Representations in Random Deep Neural Networks

Hadi Daneshmand, Amir Joudaki, Francis Bach

INRIA Paris, Sierra Team

Oct. 2021

Ínnía-



"The noisy brain"



- Microscopically neurons are random but they collectively resemble interesting stochastic processes
- Mean-field analysis of underlying stochastic processes in the brain sheds light into
 - decision-making
 - attention
 - brain dysfunctions such as schizophrenia





Inría 2

The Noisy Brain Sochastic Dynamics as a Principle of Brain Function Edmund T. Rolls & Gustavo Deco

Evolution of computational neural networks

3



Figure: Images credit: Alfredo Canziani et. al.

Noisy computational networks

- Historical neural networks do not work well with random weights
- Modern computational neural networks perform surprisingly well with random weights if the neurons are wired well together¹.

Why?

¹ Frankle, J., Schwab, D. J. & Morcos, A. S. Training batchnorm and only batchnorm: On the expressive power of random features in CNNs. *ICLR* (2021).
Hadi Daneshmand, Amir Joudaki, Francis Bach | Batch Normalization Orthogonalizes Representations in Deep Random Networks

Historical single-layer MLPs

$$H_1 = rac{1}{\sqrt{ ext{batchsize}}}F(W_0H_0)$$

- ► $H_0 \in \mathbb{R}^{\text{width} \times \text{batchsize}}$ is a deterministic matrix
- $W_0 \in \mathbb{R}^{\text{width} \times \text{width}}$ is a random Gaussian matrix
- Empirical eigenvalue distribution (e.e.d): $\frac{1}{\text{width}} \sum_{i=1}^{\text{width}} \delta(\lambda_i(H_1^\top H_1))$
- Given the first two moments of H₁,² characterizes e.e.d. of H₁ as batchsize and width tends to ∞ denoted by p

²Louart, C., Liao, Z., Couillet, R., *et al.* A random matrix approach to neural networks. *The Annals of Applied Probability* (2018).

Historical single-layer MLPs

Stieltjes transformation of density p on interval I:

$$S_p(z) = \int_I rac{p(t)dt}{z-t}, \quad z \in C \setminus I$$

• Given $G = \mathbf{E} \left[H_1^\top H_1 \right]$, the following holds³

$$S_p(z) = rac{1}{ ext{batchsize}} \operatorname{Tr}\left(\underbrace{rac{ ext{width}}{ ext{batchsize}} rac{G}{1+s(z)} - zI}_{M(s)}
ight)^{-1}$$

s(z) is the solution of

$$s(z) = rac{1}{ ext{batchsize}} \operatorname{Tr} \left(GM^{-1}(s(z)) \right)$$

³Louart, C., Liao, Z., Couillet, R., et al. A random matrix approach to neural networks. The Annals of Applied Probability (2018).
Hadi Daneshmand, Amir Joudaki, Francis Bach | Batch Normalization Orthogonalizes Representations in Deep Random Networks

Historical random deep networks

- ▶ Let $x_{\ell} \in \mathbb{R}^{\text{width}}$ be representation of input x_0 at layer ℓ
- The representations make a Markov chain as:

$$x_{\ell+1} = \frac{W_\ell x_\ell}{\|W_\ell x_\ell\|}$$

Suppose the elements of $W_{\ell} \in \mathbb{R}^{\text{width} \times \text{width}}$ are i.i.d. Gaussian.

Coupled representations

$$x_{\ell+1} = \frac{W_{\ell} x_{\ell+1}}{\|W_{\ell} x_{\ell+1}\|}, \qquad y_{\ell+1} = \frac{W_{\ell} y_{\ell+1}}{\|W_{\ell} y_{\ell+1}\|}$$

Inría 8

Coupled representations

$$x_{\ell+1} = \frac{W_{\ell} x_{\ell+1}}{\|W_{\ell} x_{\ell+1}\|}, \qquad y_{\ell+1} = \frac{W_{\ell} y_{\ell+1}}{\|W_{\ell} y_{\ell+1}\|}$$

Inría (8

The chains contracts to a random directions independent from the starting state.

Product of random matrices

- Consider the product of Gaussian matrix as $S_{\ell} = W_{\ell} \dots W_1$
- ► Claim: $S_{\ell}/||S_{\ell}||$ becomes rank one in limit.
- ▶ Therefore, $(S_{\ell}x)/||S_{\ell}||$ becomes independent from *x* as $\ell \to \infty$.

Inría

Lyapunov exponents

Definition:

$$\lambda_k = \lim_{\ell \to \infty} \frac{1}{2\ell} \log \left(k^{\text{th}} \text{ largest eigenvalue of } S_\ell^\top S_\ell \right)$$

10

Inría

• Computation⁴: $\lambda_k = \frac{1}{2} (\log(2) + \Psi(\frac{d-k+1}{2}))$



▶ $\lambda_1 - \lambda_2 < 0$ implies $S_{\ell} / ||S_{\ell}||$ becomes rank one in limit.

⁴Newman, C. M. The distribution of Lyapunov exponents: Exact results for random matrices. *Communications in mathematical physics* (1986).
Hadi Daneshmand, Amir Joudaki, Francis Bach, I Batch Normalization Orthogonalizes Representations in Deep Random Networks

Modern NN with Batch normalization (BN)

BN is one of the main building block of modern neural networks⁵

Representation H_{ℓ} :width \times batchsize.

$$H_{\ell+1} = F(BN_{\alpha,\beta}(W_{\ell}H_{\ell}))$$
(1)

(nría_

 $\blacktriangleright BN_{\alpha,\beta}: \mathbb{R}^{\textit{width} \times \textit{batchsize}} \rightarrow \mathbb{R}^{\textit{width} \times \textit{batchsize}}$

 $[BN_{\alpha,\beta}(M)]_{:} = \alpha_i \text{centered}(M_{i:}) + \beta_i$

Learning only parameters α and β (per unit) leads to surprisingly good performance⁶

⁵loffe, S. & Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. in *ICML* (2015).

⁶Frankle, J., Schwab, D. J. & Morcos, A. S. Training batchnorm and only batchnorm: On the expressive power of random features in CNNs. *ICLR* (2021).

The Markov chain of representations

(nría_

We study the following Markov chain of matrices⁷⁸.

- BN(M) normalizes M row-wise
- Representations:

$$H_{\ell+1} = \left(rac{1}{\sqrt{\mathit{width}}}
ight) \mathit{BN}(W_\ell H_\ell)$$

W_ℓ: (*width* × *width*) with Gaussian elements

⁷Daneshmand, H., Joudaki, A. & Bach, F. Batch Normalization Orthogonalizes Representations in Deep Random Networks. *NeurIPS21*. Spotlight presentation (among top %3 of submissions).

⁸Daneshmand, H. *et al.* Batch Normalization Provably Avoids Rank Collapse for Randomly Initialised Deep Networks. *NeurIPS20.*

Theoretical results

► **E**
$$\left[\text{orthogonality gap}(H_{\ell}) \right] = \mathcal{O}\left((1 - \alpha)^{\ell} + \frac{\text{batchsize}}{\alpha \sqrt{\text{width}}} \right)$$

► Wasser.₂($W_{\ell}H_{\ell}$, Gaussian)² = $\mathcal{O}\left((1 - \alpha)^{\ell} (\text{batchsize}) + \frac{(\text{batchsize})^{2}}{\alpha \sqrt{\text{width}}} \right)$

Inría (13)



Orthogonalization



Inría (14

- α is the minimum of smallest singular value of $\{H_1, \ldots, H_\ell\}$.
- To get a non-vacuous bound, we need an α independent from ℓ .

Modern NN vs. historical NN 15 (nría_ BN Without BN

BN Without BN $\mathbf{E}\left[\text{orth. gap}(H_{\infty})\right] = \mathcal{O}\left(\frac{\text{batch size}}{\alpha\sqrt{\text{width}}}\right)$ **E** orth. gap (H'_{∞}) = $\Theta(1)$

Inría (15

Modern NN vs. historical NN

The orthogonality influences training



16

(nría

⁹Lubana, E. S., Dick, R. P. & Tanaka, H. Beyond BatchNorm: Towards a General Understanding of Normalization in Deep Learning. *arXiv preprint arXiv:2106.05956* (2021).

Replacing BN with orthogonalization

Saving training time by starting from orthogonal representations

17

Inría



MLPs with ReLU and **without BN** for classifying CIFAR-10 Red: standard initialization with low orthogonality gaps Blue: novel initialization ensuring orthogonal representations

Gaussian approximation

Wasserstein₂(
$$W_{\ell}H_{\ell}$$
, Gaussian)² = $\mathcal{O}\left((1-\alpha)^{\ell} \left(\text{batchsize}\right) + \frac{(\text{batchsize})^2}{\alpha\sqrt{\text{width}}}\right)$

Inría (18)













Deepening our knowledge about representations in deep neural networks will allow us to design more efficient neural networks.