

A SUPERVISED LEARNING APPROACH TO PREDICTING MULTIGRID CONVERGENCE

NICOLAS NYTKO*

In collaboration with: Matthew West, Luke Olson, Scott MacLachlan

Abstract. Classical AMG solvers often require careful parameter tuning to achieve optimal convergence, and the way these parameters affect performance can be unpredictable in practice. Evidence is presented that supervised learning techniques are able to learn certain characteristics of two-level multigrid solvers, particularly the rate of convergence and optimal relaxation weight for a given coarse/fine mesh splitting. Random perturbations of C/F splittings are generated and evaluated in a multigrid solver to train a convolutional neural network (CNN) in order to predict convergence and a relaxation weight for the 1D variable coefficient Poisson equation, and to predict the convergence rate for a specific 2D convection-diffusion problem. Additionally for the 2D problem, the use of graph nets is explored for use on general finite-element meshes.

Key words. convergence, relaxation weight, machine learning, supervised learning, graph net, cnn, convolutional network

1. Introduction. Multigrid methods are some of the most widely used linear solvers for sparse systems, particularly ones that arise from the discretization of partial differential equations. By “downsampling” the original problem to a series of smaller and smaller problems, fast convergence can be obtained with minimal work. The question becomes, what is the optimal way to downsample each problem and obtain a splitting of coarse/fine points? Classical Algebraic Multigrid (AMG) techniques attempt to answer this by exploiting the structure of the original system to create a set of interpolation operators. However, these methods often require careful parameter tweaking to obtain sensible splittings and often the best way to test the parameter space is to actually run the multigrid iterations.

In this paper, a method of training neural networks to predict convergence rates for a specific PDE or class of equations is introduced. The neural networks take as input a C/F splitting of a mesh, the output of an AMG aggregation/splitting method, and give a rate of convergence between 0 and 1, inclusive. Lower values of convergence indicate faster methods – more of the error is dissipated at each iteration. Two problems are discussed here: (1) a 1D variable coefficient Poisson equation, and (2) a 2D recirculating flow convection-diffusion problem with specific parameters. Sections 2.1 and 2.3 describe how random grid splittings are generated for the two problems and evaluated in a multigrid solver.

For the 1D poisson case, convolutional neural networks (CNNs) are trained to predict the convergence rate and optimal relaxation weight (Section 2.2). For the 2D convection-diffusion problem, a CNN was first trained to predict convergence (Section 2.4), and then two different graph nets were trained to extrapolate to differently sized inputs and non-grid meshes (Sections 2.5, 2.6). Numerical results consisting of the network error performance and sample predictions are shown in Section 3. Finally, further discussion and possible directions of future work are given in Section 4.

2. Methods. The methods described here will be split into two subsections describing each problem that was explored. First, the Poisson case and the data

*University of Illinois at Urbana-Champaign, (nnytko2@illinois.edu).

generation is discussed (Section 2.1) followed by the respective convolutional neural network that was trained (Section 2.2). Afterwards, an overview of the data generation for the 2D convection-diffusion problem is given (Section 2.3) and then followed by a description of the convolutional network that was trained (Section 2.4). The use of two graph nets is also described: one that performs a simple edge convolution for each node (Section 2.5), and another that implements “message passing” and learns optimal edge weighting of the convolution (Section 2.6).

2.1. 1D Poisson. Formally, the problem being solved is the Poisson equation in one dimension with variable coefficients

$$(2.1) \quad -\nabla \cdot (k(\mathbf{x}) \nabla \mathbf{u}) = f,$$

with the right-hand side $f(\mathbf{x}) = \mathbf{x}$ being arbitrarily chosen. Eqn. (2.1) is discretized using finite differences on a grid of $N = 31$ internal points on the domain $\Omega = [-1, 1]$ and Dirichlet boundary conditions, $\partial\Omega = 0$. To preserve the symmetric, positive definite properties of the resulting linear system, the $k(\mathbf{x})$ function is discretized on the grid *midpoints*[1].

Before the neural network can be trained, a dataset of approximately 300,000 random C/F grid splittings and their computed convergence rate and optimal relaxation weight was generated. A set of 6 reference grids were first created according to various “coarsening” factors:

$$(2.2) \quad r = \{2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9\}.$$

Each of the values in 2.2, say $r_i = j$ refers to a grid in which each j th point is coarse, and the rest fine. So, a coarsening by 3 is a grid that has roughly $\frac{1}{3}$ of its distribution of points as coarse and all others fine. Each of these reference grids was randomly permuted such that each grid point had a random probability of being flipped to the opposite value – coarse point to fine and fine point to coarse. For each reference grid, approximately 1000 random trials of each following probability were run:

$$(2.3) \quad p = \{0.01 \ 0.05 \ 0.1 \ 0.25 \ 0.5 \ 0.75\}$$

Additionally, each trial was run with a randomly-chosen function for the variable coefficients $k(\mathbf{x})$. The function is randomly chosen from

$$(2.4) \quad k(\mathbf{x}) = \begin{cases} \alpha & 0 < \alpha < 10 \\ \text{rand}() (\alpha + 1) & 0 < \alpha < 10 \\ \alpha \cos(\pi x \beta) + \gamma & 0 < \alpha < 10, 0 < \beta < 10, \alpha < \gamma < 10 \\ \left| \sum_{i=1}^5 \alpha_i x^i \right| + 0.01 & -10 < \alpha < 10 \end{cases},$$

with coefficient values specifically picked to prevent the function from being non-positive at any value in Ω .

The full set of C/F splittings was then given to a 2 level V-cycle Multigrid solver that was run for a maximum of 15 iterations. The solver is composed of one round of weighted Jacobi pre-smoothing, a coarse error correction, then another round of Jacobi post-smoothing. At each iteration, the absolute error between the approximation and the “exact” solution (pre-computed via sparse linear solve) was found and saved: $e_i = \|\mathbf{A}^{-1} \mathbf{f} - \mathbf{u}_i\|$. This sequence of errors was then used to compute the average convergence rate. To acquire the optimal relaxation weight, the smoother was run

through a bracketed numerical optimization method with the assumption that the convergence rate is unimodal as a function of relaxation weight.¹

2.2. Poisson CNN. Two separate deep convolutional networks with residual connections were trained to separately predict the optimal relaxation weight and convergence rate when given a linear C/F splitting. The architecture for both networks is identical and a textual overview is given here:

1. 6 1D CNN layers of kernel size 7, input 2 channels output 7 channels. Padding by three on each side to keep dimensions static.
2. 6 CNN layers of kernel size 5, input/output 7 channels. Padding by two on each side.
3. 6 CNN layers of kernel size 3, input/output 7 channels. Padding by 1 on each side.
4. Max-pooling layer of kernel 2, stride 2. Effectively reduces input size by half.
5. 8 Fully-connected layers to gradually reduce output to a scalar describing the convergence or relaxation weight.

Each layer is followed by an implicit ReLU nonlinear activation function. Because the convolutional layers (except for the first) keep the input and output size static, we are able to push residual values and skip layers similarly to a ResNet. Even numbered layers n take as input both the output of layers $n - 1$ and $n - 2$, while odd-numbered layers only use the output from layer $n - 1$.

The C/F splitting is remapped such that a coarse point is given a value of 1 and a fine point the value of -1 . The variable coefficients were also re-discretized to be defined on the nodal points instead of midpoints in order for the splitting and coefficients to be represented by vectors of same length. These two were then stacked into a two-channel tensor for input into the CNN. When training, input values (convergence rate and relaxation weight) are normalized to be within the range of $[0, 1]$. This normalization is undone when the output is displayed.

2.3. 2D Convection-Diffusion. This specific problem models what is called the *double glazing problem*, modeling the temperature distribution of a cavity with a single “hot” wall. This is given by the differential equation:

$$(2.5) \quad -k\nabla^2 u + \mathbf{w} \cdot \nabla u = f.$$

The wind velocity function, $\mathbf{w}(x, y)$ is defined as $\mathbf{w}(x, y) = [2y(1 - x^2) \quad 2x(1 - y^2)]$. The domain is a square of side length two centered at the origin, $\Omega = [-1, 1] \times [-1, 1]$. Dirichlet boundary conditions are defined on $\partial\Omega$, with the one “hot” wall defined on $x = 1$ with value $\partial\Omega_H = 1$. The other boundaries are “cold” walls with $\partial\Omega_c = 0$. A diffusivity constant of $k = 0.1$ is used. This problem is derived from an example by Elman, Silvester, and Wathen[2].

This PDE is discretized using finite-elements on a structured grid of 25x25 internal points. Using a grid as a basis for the discretization allows use of both CNN and more sophisticated graph convolutional techniques.

Generating a dataset of mesh splittings and convergence rates was done in an overall similar way to the 1D Poisson equation with a few notable differences. Again, a set of reference C/F splittings were generated that are later permuted. For the convection-diffusion case, the following reference splittings were used:

1. All fine points

¹Experimental testing generally asserts this to be true, however this will remain a conjecture for now.

2. All coarse points
3. Splitting as given by Ruge-Stüben coarsening ($\theta = 0.25$)
4. Coarsening in each direction by 2
5. Coarsening in each direction by 3
6. Coarsening in each direction by 4
7. Coarsening in each direction by 5

The entries of each individual reference splitting were then randomly permuted according to a defined probability. The probability values used are the same as those in the Poisson case, repeated here for convenience:

$$(2.6) \quad p = \{0.01 \quad 0.05 \quad 0.1 \quad 0.25 \quad 0.5 \quad 0.75\}.$$

Generated splittings that are unsolvable (i.e., consist of no coarse points) were rejected and their trial re-run. This generated set of C/F splittings was passed along to another 2 level V-cycle multigrid solver run for a maximum of 50 iterations. This solver performs two rounds of Jacobi pre-and-post smoothing, with a coarse error correction between the smoothing steps. The absolute error between the approximation at each iteration and the “exact” solution was computed and saved, with the full sequence used to find the average convergence rate. Note the optimal Jacobi relaxation weight was not found here, as experimentation found that the optimal weight would nearly always be 1.

The convection-diffusion problem was additionally re-discretized at 4 different mesh sizes: 15×15 , 25×25 , 35×35 , 50×50 . The above process was re-run for each mesh size to generate a new dataset for use in the graph nets. To distinguish between the two datasets, they will hereby be referred to as the *statically-sized* (only containing 25×25 mesh and splittings) and the *variably-sized* datasets.

2.4. Convection-Diffusion CNN. A 2D convolutional network was trained on the statically-sized dataset to predict convergence when given a C/F splitting for the specific recirculating flow problem. Since the input parameters are slightly less complex, a less deep (*shallower*, if you will) network was trained:

1. 3 2D CNN layers of kernel 7, input two channels output 7 channels. Padding by three on each side to keep dimensions.
2. 3 CNN layers of kernel 5, input/output 7 channels. Padding by two on each side.
3. 3 CNN layers of kernel 3, input/output 7 channels. Padding by 1 on each side.
4. 2D Max-pooling layer of kernel 2, stride 2. Effectively reduces input size by half.
5. 1 Fully-connected layer to reduce output to a scalar predicting convergence rate.

Each layer is followed by a ReLU nonlinear activation function. Odd layers are fed the output of the previous two layers, while even layers are fed the input of only the previous layer, similarly to a ResNet. C/F splittings are mapped so that coarse points have value 1 and fine points have value -1 . The CNN thus has a $25 \times 25 = 625$ -length vector as input. Interestingly, normalization of the convergence rates is unneeded as the minimum and maximum recorded values are already close to 0 and 1, respectively.

2.5. Convection-Diffusion Graph Convolutional Network (GCN). The main downside of using traditional convolutional layers is that they are useful only on structured, grid-like inputs. Grid-based convolution is ineffective on the more complex

meshes that may arise from finite-element discretizations. By treating the input mesh as a graph and using graph-based convolution techniques, we may get around this roadblock. The first graphnet that was tried uses the GCN layer introduced by Kipf and Welling[4] and is described here, the other uses an Edge-Conditioned Convolution (ECC)/message passing layer and is described in Section 2.6.

For the graphnets, the sparse FEM system was re-interpreted as a graph by assigning each row/column to be a node and defining connectivity between nodes as nonzero matrix entries. Edge weights were taken to be the entry values themselves, although normalized to be within the range of $[0, 1]$. This network was trained on the variably-sized dataset; as will be seen, the lack of any fully-connected layers allows for any arbitrary input.

The Graph Convolutional Network (GCN) layer is defined as an operator on the graph Laplacian, using the following propagation rule:

$$(2.7) \quad \mathbf{H}^{(i)} = \sigma \left(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(i-1)} \mathbf{W}^{(i)} + \mathbf{b}^{(i)} \right).$$

With $\mathbf{H}^{(i)}$ being the i 'th hidden layer, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ the adjacency matrix with added self loops, $\tilde{d}_{ii} = \sum_j \tilde{a}_{ij}$ a diagonal matrix consisting of the sum of outgoing edge weights, $\mathbf{W}^{(i)}$ the weight matrix at layer i , $\mathbf{b}^{(i)}$ a learned bias vector with same shape as $\mathbf{H}^{(i)}$, and $\sigma(\cdot)$ some nonlinear activation function. Hidden layers have dimensions $n \times f$, with n the number of graph nodes and f the number of features at each node. Thus, using only GCN layers it is not possible to reduce the number of rows in hidden layers (or the input layer for that matter).

Assume $\sigma(\cdot) = \text{ReLU}(\cdot) = \max\{\cdot, 0\}$ for all activation functions. Let \mathbf{X} be the $n \times 1$ vector containing C/F splitting values for each node. The network that was trained for this application consisted of the following architecture:

$$\begin{aligned} \mathbf{H}^{(1)} &= \sigma \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X} \mathbf{W}^{(1)} + \mathbf{b}^{(1)} \right) \in \mathbb{R}^{n \times 2} \\ \mathbf{H}^{(2)} &= \sigma \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(1)} \mathbf{W}^{(2)} + \mathbf{b}^{(2)} \right) \in \mathbb{R}^{n \times 3} \\ \mathbf{H}^{(3)} &= \sigma \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(2)} \mathbf{W}^{(3)} + \mathbf{b}^{(3)} \right) \in \mathbb{R}^{n \times 2} \\ \mathbf{H}^{(4)} &= \sigma \left(\mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(3)} \mathbf{W}^{(4)} + \mathbf{b}^{(4)} \right) \in \mathbb{R}^{n \times 1} \end{aligned}$$

$$\mathbf{R} = \begin{bmatrix} \dots & \mathbf{X}^T & \dots \\ \dots & \left(\mathbf{H}^{(1)} \right)^T & \dots \\ \dots & \left(\mathbf{H}^{(2)} \right)^T & \dots \\ \dots & \left(\mathbf{H}^{(3)} \right)^T & \dots \\ \dots & \left(\mathbf{H}^{(4)} \right)^T & \dots \end{bmatrix} \in \mathbb{R}^{9 \times n}$$

$$y = \sum_{j=1}^n \sigma \left(\sigma \left(\mathbf{r}_j \mathbf{W}^{(5)} + \mathbf{b}^{(5)} \right) \mathbf{W}^{(6)} + \mathbf{b}^{(6)} \right) \in \mathbb{R}$$

Note that the matrix \mathbf{R} is formed whose columns are the propagation history of each node; this is an attempt to emulate traditional ResNet architectures. The final scalar output is computed by summing each column of \mathbf{R} through a two layer fully-connected neural network, with the first layer taking as input 9 features and

Model	Metric	Dataset	Value
Jacobi Weight	MSE	Training	1.8331×10^{-3}
Jacobi Weight	MSE	Testing	1.8396×10^{-3}
Jacobi Weight	L1	Training	2.9254×10^{-2}
Jacobi Weight	L1	Testing	2.9271×10^{-2}
Convergence Factor	MSE	Training	1.4839×10^{-3}
Convergence Factor	MSE	Testing	1.5171×10^{-3}
Convergence Factor	L1	Training	2.3908×10^{-2}
Convergence Factor	L1	Testing	2.3865×10^{-2}

Table 3.1: Final training and testing Mean Squared Error (MSE)/ L^1 -norm loss values for the two Poisson models. Lower values correspond to higher model accuracy.

outputting 5 features, and the second layer taking as input 5 features and outputting 1 feature. This local transformation of each nodal value following by a global averaging removes any dependency for a fixed input size – instead any sized graph and splitting could (in theory) be used.

2.6. Convection-Diffusion Message Passing Network (MPNN). A downside of the GCN layer, while being simple to implement and understand, is that it does not effectively learn edge features between nodes and thus is only an approximation of a convolution on a graph. The Edge-Conditioned Convolution (ECC) layer, as defined by Simonovsky and Komodakis[5], attempts to combat this by learning optimal edge weights given an arbitrary node-edge neighborhood. This approach is also referred to as a “message-passing network”, as each connected neighbor attempts to learn a “message” that is passed to the original node.[3] The convolution operation is defined as

$$(2.8) \quad \left(\mathbf{H}^{(i)} \right)_j = \frac{1}{|\mathcal{N}(j)|} \sum_{k \in \mathcal{N}(j)} F^{(i)}(e_{j,k}) \left(\mathbf{H}^{(i-1)} \right)_k + \mathbf{b}^{(i)}$$

where $\mathbf{H}^{(i)} \in \mathbb{R}^{n \times f_i}$ is the i ’th hidden layer, $\mathcal{N}(i)$ is a map returning the neighborhood of vertex i (including itself), $F^{(i)} : E \mapsto \mathbb{R}^{f_i \times f_{i-1}}$ is some function (in our case a small fully-connected neural network) that outputs a learned weight matrix given an edge, and $\mathbf{b}^{(i)}$ is a bias term for layer i . The notation $\left(\mathbf{H}^{(i)} \right)_j$ here denotes the j ’th row of the i ’th hidden layer.

Using the ECC layer, another graph net was trained to predict the convergence rate given some C/F mesh splitting. The architecture is identical to that in Section 2.5 with the exception that the propagation layers $\mathbf{H}^{(i)} : i = 1 \dots 4$ are replaced with (2.8) and surrounded with a ReLU activation layer. Each $F^{(i)}$ is a two-layer fully-connected network taking some edge $\mathbf{e} \in \mathbb{R}^{f_i}$ and outputting a weight matrix $\Theta \in \mathbb{R}^{f_i \times f_{i-1}}$.

3. Numerical Results.

3.1. 1D Poisson. A total of 336,000 C/F “grid” splittings were randomly generated for the variable-coefficient Poisson model. Of these, a random 85% of data items were designated as a training set and the remaining 15% reserved for a testing set. The relaxation and convergence CNNs were trained for a total of 30 epochs each,

where each epoch is a total pass over all training set entries. Models were trained using an MSE loss, and final loss values for both networks are detailed in Table 3.1.

Sample convergence rate predictions on the training and testing subsets are shown in Figure 3.1; these plots display predicted rates as a function of the true rates, so error is indicated by deviation from the diagonal of the plots. It is interesting to visually see that the network has an easier time predicting more convergent grids (≈ 0.6 or so) than more degenerate grids. A lack of predictions less than ≈ 0.1 is a deficiency of the input data not having grids with those convergence values.

Predictions for relaxation weights on both training and testing subsets are shown in Figure 3.2. Like the previous convergence plots, more accurate predictions are closer to the plot diagonals. There is a positive correlation between the predictions and the true values, though unlike the convergence rates there is no obvious “split” in the prediction data. The lack of predictions below 0.1 is likely to also be caused by the input data.

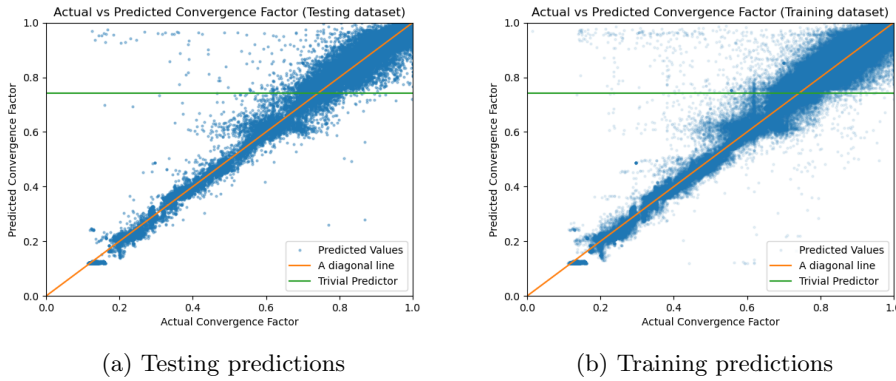


Fig. 3.1: Predicted convergence values vs. true convergence values on (3.1a) testing and (3.1b) training datasets for Poisson equation. Values closer to the diagonal represent more accurate predictions.

3.2. 2D Convection-Diffusion. Two datasets were created for the 2D recirculating flow problem: a set of *statically-sized* C/F splittings and a set of *variably-sized* splittings; both datasets consisted of 84,000 elements. The variably-sized dataset was evenly divided into the four mesh sizes: 15×15 , 25×25 , 35×35 , 50×50 . Each dataset was then split 85%-15% into training and testing subsets. All neural networks were trained for 10 epochs each on an MSE loss, final loss values are detailed in Table 3.2. The recirculating flow CNN was trained and tested on the static dataset, while the two graph nets were trained on the variable dataset and evaluated on both datasets. The CNN had the best performance for the static dataset, though the GCN network was very close despite it being trained on a completely different set of data. However, the ECC (MPNN) network had better predictive performance than the GCN network over the entire variable dataset.

Predicted convergence values by the CNN over the static dataset are displayed in Figure 3.3. The subplots detail predicted convergence rates vs their true value for each C/F splitting. Almost opposite to the phenomenon observed in the Poisson CNN

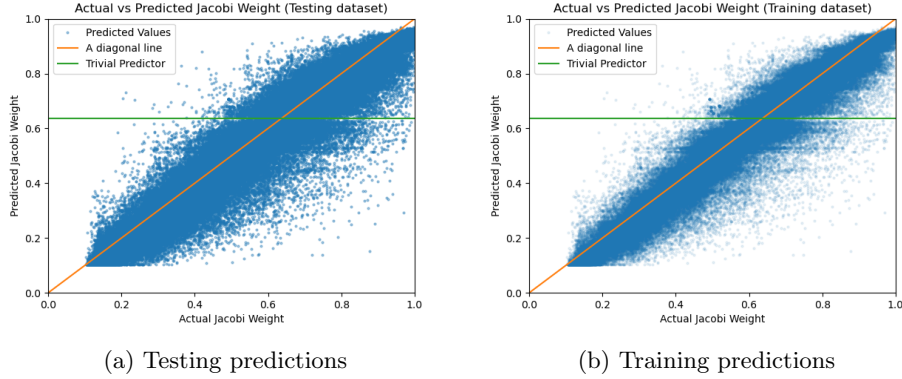


Fig. 3.2: Predicted relaxation weights vs. true relaxation weights on (3.2a) testing and (3.2b) training datasets for Poisson equation. Values closer to the diagonal represent more accurate predictions.

Dataset	Model	Metric	Value
Static	CNN	MSE	6.0332×10^{-3}
Static	CNN	L1	5.2335×10^{-2}
Static	GCN	MSE	6.7278×10^{-3}
Static	GCN	L1	5.5773×10^{-2}
Static	MPNN	MSE	8.4461×10^{-3}
Static	MPNN	L1	7.0079×10^{-2}
Variable	GCN	MSE	1.3079×10^{-2}
Variable	GCN	L1	8.4827×10^{-2}
Variable	MPNN	MSE	1.1495×10^{-2}
Variable	MPNN	L1	7.9967×10^{-2}

Table 3.2: Final testing Mean Squared Error (MSE)/ L^1 -norm loss values for trained neural networks for the recirculating flow problem. Lower values correspond to higher model accuracy.

(Figure 3.1, the network predicts more poorly converging grids accurately. This could perhaps be explained by the data more densely populated with the poor convergence grids, and has a distinct lack of optimal grids.

Convergence predictions by the ECC network over the variable-size dataset are given in Figure 3.4. Various horizontal trends of predictions are visible in the plots, with distinct subsets of convergence factors all being predicted as the same value. Further tweaking of the network architecture is likely necessary to achieve more accuracy in the predictions.

4. Discussion and Conclusions. In this paper, a method for training convolutional and graph nets was presented to predict convergence of a two-level multigrid solver when given an input C/F mesh splitting. Random mesh splittings were generated from a set of reference grids for two model problems, then evaluated in a solver to find a reference convergence rate. Additionally for the Poisson problem, a CNN to

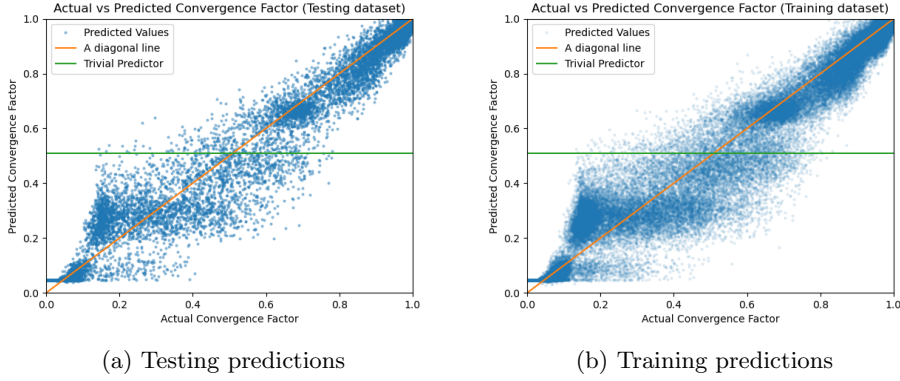


Fig. 3.3: Predicted convergence rates vs. true convergence rates on (3.3a) testing and (3.3b) training datasets for model convection-diffusion problem, using a convolutional network. Values closer to the diagonal represent more accurate predictions.

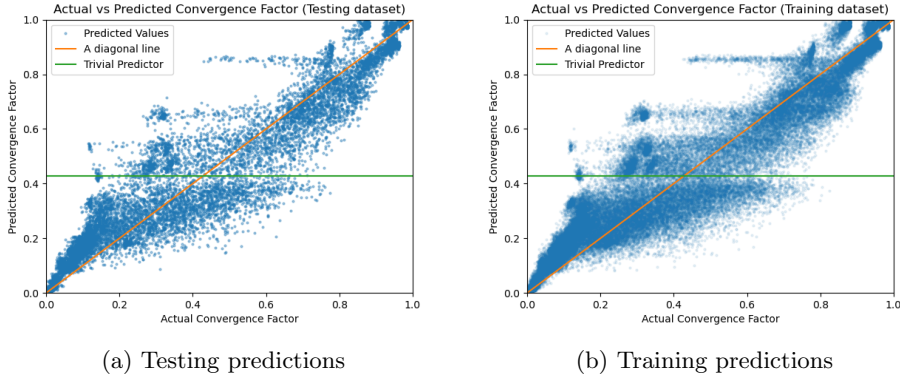


Fig. 3.4: Predicted convergence rates vs. true convergence rates on (3.4a) testing and (3.4b) training datasets for model convection-diffusion problem, using an Edge-Conditioned Convolution network. Values closer to the diagonal represent more accurate predictions.

predict optimal smoothing weight was trained.

Results indicate that neural networks can indeed predict convergence for a multigrid solver. Convolutional networks are effective at learning attributes on structured grid-like 1D or 2D meshes, and graph nets are able to generalize on variably-sized structured meshes. The graph networks could be further developed and trained on more general classes of problems, instead of the specific convection-diffusion problem or the variable-coefficient Poisson equation.

Future works could build on the networks trained here and explore their applications, using them for example to pick between different AMG aggregation methods to find the most optimal for a specific use case. Another more ambitious future work

could be to use the neural networks in an optimization method to obtain the most convergent or most work efficient C/F splitting for a problem.

REFERENCES

- [1] J. ADLER, H. D. STERCK, S. MACLACHLAN, AND L. OLSON, *Numerical Partial Differential Equations*, 2020.
- [2] H. ELMAN, D. SILVESTER, AND A. WATHEN, *Finite Elements and Fast Iterative Solvers With Applications in Incompressible Fluid Dynamics*, Oxford University Press, 2014.
- [3] J. GILMER, S. S. SCHOENHOLZ, P. F. RILEY, O. VINYALS, AND G. E. DAHL, *Neural message passing for quantum chemistry*, 2017, <https://arxiv.org/abs/1704.01212>.
- [4] T. N. KIPF AND M. WELLING, *Semi-supervised classification with graph convolutional networks*, 2017, <https://arxiv.org/abs/1609.02907>.
- [5] M. SIMONOVSKY AND N. KOMODAKIS, *Dynamic edge-conditioned filters in convolutional neural networks on graphs*, 2017, <https://arxiv.org/abs/1704.02901>.