

$$\begin{cases} \tilde{f}^1(x^i) &= W^1 x^i, \\ \tilde{f}^\ell &= W^\ell \sigma_{\text{BN}}(\tilde{f}^{\ell-1}), \quad \ell = 2, \dots, L, \end{cases} \quad (27)$$

where

$$\sigma_{\text{BN}}(f) = \sigma(\text{BN}_{\mathcal{B}_t}(f)). \quad (28)$$

Practical MgNet with batch normalization



Initialization of inputs:

$$g_1 \leftarrow \theta * g, \quad u_1 \leftarrow 0$$

Smoothing and restriction

- For $\ell = 1 : J$
 - ▶ For $i = 1 : v_\ell$

$$u^\ell \leftarrow u^\ell + \sigma_{BN} \circ B^\ell * \sigma_{BN}(g^\ell - A^\ell * u^\ell).$$

- ▶ Form restricted residual and set initial guess:

$$u^{\ell+1,0} \leftarrow \sigma_{BN} \circ \Pi_\ell^{\ell+1} *_2 u^\ell,$$

$$g^{\ell+1} \leftarrow \sigma_{BN} \circ R_\ell^{\ell+1} *_2 (g^\ell - A^\ell * u^\ell) + A^{\ell+1} * u^{\ell+1,0}.$$

- ...

Output:

$$\phi(g) \leftarrow u^J$$

Denote: $\sigma = \sigma_{BN}$ for CNNs in image classification by default.

- Ref: He, J. & Xu, J. (2019)

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

$$r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}), \quad i = 1 : v_\ell, \quad (29)$$

where

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

Proof.

$$\begin{aligned} u^{\ell,i} &= u^{\ell,i-1} + \sigma \circ B^{\ell,i} * \sigma(g^\ell - A^\ell * u^{\ell,i-1}), \\ \Rightarrow A^\ell * u^{\ell,i} &= A^\ell * u^{\ell,i-1} + A^\ell * \sigma \circ B^{\ell,i} * \sigma(g^\ell - A^\ell * u^{\ell,i-1}), \\ \Rightarrow g^\ell - A^\ell * u^{\ell,i} &= g^\ell - A^\ell * u^{\ell,i} - A^\ell * \sigma \circ B^{\ell,i} * \sigma(g^\ell - A^\ell * u^{\ell,i-1}), \\ \Rightarrow r_\ell^i &= r^{\ell,i} = r^{\ell,i-1} + A^{\ell,i} * \sigma \circ B^{\ell,i} * \sigma(r^{\ell,i-1}). \end{aligned} \quad (30)$$



Modified Pre-act ResNet, ResNet

Modified Pre-act ResNet – Pre-act ResNet- A^ℓ

$$r^{\ell,i} = r^{\ell,i-1} + A^\ell * \sigma \circ B^{\ell,i} * \sigma(r^{i-1}). \quad (31)$$

Modified ResNet – ResNet- A^ℓ

$$r^{\ell,i} = \sigma \left(r^{\ell,i-1} + A^\ell * \sigma \circ B^{\ell,i} * r^{i-1} \right) \quad (32)$$

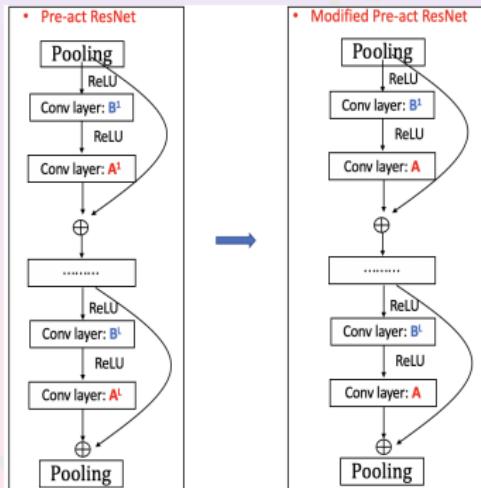


Figure: Diagram of modified models.

Modify (pre-act) ResNet numerical experiments

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

Model	CIFAR10	CIFAR100	# Parameters
ResNet18- $A^{\ell,i}$ - $B^{\ell,i}$	94.22	76.08	11M
ResNet18- A^ℓ - $B^{\ell,i}$	94.34	76.32	8.1M
ResNet18- $A^{\ell,i}$ - B^ℓ	93.95	74.23	9.7M
ResNet18- A^ℓ - B^ℓ	93.30	74.85	6.6M
pre-act ResNet18- $A^{\ell,i}$ - $B^{\ell,i}$	94.31	76.33	11M
pre-act ResNet18- A^ℓ - $B^{\ell,i}$	94.54	76.43	8.1M
pre-act ResNet18- $A^{\ell,i}$ - B^ℓ	93.96	74.45	9.7M
pre-act ResNet18- A^ℓ - B^ℓ	93.63	74.46	6.6M
ResNet34- $A^{\ell,i}$ - $B^{\ell,i}$	94.43	76.31	21M
ResNet34- A^ℓ - $B^{\ell,i}$	94.78	76.44	13M
ResNet34- $A^{\ell,i}$ - B^ℓ	93.98	74.48	15M
ResNet34- A^ℓ - B^ℓ	93.55	74.46	6.7M
pre-act ResNet34- $A^{\ell,i}$ - $B^{\ell,i}$	94.70	77.38	21M
pre-act ResNet34- A^ℓ - $B^{\ell,i}$	94.91	77.41	13M
pre-act ResNet34- $A^{\ell,i}$ - B^ℓ	94.08	75.32	15M
pre-act ResNet34- A^ℓ - B^ℓ	94.01	74.12	6.7M

MgNet: From multigrid to CNN

Multigrid:

- A^ℓ, B^ℓ, R^ℓ are all given a priori

CNN:

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

CNN versus multigrid: classic approaches versus MgNet

① 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!

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 comparison 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

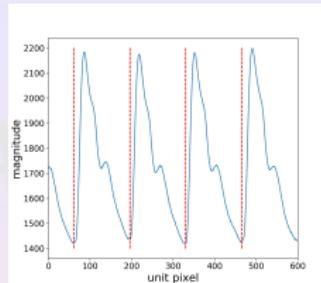
Table: The comparison of MgNet and classical CNN on Cifar100

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 comparison of MgNet and classical CNN on ImageNet.

Application: Pulse-Feeling

An ancient technique in Traditional Chinese Medication (TCM)



- It has been widely believed and claimed to be accurate
- No record of clinical trials nor quantitative studies
- it is a valid technique or it is ... ?

Deep learning for diagnosing pregnancy

model	test accuracy(%)	AUC(%)	size
ResNet	84.73	89.66	232,642
MgNet	84.68	91.04	3,450
SVM	78.08	71.32	
logistic regression	79.10	74.27	

Ref: Chen, Huang, Hao and Xu 2019

A summary of MgNet

- J. Xu, Deep Neural Networks and Multigrid Methods, (Lecture Notes at Penn State and KAUST), 2023.
- J. He, J. Xu. MgNet: A Unified Framework of Multigrid and Convolutional Neural Network. Sci China Math, 2019, 62: 1331-1354.
- Y. Chen, B. Dong, J. Xu. Meta-MgNet: Meta Multigrid Networks for Solving Parameterized Partial Differential Equations Image Classification. Journal of Computational Physics, 2022, 455.
- J. He, L. Li, J. Xu. Approximation Properties of Deep ReLU CNNs. Research in the Mathematical Sciences, 2022, 9(3).
- J. He, J. Xu, L. Zhang, J. Zhu. An interpretive constrained linear model for ResNet and MgNet. Neural Networks, 2023, 162: 384-392.

-
- 1 A uniform framework for understanding and designing CNNs: ResNet, U-Net, DenseNet ...
 - 2 Construct the new
 - ▶ reduce the free parameters 1%, .1%, .01%...?
 - ▶ increase the generalization accuracy,
 - ▶ accelerate the training speed,
 - ▶ extend to general graph models,
 - ▶ ...
 - 3 Applications:
 - ▶ image problems,
 - ▶ time series forecasting,
 - ▶ numerical PDEs, in particular for operator learning,
 - ▶ ...

Challenging Objective: MgNet vs Transformer

Can MgNet outperform Transformer?

Model	Type	Accuracy	Parameters
DeiT-Small	Transformer	79.8	22.1M
PVT-Small	Transformer	79.8	24.5M
ConvMixer	Transformer	80.2	21.1M
CrossViT-Small	Transformer	81.0	26.7M
Swin-Tiny	Transformer	81.2	28.3M
CvT-13	Transformer	81.6	20.0M
CoAtNet-0	Transformer	81.6	25.0M
CaiT-XS-24	Transformer	81.8	26.6M
ResNet-50	CNN	80.4	25.0M
MgNet-small	CNN	81.0	26.1M
MgNet	CNN	82.0	39.3M
CMT-XS	CNN+Transformer	81.8	15.2M
MgNet-CMT-XS	CNN + Transformer	82.6	17.9M
MgNet-CMT	CNN + Transformer	83.4	30.1M

Table: ImageNet results of transformers and CNNs

Observation:

MgNet has competitive performance with transformer models.

Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

Space decomposition and subspace correction

- V : Hilbert space, $A: V \rightarrow V'$: linear operator, $f \in V'$. Find $u \in V$ such that

$$Au = f.$$

- Space decomposition: $V = \sum_i V_i = \sum_i I_i V_i$:

$$u = \sum_{i=1}^J u_i = \sum_{i=1}^J I_i u_i.$$

- Subspace solvers: $R_i: V'_i \mapsto V_i$ with

$$R_i \approx A_i^{-1}, \quad (A_i u_i, v_i) = (Au_i, v_i), \quad u_i, v_i \in V_i$$

- Parallel subspace correction:

$$u \leftarrow u + B(f - Au), \quad B = \sum_{i=1}^J I_i R_i I_i^T.$$

- Successive subspace correction (SSC): $u \leftarrow u + I_i R_i I_i^T (f - Au)$, for $i = 1 : J$

Xu, J. (1992).

Space decomposition and expanded system

- Space decomposition: $V = \sum_i V_i = \sum_i I_i V_i$:

$$u = \sum_{i=1}^J u_i = \sum_{i=1}^J I_i u_i = \Pi \underline{u}$$

where

$$\Pi = (I_1, \dots, I_J), \quad \underline{u} = (u_1, \dots, u_J)^T$$

- For linear system: $Au = f$, we write

$$u = \Pi \underline{u}, \quad \underline{u} = (u_1, \dots, u_J)^T,$$

and

$$Au = f \Leftrightarrow A\Pi \underline{u} = f \Leftrightarrow \Pi' A \Pi \underline{u} = \Pi' f.$$

- Expanded system

$$\underline{A}\underline{u} = \underline{f}, \quad \text{where} \quad \underline{A} = \Pi' A \Pi, \quad \underline{f} = \Pi' f.$$

Iterative methods for expanded system

Theorem

Iterative methods for $\tilde{A}\tilde{u} = \tilde{f}$:

$$\tilde{u}^m = \tilde{u}^{m-1} + \tilde{B}(\tilde{f} - \tilde{A}\tilde{u}^{m-1}), \quad m = 1, 2, \dots$$

is equivalent to the iterative method for $Au = f$:

$$u^m = u^{m-1} + B(f - Au^{m-1}), \quad m = 1, 2, \dots$$

with

$$B = \Pi \tilde{B} \Pi'.$$

Furthermore,

$$\sigma(BA) = \sigma(\Pi \tilde{B} \Pi' A) = \sigma(\tilde{B} \Pi' A \Pi) \setminus \{0\} = \sigma(\tilde{B} \tilde{A}) \setminus \{0\},$$

and

$$\|I - BA\|_A^2 = 1 - \left(\sup_{\|v\|_A=1} \langle \tilde{B}^{-1} v, v \rangle \right)^{-1} = 1 - \left(\sup_{\|v\|_A=1} \inf_{\Pi \tilde{v} = v} \langle \tilde{B}^{-1} \tilde{v}, \tilde{v} \rangle \right)^{-1}.$$

Auxiliary space lemma

Lemma

Let \underline{V} and V be two vector spaces and let $\Pi : \underline{V} \mapsto V$ be a surjective map. Let $\underline{B} : \underline{V}' \mapsto \underline{V}$ be an SPD operator. Then $B := \Pi \underline{B} \Pi'$ is also SPD. Furthermore

$$\langle B^{-1}v, v \rangle = \inf_{\Pi \underline{v} = v} \langle \underline{B}^{-1}\underline{v}, \underline{v} \rangle.$$

PSC and SSC in the view from expanded system

Theorem

Iterative methods for $\underline{A}\underline{u} = \underline{f}$:

$$\underline{u}^m = \underline{u}^{m-1} + \underline{B}(\underline{f} - \underline{A}\underline{u}^{m-1}), \quad m = 1, 2, \dots$$

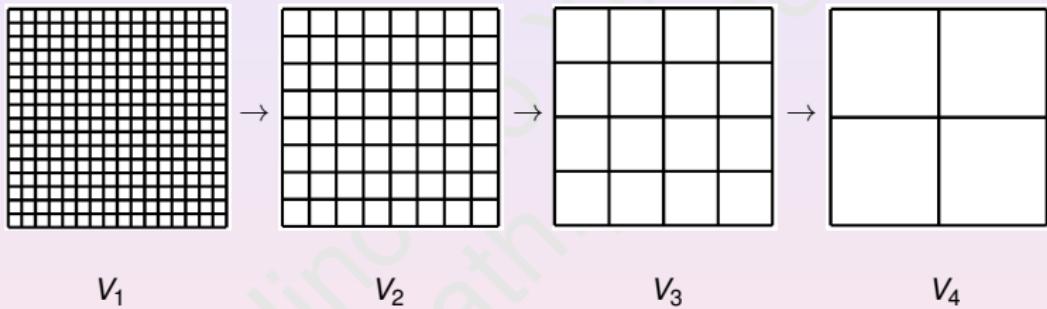
- **PSC** for $Au = f \Leftrightarrow$ modified Jacobi: $\underline{B} = \underline{R} \approx \underline{D}^{-1}$
- **SSC** for $Au = f \Leftrightarrow$ modified G-S: $\underline{B} = (\underline{R}^{-1} + \underline{L})^{-1}$.

Some history:

- X. 1992: DD, MG, Jacobi and GS \Rightarrow PSC or SSC
- Griebel 1994: MG \Leftrightarrow GS for expanded matrix in terms of multilevel nodal basis
- L. Chen 2011: PSC (SSC) \Leftrightarrow Jacobi and GS for expanded matrix (as stated above)

Examples

- Jacobi and block Jacobi methods are parallel subspace corection methods.
- Gauss–Seidel and block Gauss–Seidel methods are successive subspace correction methods.
- Multigrid methods:



$$V = \sum_{k=1}^J V_k = \sum_{k=1}^J \sum_{x_{k,i} \in N_k} \text{span}\{\phi_{k,i}\}$$

- ▶ Successive subspace correction → multigrid with Gauss–Seidel smoothers
- ▶ Parallel subspace correction → BPX preconditioner

Bramble, J.H., Pasciak, J.E., and Xu, J. (1990).

Theory: XZ-identity

Sharp convergence theory for subspace correction methods

$$u - u^n = \prod_{i=1}^J (I - T_i)(u - u^{n-1}), \quad T_i = R_i A_i P_i.$$

Theorem (Xu and Zikatanov (2002, J. AMS, 2008))

The MSC is convergent if each subspace solver is convergent:

$$\left\| \prod_{i=1}^J (I - T_i) \right\|^2 = 1 - \frac{1}{K}, \quad K = \sup_{\|v\|=1} \inf_{\sum_i v_i = v} \sum_{i=1}^J \|v_i + T_i^* \sum_{j=i+1}^J v_j\|_{R_i^{-1}}^2$$

Special case ($T_i = P_i$)

$$\left\| \prod_{i=1}^J (I - P_i) \right\|^2 = 1 - \left(\sup_{\|v\|=1} \inf_{\sum_i v_i = v} \sum_{i=1}^J \|P_i \sum_{j=i}^J v_j\|^2 \right)^{-1}$$

Convergence theory of multigrid methods

Using the XZ identity, we can obtain a uniform convergence rate of the multigrid method.

Corollary (Uniform convergence of multigrid)

The convergence rate of the multigrid method for the finite element method

$$a(u, v) = f(v), \quad \forall v \in V_h$$

has a bound independent of the mesh size h .

Convex optimization

- V : Banach space, $L: V \rightarrow \overline{\mathbb{R}}$: convex function. Find $u \in V$ such that

$$\min_{u \in V} L(u).$$

- In many applications in machine learning, L is of the form

$$L(u) = \frac{1}{N} \sum_{i=1}^N f_i(u).$$

- Gradient descent type methods
 - ▶ Full (batch) gradient descent

$$u_{t+1} = u_t - \eta_t \nabla \left(\frac{1}{N} \sum_{i=1}^N f_i(u_t) \right).$$

- ▶ Stochastic gradient descent (SGD)

$$u_{t+1} = u_t - \eta_t \nabla f_{i_t}(u_t),$$

where $\Pr(i_t = k) = \frac{1}{N}$.

Subspace correction methods for convex optimization

- V : Banach space, $L: V \rightarrow \overline{\mathbb{R}}$: convex function. Find $u \in V$ such that

$$\min_{u \in V} L(u).$$

- Space decomposition $V = \sum_{i=1}^J V_i$, $u = \sum_{i=1}^J u_i$
- Local corrections in subspaces: Find $w_i \in V_i$ such that

$$\min_{w_i \in V_i} L(u + w_i)$$

- Successive subspace correction (SSC):

$$u \leftarrow u + w_i, \text{ for } i = 1 : J$$

- Parallel subspace correction (PSC):

$$u \leftarrow u + \tau \sum_{i=1}^J w_i$$

Convergence theory

- L is M -smooth, i.e.,

$$L(u) \leq L(v) + \langle L'(v), u - v \rangle + \frac{M}{2} \|u - v\|^2, \quad \forall u, v \in V.$$

- L is μ -strongly convex, i.e.,

$$L(u) \geq L(v) + \langle L'(v), u - v \rangle + \frac{\mu}{2} \|u - v\|^2, \quad \forall u, v \in V.$$

Theorem (Tai and Xu (2002), Park (2020))

The MSC for convex optimization is convergent. Moreover, we have

$$\frac{L(u^n) - L(u)}{L(u^{n-1}) - L(u)} \leq 1 - \frac{1}{K},$$

where

$$K \approx \mu^{-1} \sup_{\|w\|=1} \inf_{w=\sum_{i=1}^J w_i} \sum_{i=1}^J \|w_i\|^2.$$

An application: Federated learning

We consider the following N -client training model:

$$\min_{\theta \in \Omega} \left\{ L(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta) \right\}$$

- N : number of clients (devices)
- f_i : loss on local data stored on the client i

Conventional training (GD)

$$\theta \leftarrow \theta - \eta \nabla L(\theta)$$

Federated learning (FL)

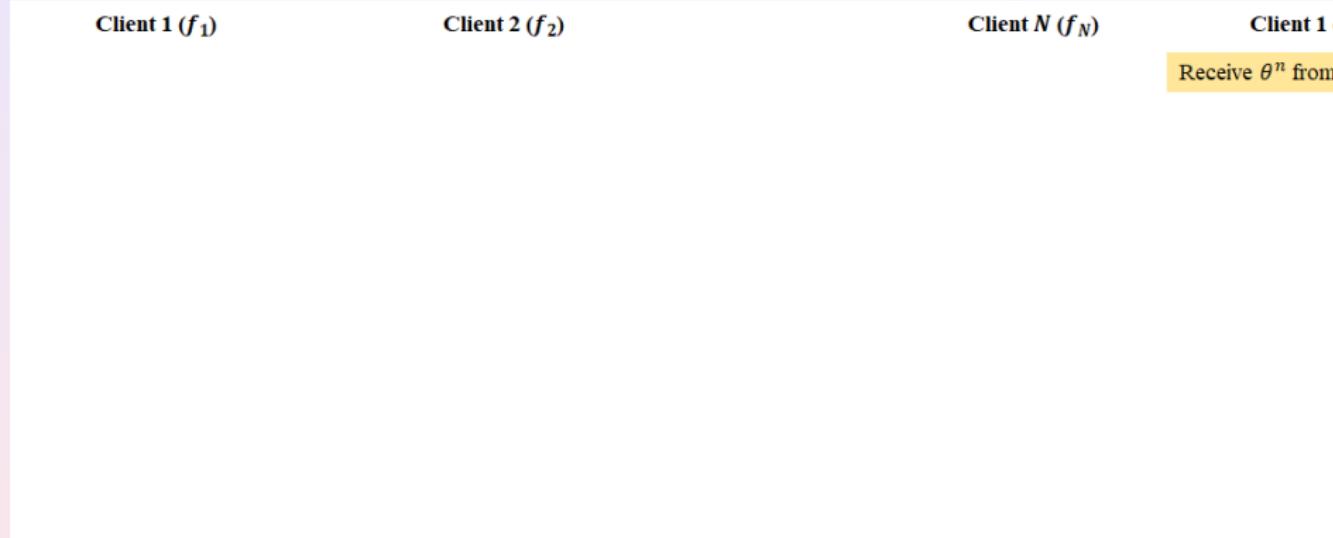
Each client performs local training (several GD steps) using its local function f_i , and the results are averaged in the server.

FedAvg: McMahan, B., Moore, E., Ramage, D., Hampson, S., and Arcas, B.A.y. (2017),
Scaffold: Karimireddy, S.P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A.T. (2020),
Scaffoldnew: Mishchenko, K., Malinovsky, G., Stich, S., and Richtarik, P. (2022),

DualFL: Park, J. and Xu, J. (2023).

An application: Federated learning

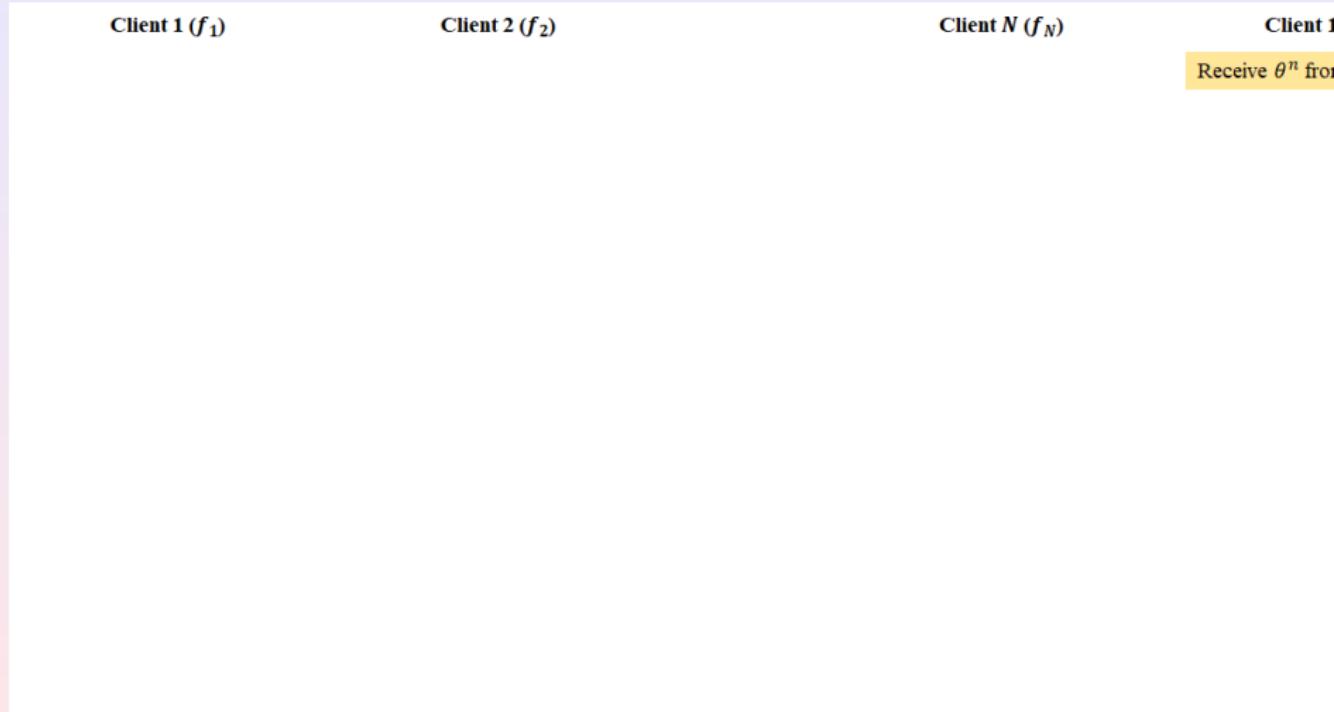
Conventional training (GD)



Communication at every gradient descent step \Rightarrow Too frequent communication!

An application: Federated learning

Federated learning (FL)



Question: By modifying local trainings and global communications, can we design federated learning algorithms with fewer communication costs?

Federated Learning \leftrightarrow Parallel Subspace Correction

Federated learning problem

$$\min_{\theta \in \Omega} \left\{ L(\theta) := \frac{1}{N} \sum_{i=1}^N f_i(\theta) \right\}$$

$$\begin{array}{c} \uparrow \\ \text{Fenchel-Rockafellar duality} \\ \downarrow \\ \theta = -\frac{1}{N\nu} \sum_{i=1}^N \xi_i, \quad \xi_i = \nabla g_i(\theta) \end{array}$$

Dual problem:

$$\min_{\xi \in \Omega^N} \left\{ L_d(\xi) := \sum_{i=1}^N g_i^*(\xi_i) + \frac{1}{2N\nu} \left\| \sum_{i=1}^N \xi_i \right\|^2 \right\}.$$

- $\nu \in (0, \mu]$, $g_i = f_i - \frac{\mu}{2} \|\cdot\|^2$
- $g_i^*: \Omega \rightarrow \overline{\mathbb{R}}$: convex conjugate of g defined by

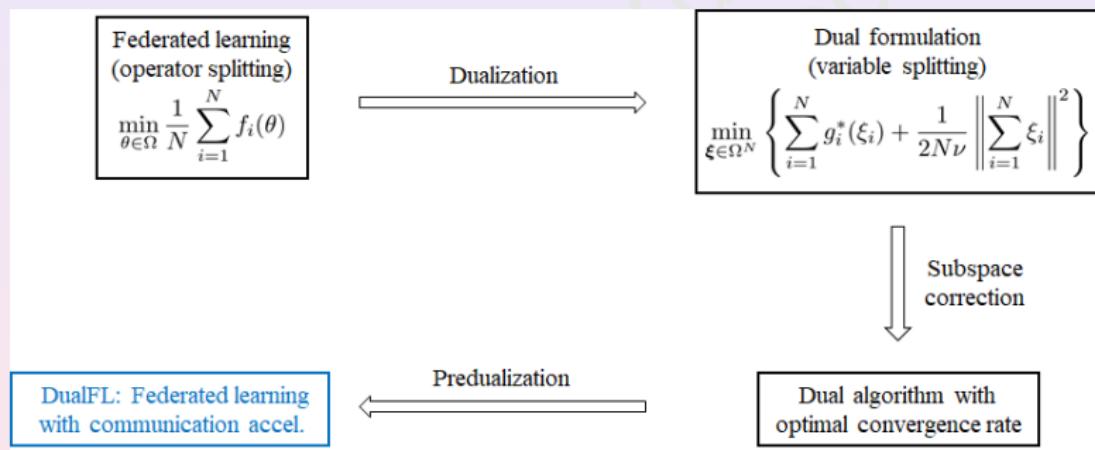
$$g_i^*(\phi) = \sup_{\theta \in \Omega} \{ \langle \phi, \theta \rangle - g_i(\theta) \}$$

$$\blacktriangleright g(x) = \frac{1}{2} x^T A x \Rightarrow g^*(y) = \frac{1}{2} y^T A^{-1} y$$

$$\blacktriangleright g(x) = e^x \Rightarrow g^*(y) = \begin{cases} y \log y - y, & \text{if } y > 0, \\ 0, & \text{if } y = 0, \\ \infty, & \text{if } y < 0. \end{cases}$$

Federated Learning \leftrightarrow Parallel Subspace Correction

By establishing a duality relation between *federated learning* and *parallel subspace correction methods*, we design a new federated learning algorithm with *optimal communication complexity*.



- Park, J., & Xu, J. (2023).

DualFL: Dualized Federated Learning

DualFL: Dualization of a PSC for the dual problem with a space decomposition

$$\Omega^N = \Omega \times \Omega \times \dots \Omega$$

Given $\rho \geq 0$ and $\nu > 0$, set $\theta^0 = \theta_j^0 = 0 \in \Omega$ ($1 \leq j \leq N$) and $\zeta^0 = \zeta^{-1} = \mathbf{0} \in \Omega^N$.

for $n = 0, 1, 2, \dots$ **do**

for each client ($1 \leq j \leq N$) **in parallel do**

$$\theta_j^{n+1} \approx \arg \min_{\theta_j \in \Omega} \left\{ L^{n,j}(\theta_j) := f_j(\theta_j) - \nu \langle \zeta_j^n, \theta_j \rangle \right\}$$

end for

$$\theta^{n+1} = \frac{1}{N} \sum_{j=1}^N \theta_j^{n+1}$$

for each client ($1 \leq j \leq N$) **in parallel do**

 Determine the local shift ζ_j^{n+1} by a linear combination of $\zeta_j^n + \theta^{n+1} - \theta_j^{n+1}$ and
 $\zeta_j^{n-1} + \theta^n - \theta_j^n$.

end for

end for

- Lee, C.-O. & Park, J. (2019), Park, J., & Xu, J. (2023)

DualFL: Dualized Federated Learning

DualFL

Client 1 (f_1)

Client 2 (f_2)

Client N (f_N)

Rec

Local shifts ζ^n obtained by dualization \Rightarrow fast and global convergence to the solution!

Communication efficiency

- Communication efficiency of conventional training (GD)

$$M = \begin{cases} \mathcal{O}\left(\frac{L}{\mu} \log \frac{1}{\epsilon}\right), & \text{if each } f_i \text{ is } \mu\text{-strongly convex and } L\text{-smooth,} \\ \mathcal{O}\left(\frac{1}{\epsilon}\right), & \text{if each } f_i \text{ is convex and } L\text{-smooth.} \end{cases}$$

- DualFL achieves communication acceleration!

Theorem (J. Park and J. Xu, 2023)

In DualFL, the number of communication rounds M to obtain an ϵ -accurate solution satisfies

$$M = \begin{cases} \mathcal{O}\left(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}\right), & \text{if each } f_i \text{ is } \mu\text{-strongly convex and } L\text{-smooth,} \\ \mathcal{O}\left(\frac{1}{\sqrt{\epsilon}}\right), & \text{if each } f_i \text{ is } \mu\text{-strongly convex,} \\ \mathcal{O}\left(\frac{1}{\sqrt{\epsilon}} \log \frac{1}{\epsilon}\right), & \text{if each } f_i \text{ is convex and } L\text{-smooth.} \end{cases}$$

- Xu, J. (1992) Tai, X.-C & Xu, J. (2002), Park, J., & Xu, J. (2023), Park, J. (2020)

Outline

- 1 Linear systems and basic iterative methods
- 2 Frequency principle
- 3 Multigrid methods
- 4 MgNet for image classification
- 5 Subspace correction and federated learning
- 6 Summary

Concluding remarks

Summary

- Iterative methods for $Au = f$
 - ▶ Richardson iteration, Jacobi method, Gauss–Seidel method
- Iterative methods for $Au = f \Leftrightarrow$ Preconditioned GD for $\min \frac{1}{2} u^T Au - f^T u$
- Frequency bias of GD
 - ▶ ReLU neural networks: Training loss decreases slowly in L^2 fitting
 - ▶ Hat neural networks: Training loss decreases fast in L^2 fitting (no frequency bias)
- Frequency bias implies non-convergence of SGD-type algorithms for NN-based PDE solvers.
- Expanded system \Leftrightarrow Over-parametrization \Leftrightarrow Multigrid methods: Remedy for GD
- Subspace correction methods: General framework for iterative methods
- An application of subspace correction methods: Federated learning

Thank You !