

Meta-Learning and Universality

Conference paper at ICLR 2018
Chelsea Finn & Sergey Levine (US Berkeley)

Youngseok Yoon

Contents

- The Universality in Meta-Learning
- Model Construct (Pre-update)
- Single gradient step and Post-update
- Show universality from model
- Experiments
- Conclusion

The Universality in Meta-Learning

- What is the universality in Meta-Learning?
- Universality in Neural Network:
Approximate any function $f: x \rightarrow y$
- Universality in Meta-Learning:
Approximate any function $f: ((x, y)_k, x^*) \rightarrow \hat{y}$
Input: The training data $(x, y)_k$, and Test data x^* .
In this presentation, only show about $k = 1$ case.
- Our model to verify Universality:
MAML: $\hat{f}(\cdot; \theta) \rightarrow \hat{f}(\cdot; \theta')$, $\theta' = \theta - \alpha \cdot \nabla_{\theta} \mathcal{L}(\text{target}, \hat{f}(\cdot; \theta))$
Update model $f(\cdot; \theta)$ using gradient update

Model Construct (Pre-update)

- Our model: Deep neural network with ReLU, MAML training.

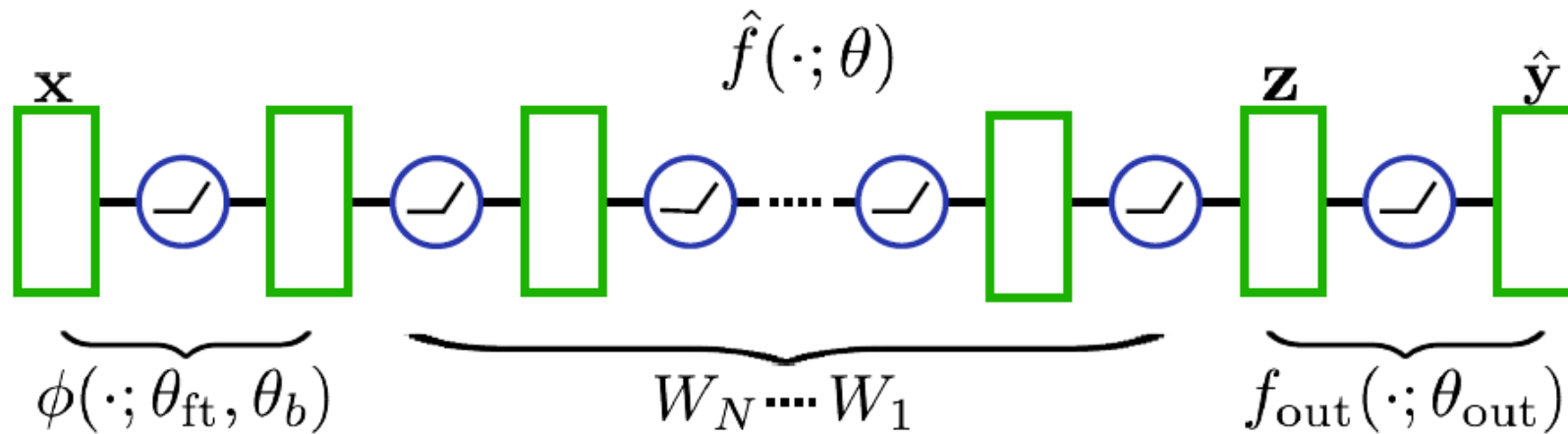
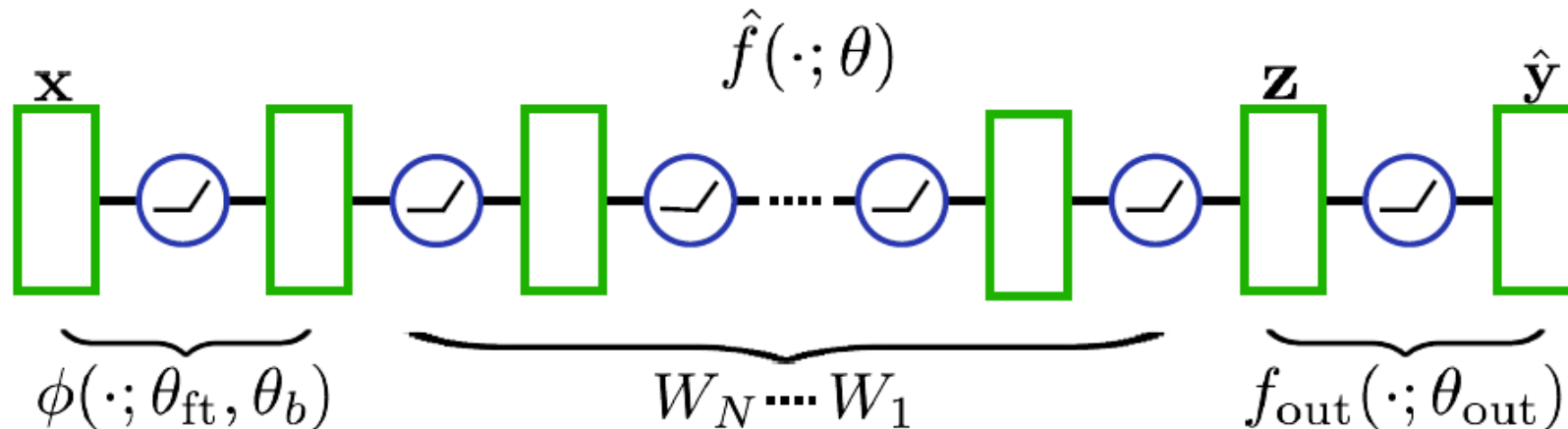


Figure 1: A deep fully-connected neural network with $N+2$ layers and ReLU nonlinearities. With this generic fully connected network, we prove that, with a single step of gradient descent, the model can approximate any function of the dataset and test input.

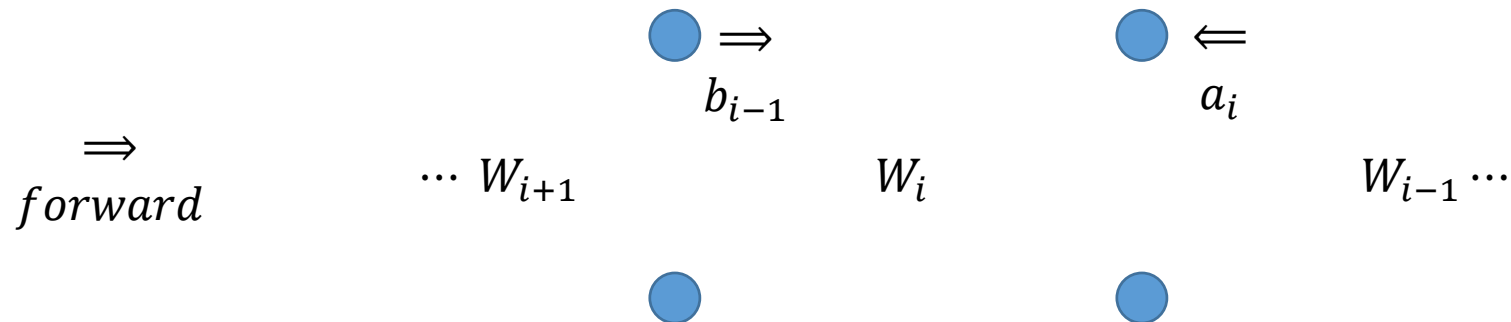
Model Construct (Pre-update)

- Construct model $\hat{f}(\cdot; \theta)$ with ReLU
- If we assume inputs and all pre-synaptic activations are non-negative, then deep ReLU act like deep linear networks. ... (1)
- $\hat{f}(\cdot; \theta) = f_{out} \left(\left(\prod_{i=1}^N W_i \right) \phi(\cdot; \theta_{ft}, \theta_b); \theta_{out} \right)$



Single gradient step and Post-update

- $\nabla_{W_i} \mathcal{L} = a_i b_{i-1}^T$
 - a_i : the error gradient with respect to the pre-synaptic activations at layer i
 - b_{i-1} : the forward post-synaptic activations at layer $i - 1$

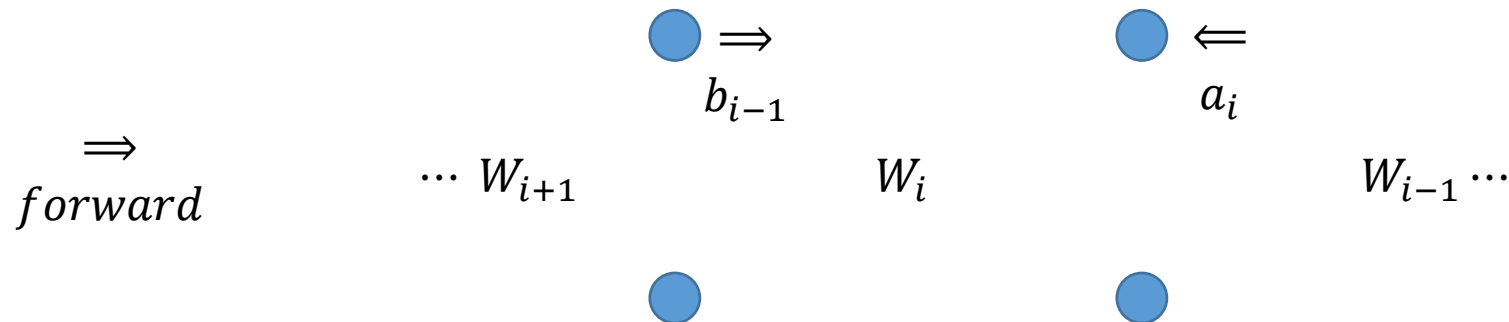


Single gradient step and Post-update

- let $\mathbf{z} = (\prod_{i=1}^N W_i) \phi(\mathbf{x}; \theta_{ft}, \theta_b)$, $\nabla_{\mathbf{z}} \mathcal{L} = e(\mathbf{x}, \mathbf{y})$

- Then,

$$\begin{aligned} \nabla_{W_i} \mathcal{L} &= a_i b_{i-1}^T \\ &= W_{i-1}^T \cdots W_1^T e(\mathbf{x}, \mathbf{y}) \cdot \left((\prod_{j=i+1}^N W_j) \phi(\mathbf{x}; \theta_{ft}, \theta_b) \right)^T \\ &= \left(\prod_{j=1}^{i-1} W_j \right)^T e(\mathbf{x}, \mathbf{y}) \cdot \phi(\mathbf{x}; \theta_{ft}, \theta_b)^T \left(\prod_{j=i+1}^N W_j \right)^T \end{aligned}$$



Single gradient step and Post-update

- let $\mathbf{z} = (\prod_{i=1}^N W_i) \phi(\mathbf{x}; \theta_{ft}, \theta_b)$, $\nabla_{\mathbf{z}} \mathcal{L} = e(\mathbf{x}, \mathbf{y})$

- Then,

$$\begin{aligned} \nabla_{W_i} \mathcal{L} &= a_i b_{i-1}^T \\ &= W_{i-1}^T \cdots W_1^T e(\mathbf{x}, \mathbf{y}) \cdot \left((\prod_{j=i+1}^N W_j) \phi(\mathbf{x}; \theta_{ft}, \theta_b) \right)^T \\ &= (\prod_{j=1}^{i-1} W_j)^T e(\mathbf{x}, \mathbf{y}) \cdot \phi(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N W_j)^T \end{aligned}$$

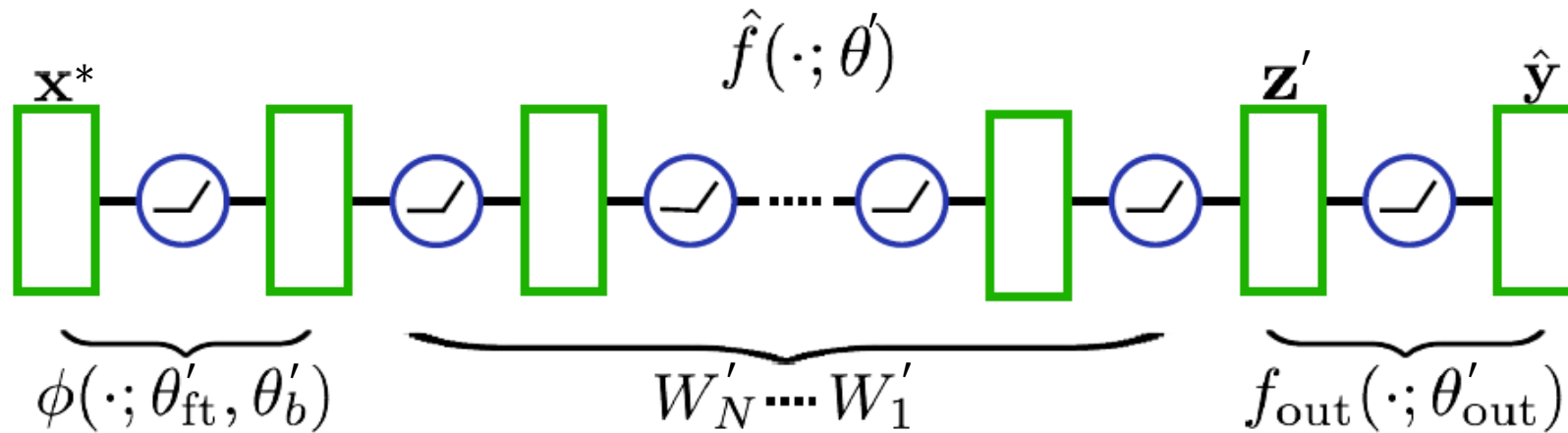
- Therefore, post-update value of $\prod_{i=1}^N W'_i = \prod_{i=1}^N (W_i - \alpha \cdot \nabla_{W_i} \mathcal{L})$
 $\prod_{i=1}^N W_i - \alpha \sum_{i=1}^N (\prod_{j=1}^{i-1} W_j) (\prod_{j=1}^{i-1} W_j)^T e(\mathbf{x}, \mathbf{y}) \cdot \phi(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N W_j)^T (\prod_{j=1}^{i-1} W_j) - O(\alpha^2)$

Single gradient step and Post-update

- Post-update value of $\prod_{i=1}^N W'_i = \prod_{i=1}^N (W_i - \alpha \cdot \nabla_{W_i} \mathcal{L})$
 $\prod_{i=1}^N W_i - \alpha \sum_{i=1}^N (\prod_{j=1}^{i-1} W_j) (\prod_{j=1}^{i-1} W_j)^T e(\mathbf{x}, \mathbf{y}) \cdot \phi(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N W_j)^T (\prod_{j=1}^{i-1} W_j) - O(\alpha^2)$
- Now, post-update value of \mathbf{z}^* when \mathbf{x}^* input into $\hat{f}(\cdot; \theta')$
 $\mathbf{z}^* = \prod_{i=1}^N W'_i \phi(\mathbf{x}^*; \theta'_{ft}, \theta'_b)$
 $= \prod_{i=1}^N W_i \phi(\mathbf{x}^*; \theta'_{ft}, \theta'_b)$
 $- \alpha \sum_{i=1}^N (\prod_{j=1}^{i-1} W_j) (\prod_{j=1}^{i-1} W_j)^T e(\mathbf{x}, \mathbf{y}) \cdot \phi(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N W_j)^T (\prod_{j=1}^{i-1} W_j) \phi(\mathbf{x}^*; \theta'_{ft}, \theta'_b)$
- And, $\hat{f}(\mathbf{x}^*; \theta') = f_{out}(\mathbf{z}^*; \theta'_{out})$

Show universality from model

- How to show universality from post-updated model?

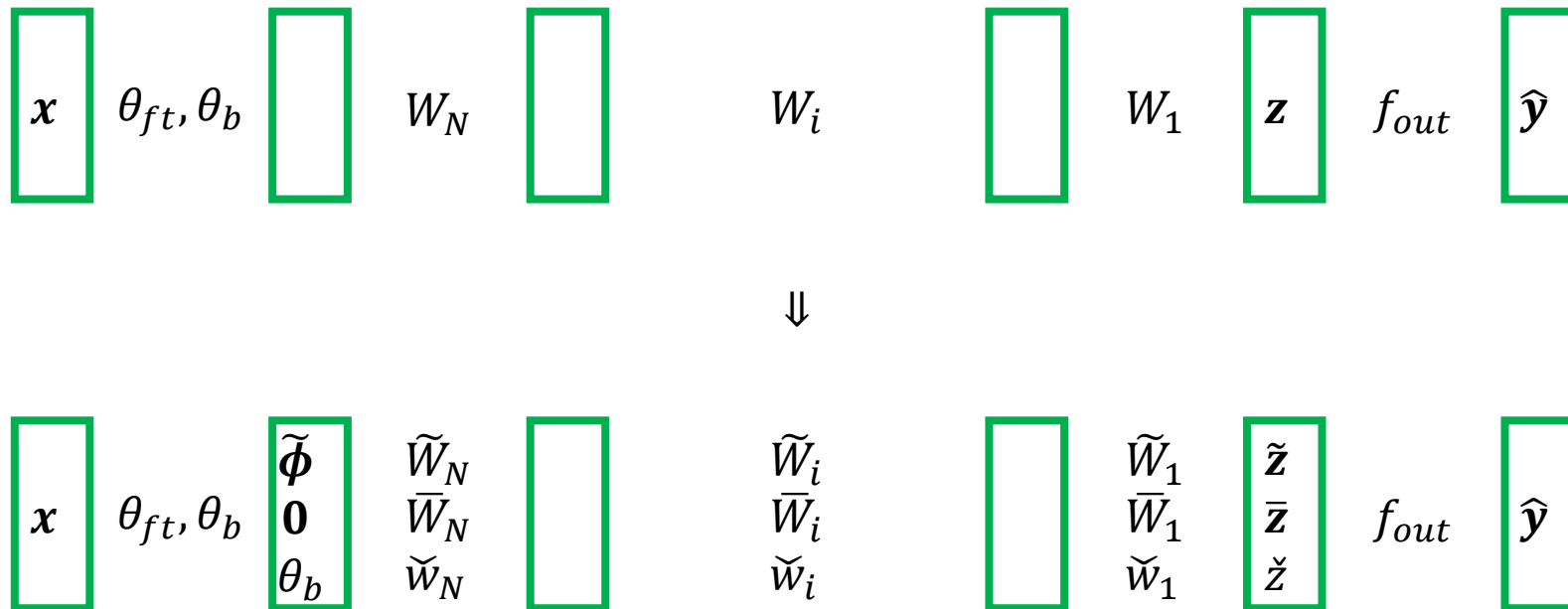


Show universality from model

- Use independently control information flow from \mathbf{x} , from \mathbf{y} , from \mathbf{x}^* by multiplexing forward information from \mathbf{x} and backward information from \mathbf{y}
- Decomposing W_i, ϕ and error gradient into three parts,

$$W_i := \begin{bmatrix} \tilde{W}_i & 0 & 0 \\ 0 & \bar{W}_i & 0 \\ 0 & 0 & \check{w}_i \end{bmatrix}, \quad \phi(\cdot; \theta_{ft}, \theta_b) := \begin{bmatrix} \tilde{\phi}(\cdot; \theta_{ft}, \theta_b) \\ \mathbf{0} \\ \theta_b \end{bmatrix}, \quad \nabla_z \mathcal{L}(y, \hat{f}(\mathbf{x}; \theta)) := \begin{bmatrix} \mathbf{0} \\ \bar{e}(y) \\ \check{e}(y) \end{bmatrix} \dots (2)$$

Show universality from model



Show universality from model

- Then, before the gradient update,

$$\mathbf{z} = \begin{bmatrix} \tilde{\mathbf{z}} \\ \bar{\mathbf{z}} \\ \check{\mathbf{z}} \end{bmatrix} = \begin{bmatrix} \prod_{i=1}^N \tilde{W}_i \tilde{\phi}(x; \theta_{ft}, \theta_b) \\ \mathbf{0} \\ 0 \end{bmatrix}$$

- Hence, represent the middle one of after gradient update, \mathbf{z}^*

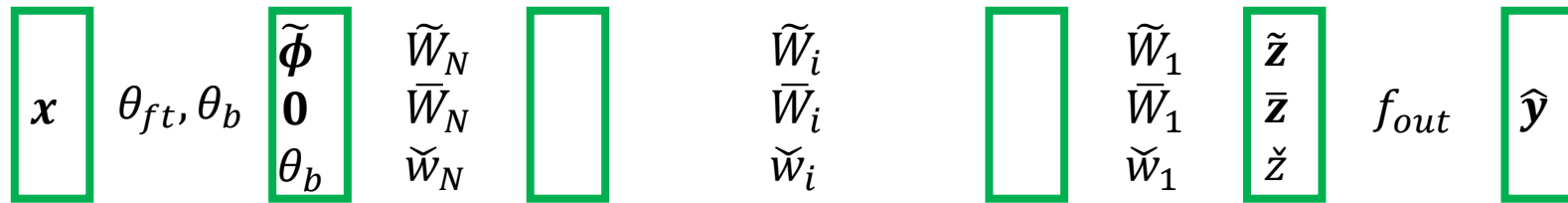
$$\mathbf{z}^* = \prod_{i=1}^N \begin{bmatrix} \tilde{W}_i & 0 & 0 \\ 0 & \bar{W}_i & 0 \\ 0 & 0 & \check{w}_i \end{bmatrix} \begin{bmatrix} \tilde{\phi}(x; \theta_{ft}, \theta_b) \\ \mathbf{0} \\ \theta_b \end{bmatrix} - \alpha \sum_{i=1}^N \left(\prod_{j=1}^{i-1} \begin{bmatrix} \tilde{W}_j & 0 & 0 \\ 0 & \bar{W}_j & 0 \\ 0 & 0 & \check{w}_j \end{bmatrix} \right) \left(\prod_{j=1}^{i-1} \begin{bmatrix} \tilde{W}_j & 0 & 0 \\ 0 & \bar{W}_j & 0 \\ 0 & 0 & \check{w}_j \end{bmatrix} \right)^T \begin{bmatrix} \mathbf{0} \\ \bar{e}(y) \\ \check{e}(y) \end{bmatrix}.$$

$$\begin{bmatrix} \tilde{\phi}(x; \theta_{ft}, \theta_b) \\ \mathbf{0} \\ \theta_b \end{bmatrix}^T \left(\prod_{j=i+1}^N \begin{bmatrix} \tilde{W}_j & 0 & 0 \\ 0 & \bar{W}_j & 0 \\ 0 & 0 & \check{w}_j \end{bmatrix} \right)^T \left(\prod_{j=1}^{i-1} \begin{bmatrix} \tilde{W}_j & 0 & 0 \\ 0 & \bar{W}_j & 0 \\ 0 & 0 & \check{w}_j \end{bmatrix} \right) \begin{bmatrix} \tilde{\phi}(x^*; \theta'_{ft}, \theta'_b) \\ \mathbf{0} \\ \theta'_b \end{bmatrix}$$

$$\therefore \bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N (\prod_{j=1}^{i-1} \bar{W}_j) (\prod_{j=1}^{i-1} \bar{W}_j)^T \bar{e}(y) \cdot \tilde{\phi}(x; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N \tilde{W}_j)^T (\prod_{j=1}^{i-1} \tilde{W}_j) \tilde{\phi}(x^*; \theta'_{ft}, \theta'_b)$$

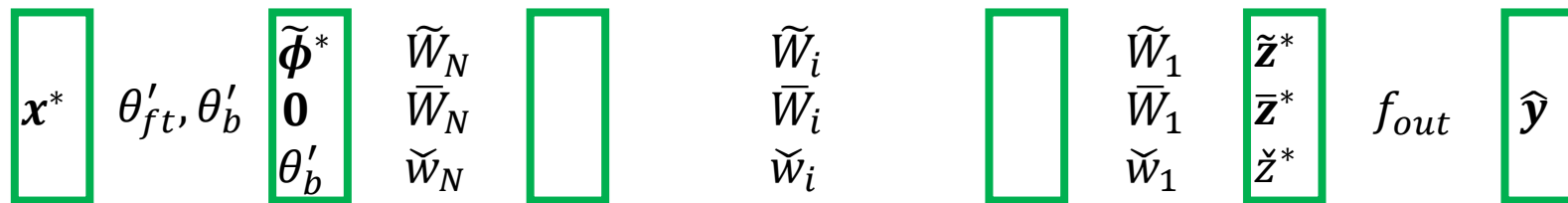
Show universality from model

- $\hat{y} = f_{out}(\mathbf{z}; \theta_{out})$ and $\bar{\mathbf{z}} = \mathbf{0}$



- $$\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N (\prod_{j=1}^{i-1} \bar{W}_j) (\prod_{j=1}^{i-1} \bar{W}_j)^T \bar{e}(\mathbf{y}) \cdot \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N \tilde{W}_j)^T (\prod_{j=1}^{i-1} \tilde{W}_j) \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)$$

$$= -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_i^T B_i \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T \dots (3)$$



Show universality from model

- Path 1)
 - Define f_{out}
- Path 2)
 - Show \bar{z}^* as new vector v , which containing informations of x , y and x^*
- At last,
 - Universality with Vector v

Show universality from model

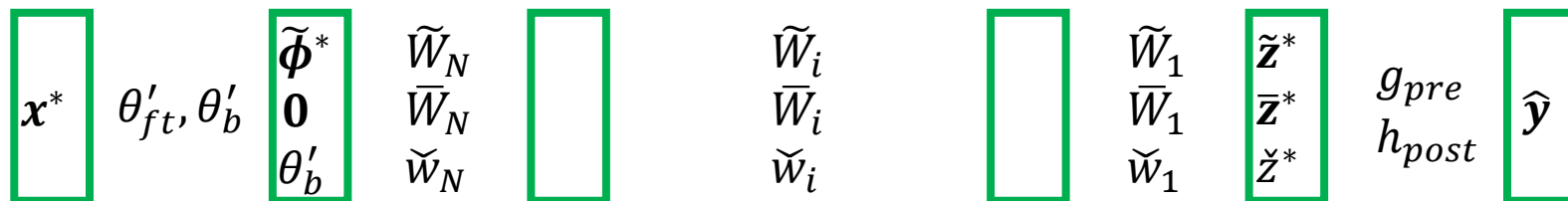
- Define f_{out} as a neural network that approximates the following multiplexer function and its derivatives

- $$f_{out} \left(\begin{bmatrix} \tilde{\mathbf{z}} \\ \bar{\mathbf{z}} \\ \check{\mathbf{z}} \end{bmatrix}; \theta_{out} \right) = \mathbf{1}_{\bar{\mathbf{z}}=0} g_{pre} \left(\begin{bmatrix} \tilde{\mathbf{z}} \\ \bar{\mathbf{z}} \\ \check{\mathbf{z}} \end{bmatrix}; \theta_g \right) + \mathbf{1}_{\bar{\mathbf{z}} \neq 0} h_{post}(\bar{\mathbf{z}}; \theta_h), \dots (4)$$

- Where g_{pre} is a linear function with parameters θ_g such that $\nabla_{\mathbf{z}} \mathcal{L} = e(\mathbf{y}) \dots (2)$
- Where h_{post} is a neural network with one or more hidden layers.

- Then, after pose-update,

- $$f_{out} \left(\begin{bmatrix} \tilde{\mathbf{z}}^* \\ \bar{\mathbf{z}}^* \\ \check{\mathbf{z}}^* \end{bmatrix}; \theta'_{out} \right) = h_{post}(\bar{\mathbf{z}}^*; \theta_h), \dots (5)^*$$



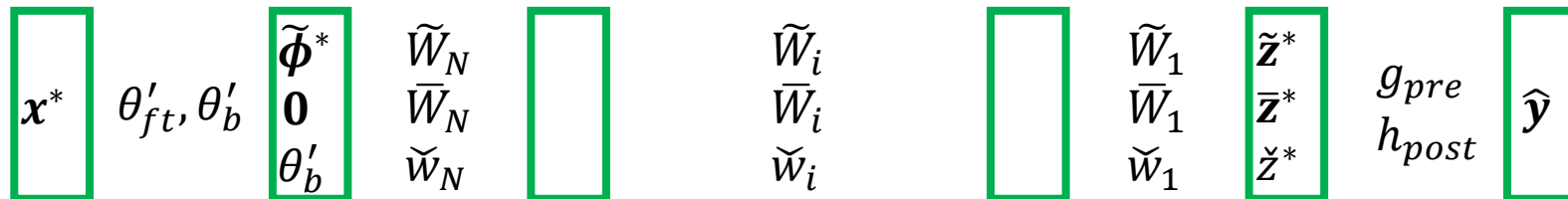
Show universality from model

- Change the form of $\bar{\mathbf{z}}^*$ with kernel form

- $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_i^T B_i \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*)$

- $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*)$ has the all information of \mathbf{x}, \mathbf{y} and \mathbf{x}^* .

- Assume $\bar{e}(\mathbf{y})$ is a linear function of \mathbf{y} which can extract the original information \mathbf{y} .



Show universality from model

- Change the form of $\bar{\mathbf{z}}^*$ with kernel form
 - $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_i^T B_i \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*)$
- $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*)$ has the all information of \mathbf{x}, \mathbf{y} and \mathbf{x}^* .
- Assume $\bar{e}(\mathbf{y})$ is a linear function of \mathbf{y} which can extract the original information \mathbf{y} .
- **Idea 1)**
 - Decompose index $i = 1, \dots, N$ to a product of $j = 0, \dots, J - 1$ and $l = 0, \dots, L - 1$, and then assign those to \mathbf{x} and \mathbf{x}^*
- **Idea 2)**
 - Discretize \mathbf{x} and \mathbf{x}^* to j and l by kernel $k(\cdot, \cdot)$
- **Idea 3)**
 - Make new vector \mathbf{v} , which contain all information of \mathbf{x}, \mathbf{x}^* and \mathbf{y} from $A_i, \bar{e}(\mathbf{y})$ and k_i

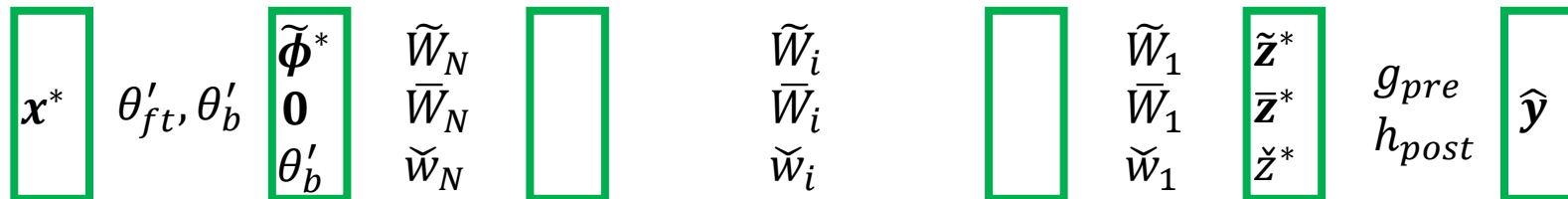
Show universality from model

- **Idea 1)**

- Decompose index $i = 1, \dots, N$ to a product of $j = 0, \dots, J - 1$ and $l = 0, \dots, L - 1$, and then assign those to \mathbf{x} and \mathbf{x}^*
- $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*) = -\alpha \sum_{j=0}^{J-1} \sum_{l=0}^{L-1} A_{jl} \bar{e}(\mathbf{y}) k_{jl}(\mathbf{x}, \mathbf{x}^*)$

- **Idea 2)**

- Discretize \mathbf{x} and \mathbf{x}^* to j and l by kernel $k(\cdot, \cdot)$
- $k_{jk}(\mathbf{x}, \mathbf{x}^*) = \begin{cases} 1 & \text{if } \text{discr}(\mathbf{x}) = e_j \text{ and } \text{discr}(\mathbf{x}^*) = e_l \\ 0 & \text{otherwise} \end{cases} \dots (6)$



Show universality from model

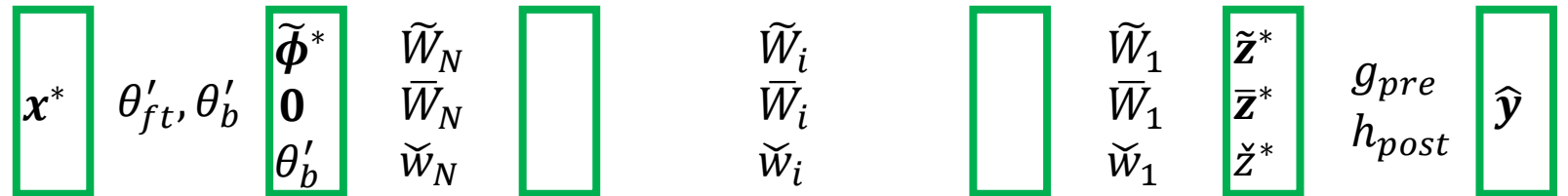
- **Idea 3)**

- Make new vector v , which contain all information of \mathbf{x}, \mathbf{x}^* and \mathbf{y} from $A_i, \bar{e}(\mathbf{y})$ and k_i
- Choose $\bar{e}(\mathbf{y})$ to be a linear function that outputs $J * L$ stacked copies of \mathbf{y}
- Define A_{jl} to be a select \mathbf{y} at the position of $(j, l) = j + J * L, A_{jl} = \begin{bmatrix} \epsilon & \dots & \epsilon \\ \vdots & 1 + \epsilon & \vdots \\ \epsilon & \dots & \epsilon \end{bmatrix}, 1 + \epsilon$ at position (j, l)

- As a result,

- $\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*) = -\alpha \sum_{j=0}^{J-1} \sum_{l=0}^{L-1} A_{jl} \bar{e}(\mathbf{y}) k_{jl}(\mathbf{x}, \mathbf{x}^*)$

$$= v(\mathbf{x}, \mathbf{x}^*, \mathbf{y}) \sim \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \mathbf{y} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$



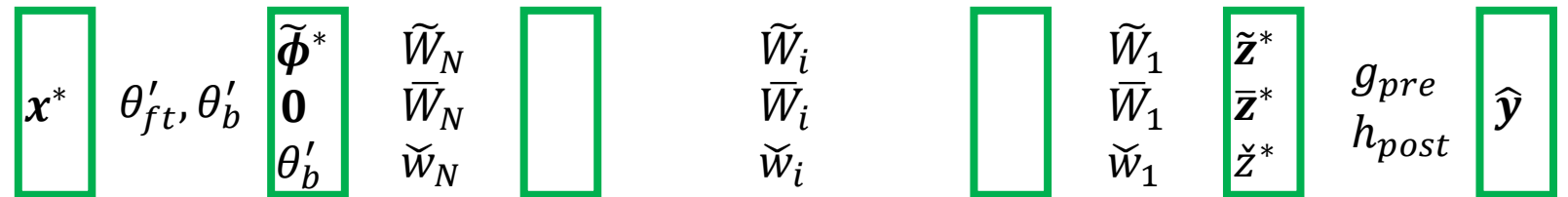
Show universality from model

- As a result,

- $$\bar{\mathbf{z}}^* = -\alpha \sum_{i=1}^N A_i \bar{e}(\mathbf{y}) k_i(\mathbf{x}, \mathbf{x}^*) = -\alpha \sum_{j=0}^{J-1} \sum_{l=0}^{L-1} A_{jl} \bar{e}(\mathbf{y}) k_{jl}(\mathbf{x}, \mathbf{x}^*) = -\alpha v(\mathbf{x}, \mathbf{x}^*, \mathbf{y}),$$

where $v(\mathbf{x}, \mathbf{x}^*, \mathbf{y}) \sim \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \mathbf{y} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

- $\therefore \hat{f}(\mathbf{x}^*; \theta') \sim h_{post}(-\alpha v(\mathbf{x}, \mathbf{x}^*, \mathbf{y}); \theta_h)$
- It is a universal function approximator with respect to its inputs $(\mathbf{x}, \mathbf{x}^*, \mathbf{y})$



Experiments

- Q1) Is there empirical benefit to using one meta-learning approach versus another, and in which case?
- A1) Empirically show the inductive bias of gradient-based vs recurrent meta-learners
 - Explore the differences between gradient-based vs recurrent
 - A learner trained with MAML, improve or start to overfit after additional gradient steps.
 - Better few-shot learning performance on tasks outside of the training distribution?

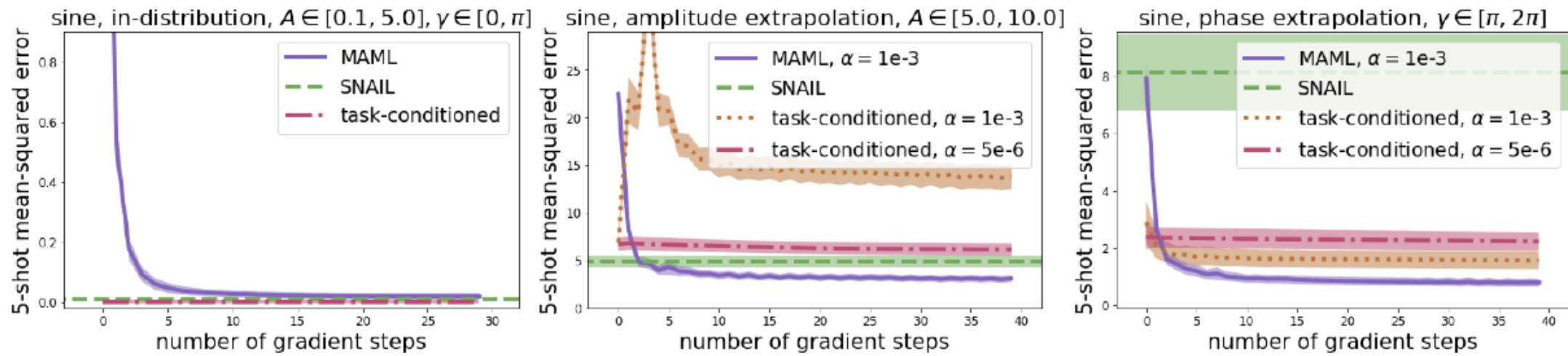


Figure 2: The effect of additional gradient steps at test time when attempting to solve new tasks. The MAML model, trained with 5 inner gradient steps, can further improve with more steps. All methods are provided with the same data – 5 examples – where each gradient step is computed using the same 5 datapoints.

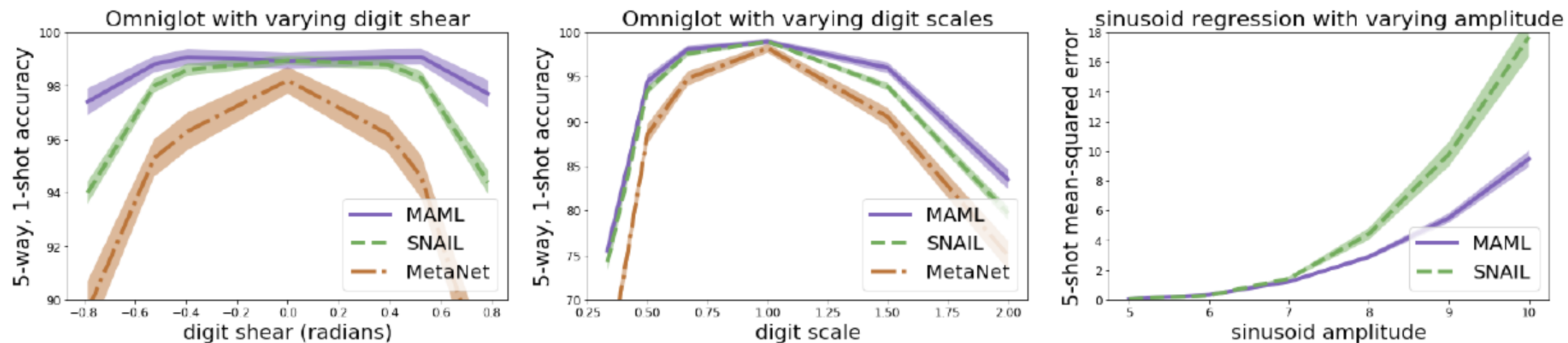


Figure 3: Learning performance on out-of-distribution tasks as a function of the task variability. Recurrent meta-learners such as SNAIL and MetaNet acquire learning strategies that are less generalizable than those learned with gradient-based meta-learning.

Experiments

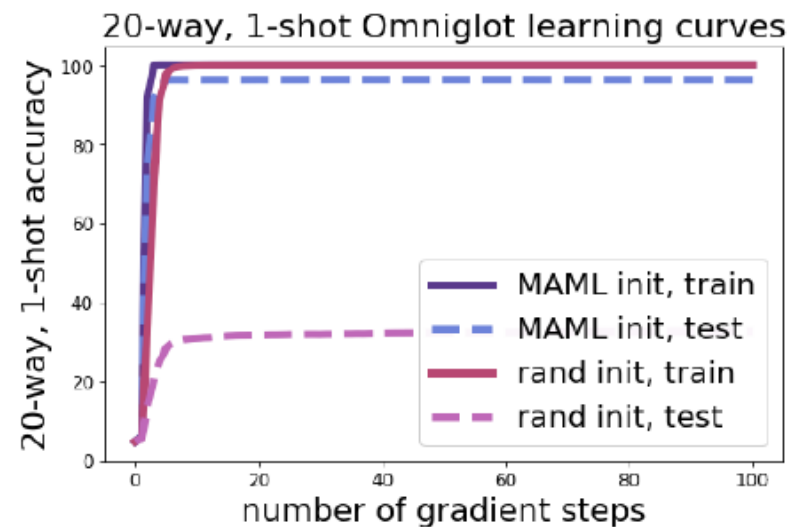


Figure 4: Comparison of finetuning from a MAML-initialized network and a network initialized randomly, trained from scratch. Both methods achieve about the same training accuracy. But, MAML also attains good test accuracy, while the network trained from scratch overfits catastrophically to the 20 examples. Interestingly, the MAML-initialized model does not begin to overfit, even though meta-training used 5 steps while the graph shows up to 100.

Experiments

- Q2) Theory suggest that deeper networks lead to increased expressive power for representing different learning procedures. Is it right?
- A1) Investigate the role of model depth in gradient-based meta-learning

Experiments

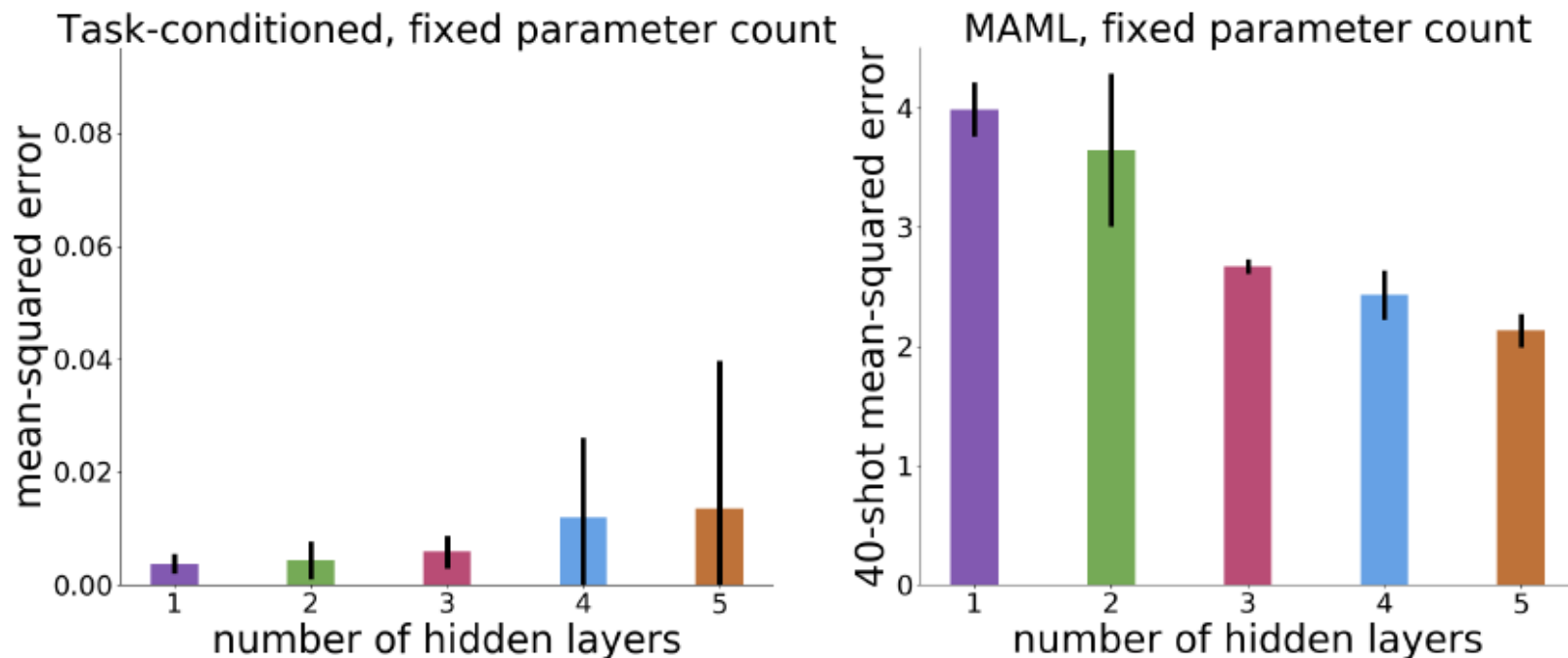


Figure 5: Comparison of depth while keeping the number of parameters constant. Task-conditioned models do not need more than one hidden layer, whereas meta-learning with MAML clearly benefits from additional depth. Error bars show standard deviation over three training runs.

Conclusion

- Meta-learners that use standard gradient descent with a sufficiently deep representation can approximate any learning procedure, and are equally expressive as recurrent learners.
- In experiments, MAML is more successful when faced with out-of-domain tasks compared to recurrent model
- We formalize what it means for a meta-learner to be able to approximate any learning algorithm in terms of its ability to represent functions of the dataset and test inputs.

Appendix

- A

- we assume inputs and all pre-synaptic activations are non-negative, then deep ReLU act like deep linear networks. ... (1)

- $$\begin{aligned}\bar{\mathbf{z}}^* &= -\alpha \Sigma_{i=1}^N (\prod_{j=1}^{i-1} \bar{W}_j) (\prod_{j=1}^{i-1} \bar{W}_j)^T \bar{e}(\mathbf{y}) \cdot \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T (\prod_{j=i+1}^N \tilde{W}_j)^T (\prod_{j=1}^{i-1} \tilde{W}_j) \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b) \\ &= -\alpha \Sigma_{i=1}^N A_i \bar{e}(\mathbf{y}) \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_i^T B_i \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T \quad \dots (3)\end{aligned}$$

- B

- $$\nabla_{\mathbf{z}} \mathcal{L}(y, \hat{f}(\mathbf{x}; \theta)) := \begin{bmatrix} \mathbf{0} \\ \bar{e}(\mathbf{y}) \\ \check{e}(\mathbf{y}) \end{bmatrix} \dots (2)$$

- $$f_{out} \left(\begin{bmatrix} \tilde{\mathbf{z}} \\ \bar{\mathbf{z}} \\ \check{\mathbf{z}} \end{bmatrix}; \theta_{out} \right) = \mathbf{1}_{\bar{\mathbf{z}}=0} g_{pre} \left(\begin{bmatrix} \tilde{\mathbf{z}} \\ \bar{\mathbf{z}} \\ \check{\mathbf{z}} \end{bmatrix}; \theta_g \right) + \mathbf{1}_{\bar{\mathbf{z}} \neq 0} h_{post}(\bar{\mathbf{z}}; \theta_h), \quad \dots (4), \quad f_{out} \left(\begin{bmatrix} \tilde{\mathbf{z}}^* \\ \bar{\mathbf{z}}^* \\ \check{\mathbf{z}}^* \end{bmatrix}; \theta'_{out} \right) = h_{post}(\bar{\mathbf{z}}^*; \theta_h), \quad \dots (5)^*$$

- C

- $$k_{jk}(\mathbf{x}, \mathbf{x}^*) = \begin{cases} 1 & \text{if } \text{discr}(\mathbf{x}) = e_j \text{ and } \text{discr}(\mathbf{x}^*) = e_l \\ 0 & \text{otherwise} \end{cases} \quad \dots (6)$$

Appendix A

$$\begin{aligned}\bar{\mathbf{z}}^* &= -\alpha \Sigma_{i=1}^N \left(\prod_{j=1}^{i-1} \bar{W}_j \right) \left(\prod_{j=1}^{i-1} \bar{W}_j \right)^T \bar{e}(\mathbf{y}) \cdot \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T \left(\prod_{j=i+1}^N \tilde{W}_j \right)^T \left(\prod_{j=1}^{i-1} \tilde{W}_j \right) \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b) \\ &= -\alpha \Sigma_{i=1}^N A_i \bar{e}(\mathbf{y}) \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_i^T B_i \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T\end{aligned}$$

- Choose all \tilde{W}_i and \bar{W}_i to be a square and full rank.
- Set $\tilde{W}_i = \tilde{M}_i \tilde{M}_{i+1}^{-1}$ and $\bar{W}_i = \bar{M}_{i-1}^{-1} \bar{M}_i$ then $\prod_{j=i+1}^N \tilde{W}_j = \tilde{M}_{i+1}$, $\prod_{j=1}^{i-1} \bar{W}_j = \bar{M}_{i-1}$ ($\tilde{M}_{N+1} = \mathbf{I}$, $\bar{M}_0 = \mathbf{I}$)
- Set $A_1 = \mathbf{I}$, $A_i = \bar{M}_{i-1} \bar{M}_{i-1}^T$ and $B_i = \tilde{M}_{i+1}$, $B_N = \mathbf{I}$

Appendix A

- All inputs are non-negative:
 - Input ϕ consist of three term, (discretization, constant 0, 0 before update and not used afterward)
 - All inputs are non-negative before and after update.
- All pre-synaptic activations are non-negative:
 - It is sufficient to show that products of W_i , $\prod_{i=j}^N W_i$ ($1 \leq j \leq N$) are positive semi-definite.
 - So, show that $\prod_{i=j}^N \tilde{W}_i$, $\prod_{i=j}^N \bar{W}_i$, $\prod_{i=j}^N \check{w}_i$ ($1 \leq j \leq N$) are positive semi-definite.
 - We define $\prod_{i=j}^N \tilde{W}_i = \tilde{M}_{j+1} = B_i$, and B_i is set to be a positive definite.
 - We define $\bar{W}_i = \bar{M}_{i-1}^{-1} \bar{M}_i$ where $A_i = \bar{M}_{i-1} \bar{M}_{i-1}^T$, so each \bar{M}_i is also symmetric positive definite and \bar{W}_i is positive definite.
 - Purpose of \check{w}_i is to provide nonzero gradient to the input θ_b , thus positive value for each \check{w}_i will suffice.

Appendix B

$$\nabla_{\mathbf{z}} \mathcal{L}(y, \hat{f}(\mathbf{x}; \theta)) := \begin{bmatrix} \mathbf{0} \\ \bar{e}(y) \\ \check{e}(y) \end{bmatrix} \dots (2) \quad \nabla_{\mathbf{z}} \mathcal{L}(y, \hat{f}(\mathbf{x}; \theta)) = \nabla_{\mathbf{z}} \hat{f}(\mathbf{x}; \theta) \cdot \nabla_{\hat{y}} \mathcal{L}(y, \hat{y}) = \begin{bmatrix} \mathbf{0} \\ \bar{e}(y) \\ \check{e}(y) \end{bmatrix}$$
$$\hat{f}(\mathbf{x}; \theta) = \hat{f}_{out}(\mathbf{z}; \theta_{out}) = g_{pre}(\mathbf{z}; \theta_g)$$

- Because we assume g_{pre} as a linear function of θ_g ,
let $g_{pre}(\mathbf{z}; \theta_g) = [\tilde{W}_g \quad \bar{W}_g \quad \check{w}_g] \mathbf{z} = \tilde{W}_g \tilde{\mathbf{z}} + \bar{W}_g \bar{\mathbf{z}}_g + \check{w}_g \check{z}$
- To make top element of $e(\mathbf{x}, \mathbf{y})$ to be $\mathbf{0}$, let $\tilde{W}_g = 0$, which make $\hat{y} = 0$.
- Then, $\bar{e}(y) = \bar{W}_g^T \cdot \nabla_{\hat{y}} \mathcal{L}(y, 0)$
- For any linear loss function for y , $\bar{e}(y) = Ay$,
to extract all information of y , A has to be invertible.
- Sufficient loss function: standard mean-squared error or softmax cross entropy.

Appendix C

$$k_{jk}(\mathbf{x}, \mathbf{x}^*) = \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_{jl}^T B_{jl} \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T = \begin{cases} 1 & \text{if } \text{discr}(\mathbf{x}) = e_j \text{ and } \text{discr}(\mathbf{x}^*) = e_l \\ 0 & \text{otherwise} \end{cases}$$

- Choose $\tilde{\phi}$ and B_{jl} :

$$\tilde{\phi}(\cdot; \theta_{ft}, \theta_b) := \begin{cases} \begin{bmatrix} [\text{discr}(\cdot)] \\ 0 \end{bmatrix} & \text{if } \theta_b = 0 \\ \begin{bmatrix} 0 \\ [\text{discr}(\cdot)] \end{bmatrix} & \text{otherwise} \end{cases} \quad B_{jl} = \begin{bmatrix} E_{jj} & E_{jl} \\ E_{lj} & 0 \end{bmatrix} + \epsilon I$$

where E_{ik} denote the matrix a 1 at (i, k) and 0 otherwise

- Then,

$$\tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_{jl}^T \sim \begin{cases} [e_j & 0] & \text{if } \text{discr}(\mathbf{x}) = e_j \\ [0 & 0] & \text{otherwise} \end{cases} \quad B_{jl} \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T \sim \begin{cases} [e_l & 0]^T & \text{if } \text{discr}(\mathbf{x}^*) = e_l \\ [0 & 0]^T & \text{otherwise} \end{cases}$$

Appendix C

$$k_{jk}(\mathbf{x}, \mathbf{x}^*) = \tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_{jl}^T B_{jl} \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T = \begin{cases} 1 & \text{if } \text{discr}(\mathbf{x}) = e_j \text{ and } \text{discr}(\mathbf{x}^*) = e_l \\ 0 & \text{otherwise} \end{cases}$$

$$\tilde{\phi}(\cdot; \theta_{ft}, \theta_b) := \begin{cases} \begin{bmatrix} \text{discr}(\cdot) \\ 0 \end{bmatrix} & \text{if } \theta_b = 0 \\ \begin{bmatrix} 0 \\ \text{discr}(\cdot) \end{bmatrix} & \text{otherwise} \end{cases}$$

$$B_{jl} = \begin{bmatrix} E_{jj} & E_{jl} \\ E_{lj} & 0 \end{bmatrix} + \epsilon I$$

$$\tilde{\phi}(\mathbf{x}; \theta_{ft}, \theta_b)^T B_{jl}^T \sim \begin{cases} \begin{bmatrix} e_j & 0 \\ 0 & 0 \end{bmatrix} & \text{if } \text{discr}(\mathbf{x}) = e_j \\ \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} & \text{otherwise} \end{cases} \quad B_{jl} \tilde{\phi}(\mathbf{x}^*; \theta'_{ft}, \theta'_b)^T \sim \begin{cases} \begin{bmatrix} e_l & 0 \\ 0 & 0 \end{bmatrix}^T & \text{if } \text{discr}(\mathbf{x}^*) = e_l \\ \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}^T & \text{otherwise} \end{cases}$$

• So,

$$k_{jk}(\mathbf{x}, \mathbf{x}^*) \sim \begin{cases} \begin{bmatrix} e_j & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} e_l \\ 0 \end{bmatrix} & \text{if } \text{discr}(\mathbf{x}) = e_j \text{ and } \text{discr}(\mathbf{x}^*) = e_l \\ 0 & \text{otherwise} \end{cases}$$

Thank you!