

Neural Architecture Growth using Tangent Space

Manon Verbockhaven, Guillaume Charpiat
`prenom.nom@inria.fr`

TAU team, LISN, INRIA Saclay, Université Paris-Saclay

Avril 2023

- Introduction, common practices
- Expressivity bottlenecks : definition and formula
- Best parameters move at each layer
- Best neurons to add
- Resultats
- Conclusion

I - Introduction

Introduction

- Medium architecture : approximate any function. (*Universal approximation theorem*)
- Huge architecture : good minimum following the gradient descent. (*Neural Tangent Kernel*)

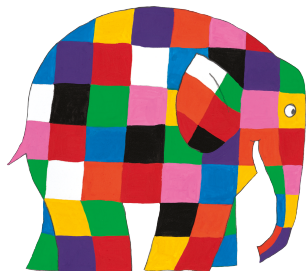
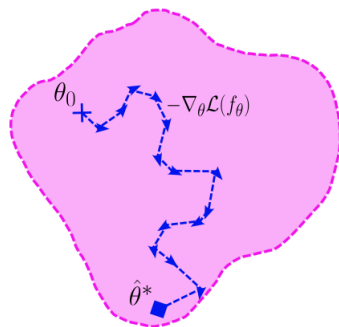


Figure: Neural Network with huge architecture

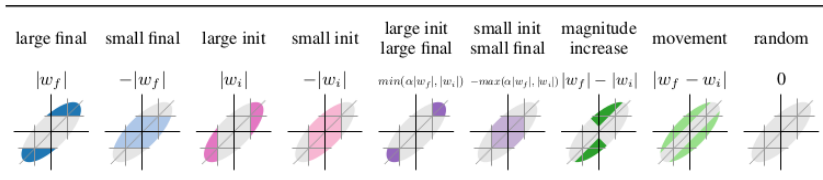
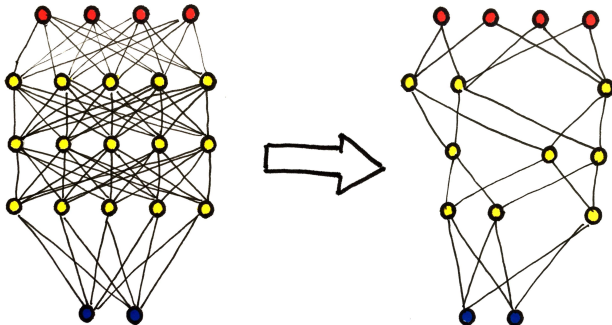
Optimizing NN



$$|\mathcal{L}(\hat{\theta}^*) - \mathcal{L}(\theta^*)| \rightarrow 0$$

Figure: Following the gradient descent on the set of parameters

pruning



NAS : Neural Architecture Search

from scratch
Human design
random search
gradient bases method ¹

meta DL
reinforcement learning ²
multi-task learning
genetic methods ³

Sensitive to exploration hyper-parameters, need a lot of computational resources.

- 1 : GradMax: Growing Neural Networks using Gradient Information, Google Research, Brain Team
- 2 : Neural Architecture Search with Reinforcement Learning, Google Brain
- 3 : Compositional pattern producing networks: A novel abstraction of development, Stanley

Starting with a small neural network

Training a small network:

- faster learning, less memory
- the solution found by gradient descent is poor

Adapt the architecture:

- search for architecture usually done in a discrete way (trial/error)

We propose: to adapt the architecture during learning:

- estimate and localize the potential **expressivity Bottlenecks**, and fill them, without trial/error.

II - Expressivity Bottlenecks, definition and formula

Definitions, Goals, Objectives

For the NN \mathcal{N} with the architecture \mathcal{A} :

- What are *expressivity bottlenecks* ?
 - Where are they ?
 - How to quantify them ?

Optimization Problem

Neural Network Settings

- Dataset $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \in (\mathbb{R}^p \times \mathbb{R}^d)^N \text{ iid} \sim \mathcal{P}$
- Neural Network $f_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^d$
- Loss function $\mathcal{L} : (\mathbb{R}^d)^2 \rightarrow \mathbb{R}^+$

Displacement

Definition (desired update for the function f_θ)

For $\eta > 0$, we define :

$$\mathbf{v}_{\text{goal}}(\cdot) := -\eta \nabla_{f_\theta(\cdot)} \mathcal{L}(f_\theta)$$

$$\mathbf{v}_{\text{goal}}(\mathbf{x}) := -\eta \nabla_{f_\theta(\mathbf{x})} \mathcal{L}(f_\theta(\mathbf{x}), \mathbf{y}) \quad \text{for } (\mathbf{x}, \mathbf{y}) \sim \mathcal{P}$$

Definition (possible update for the function f_θ)

For $\delta\theta$ an incremental direction for the parameters θ , we define :

$$\mathbf{v}(\cdot, \delta\theta) := \frac{\partial f_\theta(\cdot)}{\partial \theta} \delta\theta$$

$$\mathbf{v}(\mathbf{x}, \delta\theta) := \frac{\partial f_\theta(\mathbf{x})}{\partial \theta} \delta\theta$$

Where $\frac{\partial f}{\partial \theta}$ is the Jacobian matrix.

Example

Example (Least Square regression)

Using $\mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{y}) := \|\theta^T \mathbf{x} - \mathbf{y}\|^2$:

$$\begin{aligned} \mathbf{v}_{\text{goal}}(\mathbf{x}) &:= -\eta \nabla_{\theta^T \mathbf{x}} \mathcal{L}(\theta^T \mathbf{x}, \mathbf{y}) \\ &= -\eta 2(\theta^T \mathbf{x} - \mathbf{y}) \end{aligned}$$

$$\mathbf{v}(\mathbf{x}, \delta\theta) = (\delta\theta)^T \mathbf{x}$$

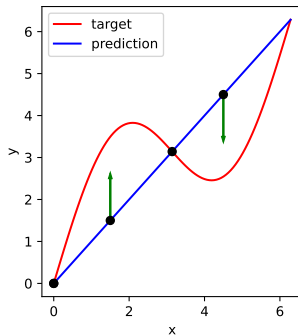
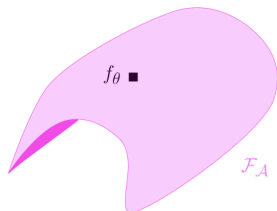


Figure: Linear prediction for a sinusoidal output

The Theory

mathematical objects

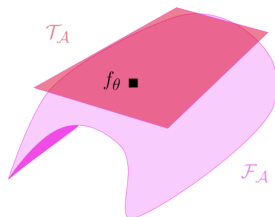
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
-



The Theory

mathematical objects

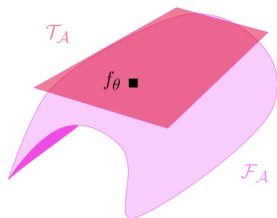
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}} := \mathcal{T}_{\mathcal{A}} = \left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta\theta \mid \text{s.t. } \delta\theta \in \Theta \right\}$ ie
Tangent space in f_{θ}



The Theory

mathematical objects

- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}} := \mathcal{T}_{\mathcal{A}} = \left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta\theta \mid \text{s.t. } \delta\theta \in \Theta \right\}$ ie
Tangent space in f_{θ}

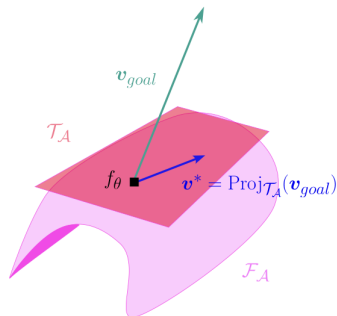


$$g \in \mathcal{T}_{\mathcal{A}} \iff \exists \delta\theta \text{ s.t. } g(\mathbf{x}) = f_{\theta}(\mathbf{x}) + \frac{\partial f(\mathbf{x})}{\partial \theta} \delta\theta$$

The Theory

mathematical objects

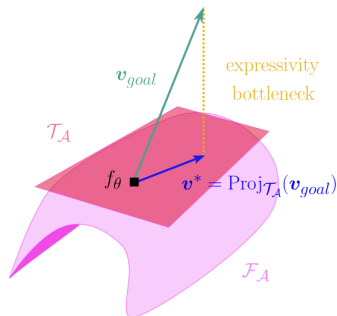
- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}} := \mathcal{T}_{\mathcal{A}} = \left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \mid \text{s.t. } \delta \theta \in \Theta \right\}$ ie
Tangent space in f_{θ}



The Theory

mathematical objects

- Architecture space : $\Theta_{\mathcal{A}}$
- $\mathcal{F}_{\mathcal{A}} = \{f_{\theta} \mid \theta \in \Theta_{\mathcal{A}}\}$
- $\mathcal{T}_{\mathcal{A}}^{f_{\theta}} := \mathcal{T}_{\mathcal{A}} = \left\{ \frac{\partial f_{\theta}}{\partial \theta} \delta \theta \mid \text{s.t. } \delta \theta \in \Theta \right\}$ ie
Tangent space in f_{θ}



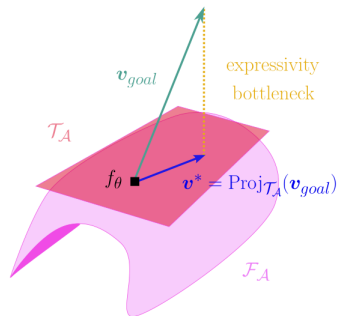
The Theory

Expressivity Bottleneck

$$\arg \min_{\mathbf{v} \in \mathcal{T}_A} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{P}} [\|\mathbf{v}_{\text{goal}}(\mathbf{x}) - \mathbf{v}(\mathbf{x})\|^2]$$

Expressivity Bottleneck approx.

$$\arg \min_{\mathbf{v} \in \mathcal{T}_A} \|\mathbf{V}_{\text{goal}} - \mathbf{V}\|_{\text{Tr}}^2$$



$$\mathbf{V}_{\text{goal}} = (\mathbf{v}_{\text{goal}}(\mathbf{x}_1) \quad \dots \quad \mathbf{v}_{\text{goal}}(\mathbf{x}_n)) \quad \text{and} \quad \mathbf{V} = (\mathbf{v}(\mathbf{x}_1) \quad \dots \quad \mathbf{v}(\mathbf{x}_n))$$

Example

Example

$$f_{\theta}(\mathbf{x}) := \sigma(\theta^T \mathbf{x}) \quad \text{and} \quad \mathcal{L}(\mathbf{x}, \mathbf{y}) := \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2$$

$$\mathbf{v}_{\text{goal}}(\mathbf{x}) = -2(f_{\theta}(\mathbf{x}) - \mathbf{y})$$

$$\mathbf{v}(\mathbf{x}, \delta\theta) = 2(\sigma(\theta^T \mathbf{x}) - \mathbf{y})\sigma'(\theta^T \mathbf{x})\mathbf{x}^T \delta\theta$$

$$\arg \min_{\delta\theta} \mathbb{E}_{\mathcal{P}}[\|\mathbf{v}_{\text{goal}}(\mathbf{x}) - \mathbf{v}(\mathbf{x}, \delta\theta)\|^2]$$

Notations

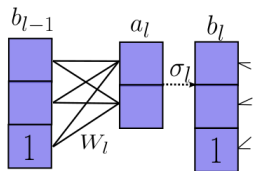


Figure: Feedforward NN

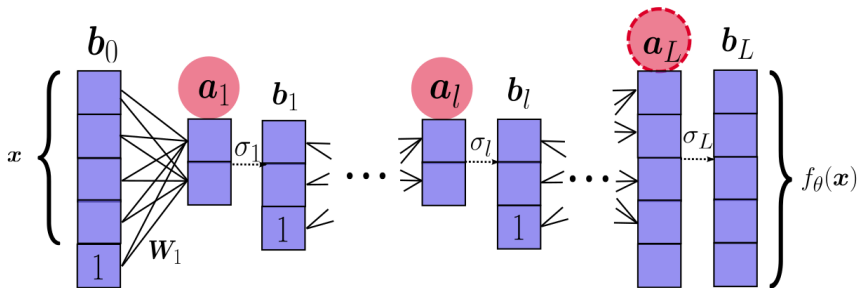
$$\mathbf{b}_0(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix}$$

$$\mathbf{a}_l(\mathbf{x}) = \mathbf{W}_l \mathbf{b}_{l-1}(\mathbf{x})$$

$$\mathbf{b}_l(\mathbf{x}) = \begin{pmatrix} \sigma_l(\mathbf{a}_l(\mathbf{x})) \\ 1 \end{pmatrix}$$

$$f_\theta(\mathbf{x}) = \sigma_L(\mathbf{a}_L(\mathbf{x}))$$

\mathbf{v} and \mathbf{v}_{goal} at layer l



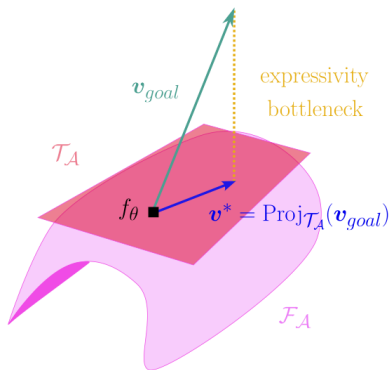
Definition ($\mathbf{v}'_{\text{goal}}$ desired update)

$$\mathbf{v}'_{\text{goal}}(\mathbf{x}) := -\eta \nabla_{\mathbf{a}_l(\mathbf{x})} \mathcal{L}(f_{\theta}(\mathbf{x}), \mathbf{y})$$

Definition (\mathbf{v}' possible update)

$$\mathbf{v}'(\mathbf{x}, \delta\theta) := \frac{\partial \mathbf{a}_l(\mathbf{x})}{\partial \theta} \delta\theta$$

III - Best parameter move



Before Fixing E.B

Non-convex optimization

Considering all the parameters :

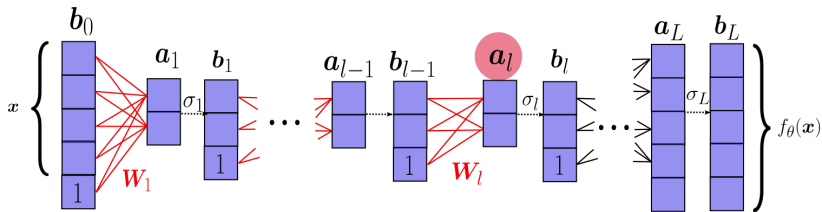
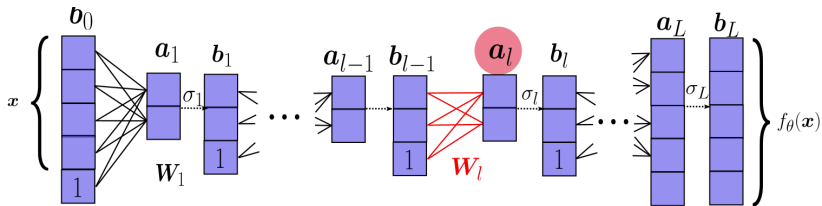
$$\delta\theta^* = \arg \min_{\delta\theta} \|\mathbf{V}_{\text{goal}}^l - \mathbf{V}(\delta\theta)\|^2$$

Convex optimization

Considering only parameters at layer l :

$$\delta\mathbf{W}_l^* = \arg \min_{\delta\theta} \|\mathbf{V}_{\text{goal}}^l - \mathbf{V}(\delta\mathbf{W}_l)\|^2$$

$\delta\theta^*$ vs δW_l^*



Convex Optimization

Linear regression from $\mathbf{b}_{l-1}(\mathbf{x})$ to predict $\mathbf{v}'_{\text{goal}}(\mathbf{x})$:

$$\delta \mathbf{W}_l^* := \arg \min_{\delta \theta} \|\mathbf{V}'_{\text{goal}} - \mathbf{V}'(\delta \mathbf{W}_l)\|^2$$

$$\delta \mathbf{W}_l^* := \arg \min_{\delta \theta} \|\mathbf{V}'_{\text{goal}} - \delta \mathbf{W}_l \mathbf{B}_{l-1}\|^2$$

$$\delta \mathbf{W}_l^* = \frac{1}{n} \mathbf{V}'_{\text{goal}} \mathbf{B}_{l-1}^T \left(\frac{1}{n} \mathbf{B}_{l-1} \mathbf{B}_{l-1}^T \right)^{-1}$$

$$\mathbf{B}_{l-1} := (\mathbf{b}_{l-1}(\mathbf{x}_1) \quad \dots \quad \mathbf{b}_{l-1}(\mathbf{x}_n))$$

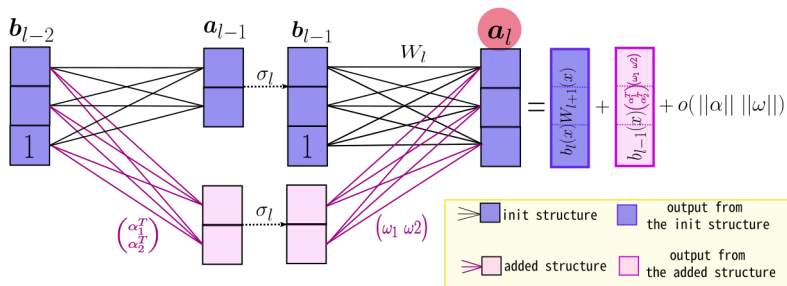
IV - Best neurons to add

Definitions, Goals, Objectives

For the NN \mathcal{N} with the architecture \mathcal{A} :

- What is *expressivity bottlenecks* ?
- Where are they ?
- How to quantify them ?
- **Can we find a method to fix them ?**

Fix E.B and changing the structure



Definition ($\mathbf{V}(\delta\theta_{l-1 \leftrightarrow l}^K)$)

$$\Delta \mathbf{a}_l(\mathbf{x}) := \sum_{k=1}^K \omega_k \sigma_l(\boldsymbol{\alpha}_k^T \mathbf{b}_{l-2}(\mathbf{x})) \quad \text{and} \quad \mathbf{v}(\mathbf{x}, \delta\theta_{l-1 \leftrightarrow l}^K) := \sum_{k=1}^K \omega_k (\boldsymbol{\alpha}_k^T \mathbf{b}_{l-2}(\mathbf{x}))$$

Full Optimization Problem

The Optimal Neurons to reduce Expressivity Bottlenecks at layer l

The optimal neurons n_1, \dots, n_K are defined by $n_i := (\alpha_i, \omega_i)$ and are the solution of the optimization problem :

$$\arg \min_{\delta \mathbf{W}_l, K, \delta \theta_{l-1 \leftrightarrow l}^K} \overbrace{\left\| \mathbf{V}_{\text{goal}}^l - \delta \mathbf{W}_l \mathbf{B}_{l-1} - \boldsymbol{\Omega} \mathbf{A}^T \mathbf{B}_{l-2} \right\|_{\text{Tr}}^2}^{\text{Expressivity Bottleneck}}$$

$$\boldsymbol{\Omega} := (\omega_1 \quad \dots \quad \omega_K) \quad \text{and} \quad \mathbf{A} := (\alpha_1 \quad \dots \quad \alpha_K)$$

Full Optimization Problem

The Optimal Neurons

The optimal neurons n_1, \dots, n_K are defined by $n_i := (\alpha_i, \omega_i)$ and are the solution of the optimization problem :

$$\arg \min_{\delta \mathbf{W}_l, \delta \theta_{l-1 \leftrightarrow l}^K} \left\| \overbrace{\mathbf{V}_{\text{goal}}^l - \delta \mathbf{W}_l \mathbf{B}_{l-1} - \Omega \mathbf{A}^T \mathbf{B}_{l-2}}^{\mathbf{V}_{\text{goal proj}}^l} \right\|_{\text{Tr}}^2$$

$$\Omega := (\omega_1 \quad \dots \quad \omega_K) \quad \text{and} \quad \mathbf{A} := (\alpha_1 \quad \dots \quad \alpha_K)$$

Theorem (Solution)

Define $\mathbf{S} := \frac{1}{n} \mathbf{B}_{l-2} \mathbf{B}_{l-2}^T$ and $\mathbf{N} := \frac{1}{n} \mathbf{B}_{l-2} \mathbf{V}_{goal\ proj}^l$ and $\mathbf{S}^{1/2^{-1}} \mathbf{N} = \sum_{c=1}^R \lambda_c \mathbf{u}_c \mathbf{v}_c^T$ (SVD), then :

maximum number of neuron to add $:= R$

initialization of new neuron k , $(\hat{\alpha}_k^, \hat{\omega}_k^*) := (\mathbf{S}^{1/2 T^{-1}} u_k, v_k)$*

Theorem

Theorem (Impact at first order on Loss)

In $\hat{\theta}_{l \leftrightarrow l+1}^{K,*} := (\hat{\alpha}_k^*, \hat{\omega}_k^*)_{k=1}^K$ with $K \leq K^* = R$, then :

$$\frac{1}{n} \|\mathbf{V}_{goal_{proj}}^{l+1} - \mathbf{V}^{l+1}(\hat{\theta}_{l \leftrightarrow l+1}^{K,*})\|_{\text{Tr}}^2 = \frac{1}{n} \|\mathbf{V}_{goal_{proj}}^{l+1}\|_{\text{Tr}}^2 - \sum_{k=1}^K \lambda_k^2$$

$$\mathcal{L}(f_{\theta \oplus \hat{\theta}_{l \leftrightarrow l+1}^{K,*}}) = \mathcal{L}(f_{\theta}) - \frac{\sigma'_l(0)}{\eta} \sum_{k=1}^K \lambda_k^2 + o(\|\hat{\theta}_{l \leftrightarrow l+1}^{K,*}\|^2)$$

Theorem (Loss variation)

Taylor at order 1 in $\mathbf{V}^{l+1}(\theta_{l \leftrightarrow l+1})$

$$\mathcal{L}(f_{\theta \oplus \theta_{l \leftrightarrow l+1}^k}) \approx \mathcal{L}(f_{\theta}) - \sigma'_l(0) \frac{1}{\eta} \frac{1}{n} \sigma'_l(0) \left\langle \mathbf{V}_{\text{goal proj}}^{l+1}(\delta \mathbf{W}_l), \mathbf{V}^{l+1}(\theta_{l \leftrightarrow l+1}) \right\rangle_{\text{Tr}}$$

At layer l we have :

- we evaluate the expressivity bottleneck $\|\mathbf{V}_{\text{goal}}^l - \mathbf{V}^l\|_{\text{Tr}}$
- we compute the best parameter move without changing the architecture
- we compute the best neurons to add by changing the architecture (finite quantity of neurons)

Neurons selection

- $\lambda_k > \text{seuil}$
- statistical hypothesis test
- Numerical stability
- batchsize for estimation

V - Results

Random neurons vs our approach

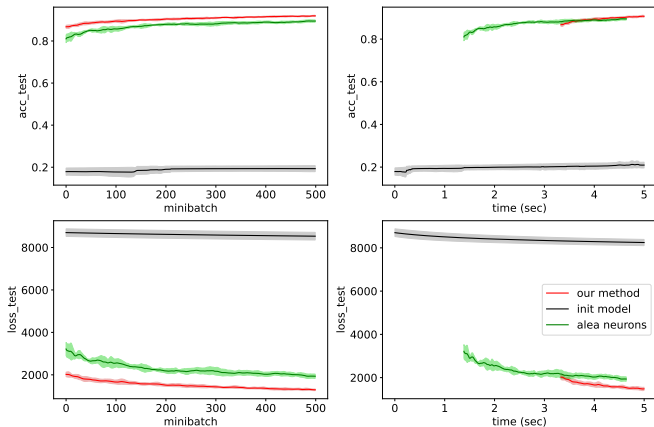


Figure: Acc and Loss on test set for the MNIST dataset. Comparison random add / our method.

init architecture [1, 1]
alea/our architecture [110, 61]

MNIST

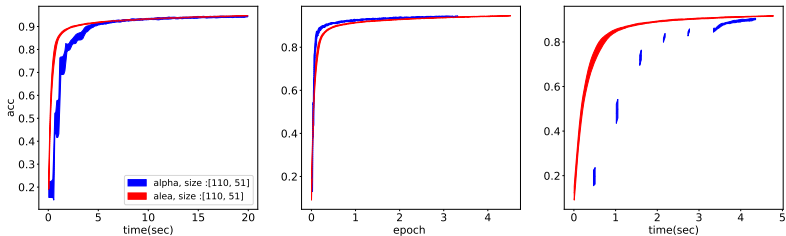


Figure: Accuracy on test set for the MNIST dataset, **model** : usual training with architecture [110, 51] a time 0. **model** : our method with architecture [110, 51].

MNIST

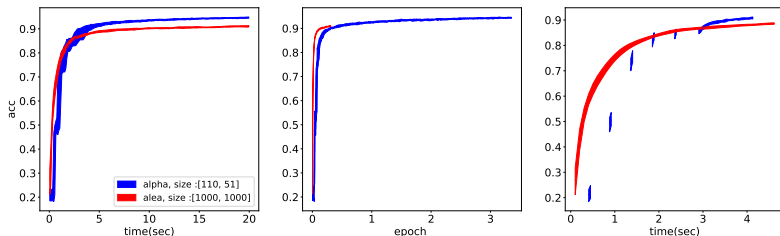
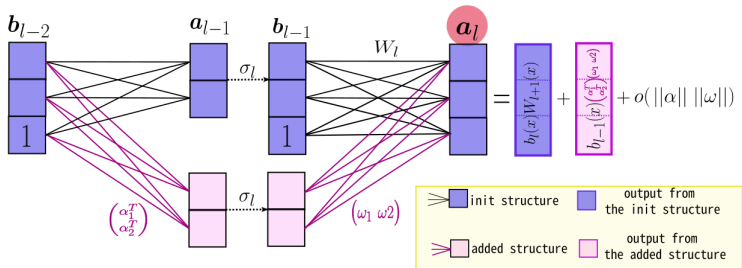


Figure: Accuracy on test set for the MNIST dataset, **model** : usual training with architecture [1000, 1000] a time 0. **model** : our method with architecture [110, 51].

comparaison 30 lancements rouge contre une bleue

Discussion

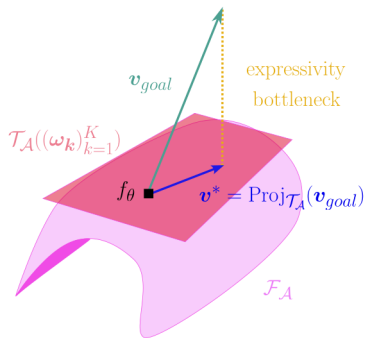
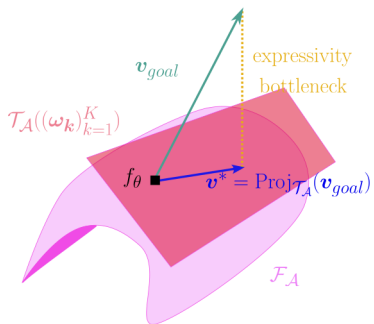
- looking for linear correlation between B_{l-2} and V_{goal}^l
 - connecting to others layers
 - can be extended to higher order moments (further Taylor expansion)
- greedy approach
- future work : convolution, extension to addition of layers (any DAG)



https://www.lri.fr/~gcharpia/Expressivity_bottlenecks_preprint.pdf

Questions ?

?



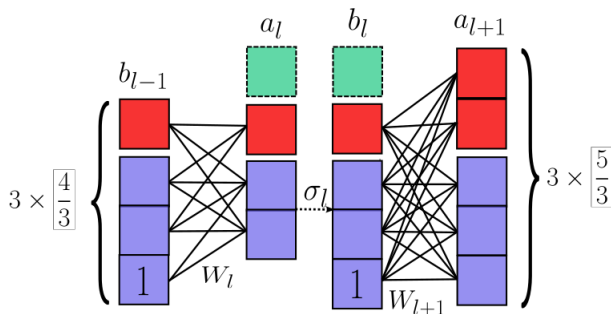
$$\mathbf{a}_{l+1}(\mathbf{x}) = \mathbf{a}_{l+1}(\mathbf{x})|_{(\boldsymbol{\alpha}_k, \boldsymbol{\omega}_k)_{k=1}^K=0} + \sum_{k=1}^K \boldsymbol{\omega}_k \mathbf{b}_{l-1}(\mathbf{x})^T \boldsymbol{\alpha}_k + o(\left(\|(\boldsymbol{\alpha}_k)_{k=1}^K\| + \|(\boldsymbol{\omega}_k)_{k=1}^K\|\right)^2)$$

Estimating the best neurons/update

$$\left(\overbrace{(\hat{\alpha}_k^*, \hat{\omega}_k^*)_{k=1}^{K^*}}^{\hat{\theta}_{l \leftrightarrow l+1}^{K^*, *}}, K^* \right) := \arg \min_{(\alpha_k, \omega_k)_{k=1}^K, K} \frac{1}{\boxed{n}} \left\{ \left\| \mathbf{V}_{\text{goal}_{proj}}^{l+1} - \mathbf{V}^{l+1}(\theta_{l \leftrightarrow l+1}^K) \right\|_{\text{Tr}}^2 \right\}$$

With $\mathbf{V} := (v_{\text{goal}}(\mathbf{x}_1) \quad \dots \quad v_{\text{goal}}(\mathbf{x}_n))$

Estimating the best neurons and the naturel gradient



$$n \leftarrow n \times \max\left(\sqrt{\frac{4}{3}}, \sqrt{\frac{5}{3}}\right)$$

Defining the updates

$\theta^{\delta(l+1)}(\gamma_0) = (\mathbf{W}_1, \dots, \mathbf{W}_{l+1} + \gamma_0 \delta \mathbf{W}_{l+1}^*, \dots, \mathbf{W}_L)$ and
 $\hat{\theta}_{l \leftrightarrow l+1}^{K,*}(\gamma) = (\gamma \hat{\alpha}_k^*, \hat{\omega}_k^*)_{k=1}^K$, we apply a line search algorithm to find a local minimum of the loss function :

$$\gamma_0^* := \arg \min_{\gamma_0 \in \mathcal{V}(0)} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f_{\theta^{\delta(l+1)}(\gamma_0)}(X_i), Y_i) \quad (1)$$

$$\gamma^* := \arg \min_{\gamma \in \mathcal{V}(0)} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f_{\theta^{\delta(l+1)}(\gamma_0^*) \oplus \hat{\theta}_{l \leftrightarrow l+1}^{K,*}(\gamma)}(X_i), Y_i) \quad (2)$$

where $\mathcal{V}(0)$ is a positive neighbourhood of zero.