Scaling ResNets in the large-depth regime

NON-LINEAR AND HIGH DIMENSIONAL INFERENCE, IHP, OCTOBER 3RD 2022

Adeline Fermanian





Joint work with



Gérard Biau Sorbonne University



Pierre Marion SORBONNE UNIVERSITY



Jean-Philippe Vert OWKIN



Learning with ResNets

Scaling deep ResNets

Scaling in the continuous-time setting

Beyond initialization



Learning with ResNets

Scaling deep ResNets

Scaling in the continuous-time setting

Beyond initialization

How most people see the supervised learning problem

Learn how to build an image-recognizing convolutional neural network with Python and Keras in less than 15minutes!



Fabian Bosler Oct 5, 2019 · 10 min read *





https://towardsdatascience.com/cat-dog-or-elon-musk-145658489730

How machine learners see the supervised learning problem



https://medium.datadriveninvestor.com/depth-estimation-with-deep-neural-networks-part-2-81ee374888eb

Sol: understand the relationship between $x \in \mathbb{R}^{n_{in}}$ and $y \in \mathbb{R}^{n_{out}}$.

- Soal: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.

- Sol: understand the relationship between $x \in \mathbb{R}^{n_{in}}$ and $y \in \mathbb{R}^{n_{out}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- **>** Model: $\{F_{\pi} : \mathbb{R}^{n_{\text{in}}} \mapsto \mathbb{R}^{n_{\text{out}}}, \pi \in \Pi\}.$

- Sol: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- $\mathbf{E} \quad \mathbf{Model:} \ \{F_{\pi}: \mathbb{R}^{n_{\text{in}}} \mapsto \mathbb{R}^{n_{\text{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.

- Soal: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- $\mathbf{E} \quad \mathbf{Model:} \ \{F_{\pi}: \mathbb{R}^{n_{\mathsf{in}}} \mapsto \mathbb{R}^{n_{\mathsf{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.
- **Regression:** $\ell(F_{\pi}(x), y) = (y F_{\pi}(x))^2$

- Soal: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- $\mathbf{E} \quad \mathbf{Model:} \ \{F_{\pi}: \mathbb{R}^{n_{\mathsf{in}}} \mapsto \mathbb{R}^{n_{\mathsf{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.
- **>** Regression: $\ell(F_{\pi}(x), y) = (y F_{\pi}(x))^2$ Binary classification: $\ell(F_{\pi}(x), y) = \mathbb{1}_{[yF_{\pi}(x) \leq 0]}$.

- Soal: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- $\mathbf{E} \quad \mathbf{Model:} \ \{F_{\pi}: \mathbb{R}^{n_{\mathsf{in}}} \mapsto \mathbb{R}^{n_{\mathsf{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.
- **>** Regression: $\ell(F_{\pi}(x), y) = (y F_{\pi}(x))^2$ Binary classification: $\ell(F_{\pi}(x), y) = \mathbb{1}_{[yF_{\pi}(x) \leq 0]}$.

Theoretical risk minimization: choose

$$\pi^{\star} \in \operatorname*{argmin}_{\pi \in \Pi} \mathscr{L}(\pi) = \mathbb{E}(\ell(F_{\pi}(x), y)).$$

- Soal: understand the relationship between $x \in \mathbb{R}^{n_{\text{in}}}$ and $y \in \mathbb{R}^{n_{\text{out}}}$.
- **>** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- $\mathbf{E} \quad \mathbf{Model:} \ \{F_{\pi}: \mathbb{R}^{n_{\mathsf{in}}} \mapsto \mathbb{R}^{n_{\mathsf{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.
- **>** Regression: $\ell(F_{\pi}(x), y) = (y F_{\pi}(x))^2$ Binary classification: $\ell(F_{\pi}(x), y) = \mathbb{1}_{[yF_{\pi}(x) \leq 0]}$.

Theoretical risk minimization: choose

$$\pi^{\star} \in \operatorname*{argmin}_{\pi \in \Pi} \mathscr{L}(\pi) = \mathbb{E}(\ell(F_{\pi}(x), y)).$$

Empirical risk minimization: choose

$$\pi_n \in \operatorname*{argmin}_{\pi \in \Pi} \mathscr{L}_n(\pi) = \frac{1}{n} \sum_{i=1}^n \ell(F_\pi(x_i), y_i).$$

- ig> Goal: understand the relationship between $x\in \mathbb{R}^{n_{ ext{in}}}$ and $y\in \mathbb{R}^{n_{ ext{out}}}.$
- **)** Data: $(x_1, y_1), \ldots, (x_n, y_n) \in \mathbb{R}^{n_{\text{in}}} \times \mathbb{R}^{n_{\text{out}}}$, i.i.d. $\sim (x, y)$.
- **>** Model: $\{F_{\pi} : \mathbb{R}^{n_{\text{in}}} \mapsto \mathbb{R}^{n_{\text{out}}}, \pi \in \Pi\}.$
- **>** Loss function $\ell : \mathbb{R}^{n_{\text{out}}} \times \mathbb{R}^{n_{\text{out}}} \to \mathbb{R}_+$.
- **Regression**: $\ell(F_{\pi}(x), y) = (y F_{\pi}(x))^2$ Binary classification: $\ell(F_{\pi}(x), y) = \mathbb{1}_{[yF_{\pi}(x) \leq 0]}$.

> Theoretical risk minimization: choose

$$\pi^* \in \operatorname*{argmin}_{\pi \in \Pi} \mathscr{L}(\pi) = \mathbb{E}(\ell(F_{\pi}(x), y)).$$

Empirical risk minimization: choose

$$\pi_n \in \operatorname*{argmin}_{\pi \in \Pi} \mathscr{L}_n(\pi) = \frac{1}{n} \sum_{i=1}^n \ell(F_\pi(x_i), y_i).$$

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

$$h_0 = Ax, \quad h_{k+1} = h_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$$

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

$$h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$$

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

 $h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$

> Different forms for $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$ = different architectures.

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

$$h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$$

> Different forms for $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$ = different architectures.

Original Parametric Simple General ResNet

$$f(\mathbf{h}_{k}, \theta_{k+1}) = V_{k+1} \operatorname{ReLU}(W_{k+1}\mathbf{h}_{k} + b_{k+1})$$

▷ ReLU(x) = max(x, 0) = activation function ▷ $\theta_k = (W_k, b_k)$ = weight matrice + bias ▷ $\pi = (A, B, (V_k)_{1 \le k \le L}, (\theta_k)_{1 \le k \le L})$





Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

$$h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$$

> Different forms for $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$ = different architectures.

Original Parametric Simple General ResNet

$$f(\mathbf{h}_{k},\theta_{k+1}) = V_{k+1}\sigma(W_{k+1}\mathbf{h}_{k} + b_{k+1})$$

 $\triangleright \sigma =$ activation function

- $\triangleright \theta_k = (W_k, b_k) =$ weight matrice + bias
- $\triangleright \ \pi = (A, B, (V_k)_{1 \leq k \leq L}, (\theta_k)_{1 \leq k \leq L})$





Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

$$h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$$

> Different forms for $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$ = different architectures.

Original Parametric Simple General ResNet

$$f(\mathbf{h}_{\mathbf{k}}, \theta_{k+1}) = V_{k+1}\sigma(\mathbf{h}_{\mathbf{k}})$$

$\triangleright \sigma = ext{activation function}$

$$\begin{array}{l} \triangleright \ \theta_k = \emptyset \\ \triangleright \ \pi = (A, B, (V_k)_{1 \leqslant k \leqslant L}) \end{array}$$



He et al. (2016)

Sequence of hidden states $h_0, \ldots, h_L \in \mathbb{R}^d$ defined by recurrence:

 $h_0 = Ax, \quad h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1}), \quad F_{\pi}(x) = Bh_L.$

> Different forms for $f : \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$ = different architectures.

Original Parametric Simple General ResNet

$$f(\mathbf{h}_{\mathbf{k}}, \theta_{k+1}) = V_{k+1}g(\mathbf{h}_{\mathbf{k}}, \theta_{k+1})$$

 $\triangleright \ g: \mathbb{R}^d \times \mathbb{R}^p \to \mathbb{R}^d$

 $\triangleright \theta_k = \text{parameters}$

 $\triangleright \ \pi = (A, B, (V_k)_{1 \leq k \leq L}, (\theta_k)_{1 \leq k \leq L})$







The revolution of ResNets



Examples from the ImageNet dataset

https://blog.roboflow.com/introduction-to-imagenet

The revolution of ResNets



ImageNet performance over time

https://semiengineering.com/new-vision-technologies-for-real-world-applications

The revolution of ResNets



ImageNet performance over time

https://semiengineering.com/new-vision-technologies-for-real-world-applications



Deep learning \rightarrow neural ODE \leftarrow ODE

> Traditional neural networks

 $h_{k+1} = f(h_k, \theta_{k+1})$

> Traditional neural networks

 $h_{k+1} = f(h_k, \theta_{k+1})$

Residual neural networks (He et al., 2016)

 $h_{k+1} = \mathbf{h}_k + f(h_k, \theta_{k+1})$

> Traditional neural networks

 $h_{k+1} = f(h_k, \theta_{k+1})$

Residual neural networks (He et al., 2016)

$$h_{k+1} = \mathbf{h}_{k} + \frac{1}{L}f(h_k, \theta_{k+1})$$

> Traditional neural networks

 $h_{k+1} = f(h_k, \theta_{k+1})$

Residual neural networks (He et al., 2016)

$$h_{k+1} = \mathbf{h}_{k} + \frac{1}{L}f(h_k, \theta_{k+1})$$

> Neural ODE (Chen et al., 2018)

 $dH_t = f(H_t, \Theta_t) dt$

> Traditional neural networks

 $h_{k+1} = f(h_k, \theta_{k+1})$

Residual neural networks (He et al., 2016)

 $h_{k+1} = oldsymbol{h}_k + rac{1}{L}f(h_k, heta_{k+1})$

> Neural ODE (Chen et al., 2018)

 $dH_t = f(H_t, \Theta_t) dt$



New network architectures: Runge-Kutta networks



Benning et al. (2019)

New network architectures: antisymmetric networks



Chang et al. (2019)

In summary

ResNet
 $h_0 = Ax$ Neural ODE
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{L}f(h_k, \theta_{k+1})$ $H_0 = Ax$ $F_{\pi}(x) = Bh_T$ $dH_t = f(H_t, \Theta_t)dt$ $F_{\pi}(x) = Bh_T$ $F_{\Pi}(x) = BH_1$ $f(h, \theta) = V\sigma(Wh + b)$



Learning with ResNets

Scaling deep ResNets

Scaling in the continuous-time setting

Beyond initialization
Stability at initialization

> Original ResNet:

$$egin{aligned} h_0 &= Ax\ h_{k+1} &= h_k + V_{k+1} \operatorname{ReLU}(W_{k+1}h_k)\ F_\pi(x) &= Bh_L. \end{aligned}$$

> Original ResNet:

$$egin{aligned} h_0 &= Ax\ h_{k+1} &= h_k + V_{k+1} \operatorname{ReLU}(W_{k+1}h_k)\ F_\pi(x) &= Bh_L. \end{aligned}$$

At initialization: A, B, $(V_k)_{1 \le k \le L}$, and $(W_k)_{1 \le k \le L}$ are i.i.d. Gaussian matrices.

> Original ResNet:

$$egin{aligned} h_0 &= Ax\ h_{k+1} &= h_k + V_{k+1} \operatorname{ReLU}(W_{k+1}h_k)\ F_\pi(x) &= Bh_L. \end{aligned}$$

At initialization: A, B, $(V_k)_{1 \le k \le L}$, and $(W_k)_{1 \le k \le L}$ are i.i.d. Gaussian matrices.



> Original ResNet:

$$egin{aligned} h_0 &= Ax\ h_{k+1} &= h_k + V_{k+1} \operatorname{ReLU}(W_{k+1}h_k)\ F_\pi(x) &= Bh_L. \end{aligned}$$

At initialization: A, B,
$$(V_k)_{1 \le k \le L}$$
, and $(W_k)_{1 \le k \le L}$ are i.i.d. Gaussian matrices.

Solution: batch normalization or scaling.



Scaling ResNets

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

Scaling ResNets

A scaling factor $1/L^{\beta}$:

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

> Question: choice of β .

Scaling ResNets

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

- **>** Question: choice of β .
- $\beta = 0$ (original ResNets)?

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

- **>** Question: choice of β .
- $\beta = 0$ (original ResNets)? $\beta = 1$ (neural ODE)?

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

- **>** Question: choice of β .
- $\beta = 0$ (original ResNets)? $\beta = 1$ (neural ODE)?
- > Many empirical studies, no consensus.

$$h_{k+1} = h_k + \frac{1}{L^{\beta}} V_{k+1} \operatorname{ReLU}(W_{k+1}h_k).$$

- **>** Question: choice of β .
- $\beta = 0$ (original ResNets)? $\beta = 1$ (neural ODE)?
- > Many empirical studies, no consensus.
- > Our approach: mathematical analysis at initialization.



(a) $\|h_L - h_0\| / \|h_0\|, \beta = 1$ (b) $\|h_L - h_0\| / \|h_0\|, \beta = 0.25$ (c) $\|h_L - h_0\| / \|h_0\|, \beta = 0.5$



With an i.i.d. initialization, the critical value for scaling is $\beta = 1/2$.



- With an i.i.d. initialization, the critical value for scaling is $\beta = 1/2$.
- > Not the ODE scaling! 😌

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

- 1. If $\beta > 1/2$
- 2. If $\beta < 1/2$
- 3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}} 1_{L \to \infty} 0$.
2. If $\beta < 1/2$

3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}}_{L \to \infty} 0.$ \rightarrow identity

2. If $\beta < 1/2$

3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}}{L \to \infty} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}}{L \to \infty} \infty.$
3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}}{L \to \infty} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}}{L \to \infty} \infty.$ \rightarrow explosion
3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}} 1_{L \to \infty} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}} 1_{L \to \infty} \infty.$ \rightarrow explosion

3. If $\beta = 1/2$ then, with probability at least $1 - \delta$,

$$\exp\left(\frac{3}{8} - \sqrt{\frac{22}{d\delta}}\right) - 1 < \frac{\|h_L - h_0\|^2}{\|h_0\|^2} < \exp\left(1 + \sqrt{\frac{10}{d\delta}}\right) + 1.$$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}} 1_{L \to \infty} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $||h_L - h_0|| / ||h_0|| \xrightarrow{\mathbb{P}} 1_{L \to \infty} \infty.$ \rightarrow explosion

3. If $\beta = 1/2$ then, with probability at least $1 - \delta$,

$$\exp\left(\frac{3}{8} - \sqrt{\frac{22}{d\delta}}\right) - 1 < \frac{\|h_L - h_0\|^2}{\|h_0\|^2} < \exp\left(1 + \sqrt{\frac{10}{d\delta}}\right) + 1. \quad \to \mathsf{stability}$$



> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

> Target: $||p_0 - p_L|| / ||p_L||$ when L is large.

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

- Target: $||p_0 p_L|| / ||p_L||$ when L is large.
- **Backpropagation** formula:

$$p_k = p_{k+1} + rac{1}{L^eta} rac{\partial g(h_k, heta_{k+1})^ op}{\partial h} V_{k+1}^ op p_{k+1}$$

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

- Target: $||p_0 p_L|| / ||p_L||$ when L is large.
- **Backpropagation** formula:

$$p_k = p_{k+1} + rac{1}{L^eta} rac{\partial g(h_k, heta_{k+1})^ op}{\partial h} V_{k+1}^ op p_{k+1} \quad o ext{ wrong way}.$$

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

- Target: $||p_0 p_L|| / ||p_L||$ when L is large.
- **Backpropagation** formula:

$$p_k = p_{k+1} + rac{1}{L^eta} rac{\partial g(h_k, heta_{k+1})^ op}{\partial h} V_{k+1}^ op p_{k+1} \quad o ext{wrong way}.$$

> Our approach: with $q_k(z) = \frac{\partial h_k}{\partial h_0} z$,

$$q_{k+1}(z) = q_k(z) + \frac{1}{L^{\beta}} V_{k+1} \frac{\partial g(h_k, \theta_{k+1})}{\partial h} q_k(z)$$

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

- Target: $||p_0 p_L|| / ||p_L||$ when L is large.
- > Backpropagation formula:

$$p_k = p_{k+1} + rac{1}{L^eta} rac{\partial g(h_k, heta_{k+1})^ op}{\partial h} V_{k+1}^ op p_{k+1} \quad o ext{wrong way}.$$

> Our approach: with $q_k(z) = \frac{\partial h_k}{\partial h_0} z$,

$$q_{k+1}(z) = q_k(z) + \frac{1}{L^{\beta}} V_{k+1} \frac{\partial g(h_k, \theta_{k+1})}{\partial h} q_k(z) \quad \to \text{flow of information} = \checkmark.$$

> Objective: assess the backwards dynamics of the gradients $p_k = \frac{\partial \mathscr{L}_n}{\partial h_k} \in \mathbb{R}^d$.

- Target: $||p_0 p_L|| / ||p_L||$ when L is large.
- > Backpropagation formula:

$$p_k = p_{k+1} + rac{1}{L^eta} rac{\partial g(h_k, heta_{k+1})^ op}{\partial h} V_{k+1}^ op p_{k+1} \quad o ext{wrong way}.$$

> Our approach: with $q_k(z) = \frac{\partial h_k}{\partial h_0} z$,

$$q_{k+1}(z) = q_k(z) + \frac{1}{L^{\beta}} V_{k+1} \frac{\partial g(h_k, \theta_{k+1})}{\partial h} q_k(z) \quad \to \text{flow of information} = \checkmark.$$

Conclusion with

$$\frac{\|p_0\|^2}{\|p_L\|^2} = \mathbb{E}_{z \sim \mathcal{N}(0, I_d)} \left(\left| \left(\frac{p_L}{\|p_L\|} \right)^\top q_L(z) \right|^2 \right).$$



Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

- 1. If $\beta > 1/2$
- 2. If $\beta < 1/2$
- 3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||p_0 - p_L|| / ||p_L|| \xrightarrow{\mathbb{P}} L \to \infty 0$.
2. If $\beta < 1/2$

3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||p_0 - p_L|| / ||p_L|| \xrightarrow{\mathbb{P}}_{L \to \infty} 0.$ \rightarrow identity
2. If $\beta < 1/2$

3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $\|p_0 - p_L\| / \|p_L\| \xrightarrow{\mathbb{P}} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $\mathbb{E}(\|p_0 - p_L\| / \|p_L\|) \xrightarrow{L \to \infty} \infty.$
3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||p_0 - p_L|| / ||p_L|| \xrightarrow{\mathbb{P}} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $\mathbb{E}(||p_0 - p_L|| / ||p_L||) \xrightarrow{L \to \infty} \infty.$ \rightarrow explosion

3. If $\beta = 1/2$

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||p_0 - p_L|| / ||p_L|| \xrightarrow{\mathbb{P}} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $\mathbb{E}(||p_0 - p_L|| / ||p_L||) \xrightarrow{L \to \infty} \infty.$ \rightarrow explosion

3. If $\beta = 1/2$ then

$$\exp\left(\frac{1}{2}\right) - 1 \leqslant \mathbb{E}\left(\frac{\|p_0 - p_L\|^2}{\|p_L\|^2}\right) \leqslant \exp(4) - 1.$$
Scaling with standard initialization – Gradients

Theorem

Assumption: the entries of $\sqrt{d} V_k$ and $\sqrt{d} W_k$ are symmetric i.i.d. sub-Gaussian.

1. If
$$\beta > 1/2$$
 then $||p_0 - p_L|| / ||p_L|| \xrightarrow{\mathbb{P}} 0.$ \rightarrow identity
2. If $\beta < 1/2$ then $\mathbb{E}(||p_0 - p_L|| / ||p_L||) \xrightarrow{L \to \infty} \infty.$ \rightarrow explosion

3. If $\beta = 1/2$ then

$$\exp\left(\frac{1}{2}\right) - 1 \leqslant \mathbb{E}\left(\frac{\|p_0 - p_L\|^2}{\|p_L\|^2}\right) \leqslant \exp(4) - 1. \quad \to \text{stability}$$

Stability – output/gradients



(a) Distribution of $\|h_L\|/\|h_0\|$



Simple ResNet: $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$.

- Simple ResNet: $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$.
- > The entries of V_k are i.i.d. $\mathcal{N}(0, 1/d)$.

- Simple ResNet: $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$.
- The entries of V_k are i.i.d. $\mathcal{N}(0, 1/d)$.

> For $\mathbf{B}: [0,1] \to \mathbb{R}^{d \times d}$ a $(d \times d)$ -dimensional Brownian motion

- Simple ResNet: $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$.
- The entries of V_k are i.i.d. $\mathcal{N}(0, 1/d)$.

> For $\mathbf{B}: [0,1] \to \mathbb{R}^{d \times d}$ a $(d \times d)$ -dimensional Brownian motion

$$\mathbf{B}_{(k+1)/L,i,j} - \mathbf{B}_{k/L,i,j} \sim \mathcal{N}\left(0, \frac{1}{L}\right).$$

- Simple ResNet: $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$.
- The entries of V_k are i.i.d. $\mathcal{N}(0, 1/d)$.
- **>** For $\mathbf{B}: [0,1] \to \mathbb{R}^{d \times d}$ a $(d \times d)$ -dimensional Brownian motion

$$\mathbf{B}_{(k+1)/L,i,j} - \mathbf{B}_{k/L,i,j} \sim \mathcal{N}\left(0,\frac{1}{L}\right).$$

> Consequence:

$$h_0 = Ax, \quad h_{k+1}^{\top} = h_k^{\top} + \frac{1}{\sqrt{d}}\sigma(h_k^{\top})(\mathbf{B}_{(k+1)/L} - \mathbf{B}_{k/L}), \quad 0 \leq k \leq L-1.$$

SDE regime

ResNetNeural SDE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{\sqrt{L}}V_{k+1}\sigma(h_k)$ $dH_t^\top = \frac{1}{\sqrt{d}}\sigma(H_t^\top) dB_t$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

SDE regime

ResNetNeural SDE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$ $dH_t^\top = \frac{1}{\sqrt{d}} \sigma(H_t^\top) dB_t$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

Proposition

Assumption: the entries of V_k are i.i.d. Gaussian $\mathcal{N}(0, 1/d)$ and σ is Lipschitz continuous.

SDE regime

ResNetNeural SDE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{\sqrt{L}} V_{k+1} \sigma(h_k)$ $dH_t^\top = \frac{1}{\sqrt{d}} \sigma(H_t^\top) dB_t$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

Proposition

Assumption: the entries of V_k are i.i.d. Gaussian $\mathcal{N}(0, 1/d)$ and σ is Lipschitz continuous. Then the SDE has a unique solution H and, for any $0 \le k \le L$,

$$\mathbb{E} ig(\|H_{k/L} - h_k\| ig) \leqslant rac{C}{\sqrt{L}}.$$



For deep ResNets with i.i.d. initialization:

Summary so far

For deep ResNets with i.i.d. initialization:

- \triangleright the critical value for scaling is $\beta = 1/2$
- $\triangleright\,$ this value corresponds in the deep limit to a SDE.

Summary so far

For deep ResNets with i.i.d. initialization:

- \triangleright the critical value for scaling is $\beta = 1/2$
- $\triangleright\,$ this value corresponds in the deep limit to a SDE.

Remaining questions:

- ▷ Can we obtain other limits? For example ODEs?
- ▷ Do they correspond to the same critical value?

Summary so far

For deep ResNets with i.i.d. initialization:

- \triangleright the critical value for scaling is $\beta = 1/2$
- \triangleright this value corresponds in the deep limit to a SDE.

Remaining questions:

- ▷ Can we obtain other limits? For example ODEs?
- Do they correspond to the same critical value?

Key: link between β and the weight distributions.



Learning with ResNets

Scaling deep ResNets

Scaling in the continuous-time setting

Beyond initialization

ldea: the weights $(V_k)_{1 \leq k \leq L}$ and $(\theta_k)_{1 \leq k \leq L}$ are discretizations of smooth functions.

> Idea: the weights $(V_k)_{1 \le k \le L}$ and $(θ_k)_{1 \le k \le L}$ are discretizations of smooth functions. **>** $(V_k)_{1 \le k \le L} \hookrightarrow \mathscr{V} : [0, 1] \to \mathbb{R}^{d \times d}$ > Idea: the weights (V_k)_{1≤k≤L} and (θ_k)_{1≤k≤L} are discretizations of smooth functions.
 > (V_k)_{1≤k≤L} → 𝒴 : [0,1] → ℝ^{d×d} (θ_k)_{1≤k≤L} → Θ : [0,1] → ℝ^p.

> Idea: the weights (V_k)_{1≤k≤L} and (θ_k)_{1≤k≤L} are discretizations of smooth functions.
> (V_k)_{1≤k≤L} → 𝒴 : [0,1] → ℝ^{d×d} (θ_k)_{1≤k≤L} → Θ : [0,1] → ℝ^p.
> Model:

$$h_0 = Ax, \quad h_{k+1} = h_k + \frac{1}{L} V_{k+1} g(h_k, \theta_{k+1}), \quad 0 \le k \le L - 1,$$

where $V_k = \mathscr{V}_{k/L}$ and $\theta_k = \Theta_{k/L}$.

> Idea: the weights (V_k)_{1≤k≤L} and (θ_k)_{1≤k≤L} are discretizations of smooth functions.
> (V_k)_{1≤k≤L} → 𝒴 : [0,1] → ℝ^{d×d} (θ_k)_{1≤k≤L} → Θ : [0,1] → ℝ^p.
> Model:

$$h_0 = Ax, \quad h_{k+1} = h_k + \frac{1}{L} V_{k+1} g(h_k, \theta_{k+1}), \quad 0 \le k \le L - 1,$$

where $V_k = \mathscr{V}_{k/L}$ and $\theta_k = \Theta_{k/L}$.

Assumption: the stochastic processes \mathscr{V} and Θ are a.s. Lipschitz continuous and bounded.

> Idea: the weights (V_k)_{1≤k≤L} and (θ_k)_{1≤k≤L} are discretizations of smooth functions.
> (V_k)_{1≤k≤L} → 𝒴 : [0,1] → ℝ^{d×d} (θ_k)_{1≤k≤L} → Θ : [0,1] → ℝ^p.
> Model:

$$h_0 = Ax, \quad h_{k+1} = h_k + rac{1}{L} V_{k+1} g(h_k, heta_{k+1}), \quad 0 \leqslant k \leqslant L - 1,$$

where $V_k = \mathscr{V}_{k/L}$ and $\theta_k = \Theta_{k/L}$.

Assumption: the stochastic processes \mathscr{V} and Θ are a.s. Lipschitz continuous and bounded.

Example: the entries of \mathscr{V} and Θ are independent Gaussian processes with zero expectation and covariance $K(x, x') = \exp(-\frac{(x-x')^2}{2\ell^2})$.





Scaling and weight regularity



<mark>(a)</mark> l.i.d.

Scaling and weight regularity



ODE regime

ResNetNeural ODE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{L} V_{k+1} g(h_k, \theta_{k+1})$ $dH_t = \mathscr{V}_t g(H_t, \Theta_t) dt$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

ODE regime

ResNetNeural ODE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{L} V_{k+1} g(h_k, \theta_{k+1})$ $dH_t = \mathscr{V}_t g(H_t, \Theta_t) dt$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

Proposition

Assumption: the function *g* is Lipschitz continuous on compact sets.

ODE regime

ResNetNeural ODE
$$h_0 = Ax$$
 $H_0 = Ax$ $h_{k+1} = h_k + \frac{1}{L} V_{k+1} g(h_k, \theta_{k+1})$ $dH_t = \mathscr{V}_t g(H_t, \Theta_t) dt$ $F_{\pi}(x) = Bh_L$ $F_{\Pi}(x) = BH_1$

Proposition

Assumption: the function *g* is Lipschitz continuous on compact sets.

Then the ODE has a unique solution H and, a.s., for any $0 \leq k \leq L$,

$$\|H_{k/L}-h_k\|\leqslant rac{c}{L}.$$





> Again <u>3 cases</u>: identity/explosion/stability.



Again 3 cases: identity/explosion/stability.
 With a smooth initialization, the critical scaling is β = 1.



- Again <u>3 cases</u>: identity/explosion/stability.
- With a smooth initialization, the critical scaling is $\beta = 1$.
- It is the scaling that corresponds in the deep limit to an ODE.

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

- 1. If $\beta > 1$
- 2. If $\beta = 1$

3. If $\beta < 1$

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \xrightarrow{L \to \infty} 0$.
2. If $\beta = 1$

3. If $\beta < 1$

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \xrightarrow{L \to \infty} 0.$ \rightarrow identity
2. If $\beta = 1$

3. If $\beta < 1$
Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $\|h_L - h_0\| / \|h_0\| \xrightarrow{L \to \infty} 0.$ \rightarrow identity

2. If
$$\beta = 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \leq c$.

3. If $\beta < 1$

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \xrightarrow{L \to \infty} 0.$ \rightarrow identity

2. If
$$\beta = 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \leq c$. \rightarrow stability

3. If $\beta < 1$

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $\|h_L - h_0\| / \|h_0\| \xrightarrow{L \to \infty} 0.$ \rightarrow identity

2. If
$$\beta = 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \leq c$. \rightarrow stability

3. If
$$\beta < 1$$
 + assumptions, then $\max_k \frac{\|h_k - h_0\|}{\|h_0\|} \xrightarrow{L \to \infty} \infty$.

Theorem

Assumption: $\mathscr V$ and Θ are a.s. Lipschitz continuous and bounded.

1. If
$$\beta > 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \xrightarrow{L \to \infty} 0.$ \rightarrow identity

2. If
$$\beta = 1$$
 then, a.s., $||h_L - h_0|| / ||h_0|| \leq c$. \rightarrow stability

3. If
$$\beta < 1$$
 + assumptions, then $\max_k \frac{\|h_k - h_0\|}{\|h_0\|} \xrightarrow{L \to \infty} \infty$. \to explosion

Intermediate regimes

Challenge: describe the transition between the i.i.d. and smooth cases.

- **Challenge:** describe the transition between the i.i.d. and smooth cases.
- We initialize the weights as increments of a fractional Brownian motion $(B_t^H)_{t \in [0,1]}$.

- **Challenge:** describe the transition between the i.i.d. and smooth cases.
- We initialize the weights as increments of a fractional Brownian motion $(B_t^H)_{t \in [0,1]}$.
- **>** Recall: B^H is Gaussian, starts at zero, has zero expectation, and covariance function

$$\mathbb{E}(B_s^H B_t^H) = \frac{1}{2}(|s|^{2H} + |t|^{2H} - |t - s|^{2H}), \quad 0 \le s, t \le 1.$$

- Challenge: describe the transition between the i.i.d. and smooth cases.
- We initialize the weights as increments of a fractional Brownian motion $(B_t^H)_{t \in [0,1]}$.
- **>** Recall: B^H is Gaussian, starts at zero, has zero expectation, and covariance function

$$\mathbb{E}(B_s^H B_t^H) = \frac{1}{2}(|s|^{2H} + |t|^{2H} - |t - s|^{2H}), \quad 0 \le s, t \le 1.$$

The Hurst index $H \in (0,1)$ describes the raggedness of the process.





 \triangleright H = 1/2: standard Brownian motion (SDE regime).



 \triangleright H = 1/2: standard Brownian motion (SDE regime).

- \triangleright H < 1/2: the increments are negatively correlated.
- \triangleright H > 1/2: the increments are positively correlated.



 \triangleright H = 1/2: standard Brownian motion (SDE regime).

- \triangleright H < 1/2: the increments are negatively correlated.
- \triangleright H > 1/2: the increments are positively correlated.
- \triangleright When $H \rightarrow 1$: the trajectories converge to linear functions (ODE regime).

A continuum of intermediate regularities



A continuum of intermediate regularities





Learning with ResNets

Scaling deep ResNets

Scaling in the continuous-time setting

Beyond initialization



l.i.d. initialization, $\beta = 1/2$



Smooth initialization, $\beta = 1$



I.i.d. initialization, $\beta = 1$



> The weights after training still exhibit a strong structure as functions of the layer.



l.i.d. initialization, $\beta = 1$

The weights after training still exhibit a strong structure as functions of the layer.
Their regularity is influenced by both the initialization and the choice of β.

Performance after training





(b) On CIFAR-10



> Deep limits allow to understand scaling and initialization strategies for ResNets.



- > Deep limits allow to understand scaling and initialization strategies for ResNets.
- With standard initialization the correct scaling is $\beta = 1/2$.

- > Deep limits allow to understand scaling and initialization strategies for ResNets.
- With standard initialization the correct scaling is $\beta = 1/2$.
- > To train very deep ResNets, it is important to adapt scaling to the weight regularity.

- > Deep limits allow to understand scaling and initialization strategies for ResNets.
- With standard initialization the correct scaling is $\beta = 1/2$.
- > To train very deep ResNets, it is important to adapt scaling to the weight regularity.
- Perspectives: what about training? how should we choose the regularity for a given problem?

- > Deep limits allow to understand scaling and initialization strategies for ResNets.
- With standard initialization the correct scaling is $\beta = 1/2$.
- > To train very deep ResNets, it is important to adapt scaling to the weight regularity.
- Perspectives: what about training? how should we choose the regularity for a given problem?
- **>** To know more: arXiv:2206.06929.

Thank you!



adeline.fermanian@mines-paristech.fr



https://afermanian.github.io