

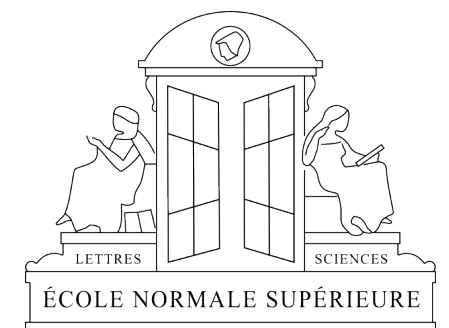
The Gaussian Equivalence of generative models for learning with two-layer neural networks

Sebastian Goldt (ENS Paris)

July 3rd, 2020 • Youth in high dimensions

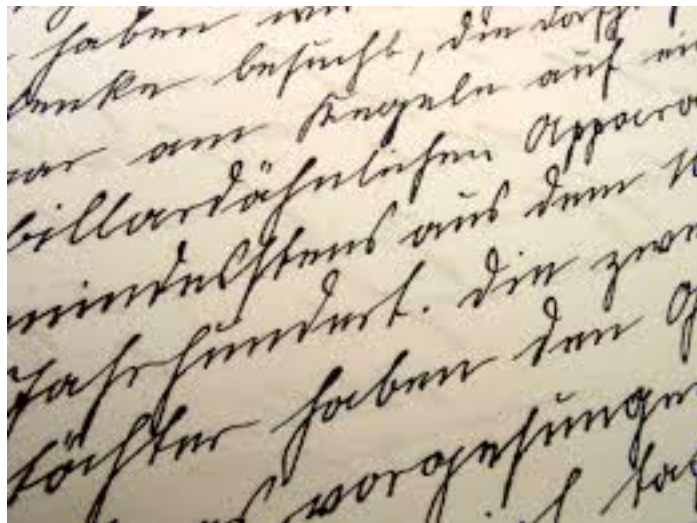
Joint work w/

F. Krzakala (ENS), M. Mézard (ENS),
G. Reeves (Duke), L. Zdeborová (CEA Saclay)



The problem of **data structure**

The data sets we care about in machine learning contain a lot of structure.



Written text (NLP)



Images

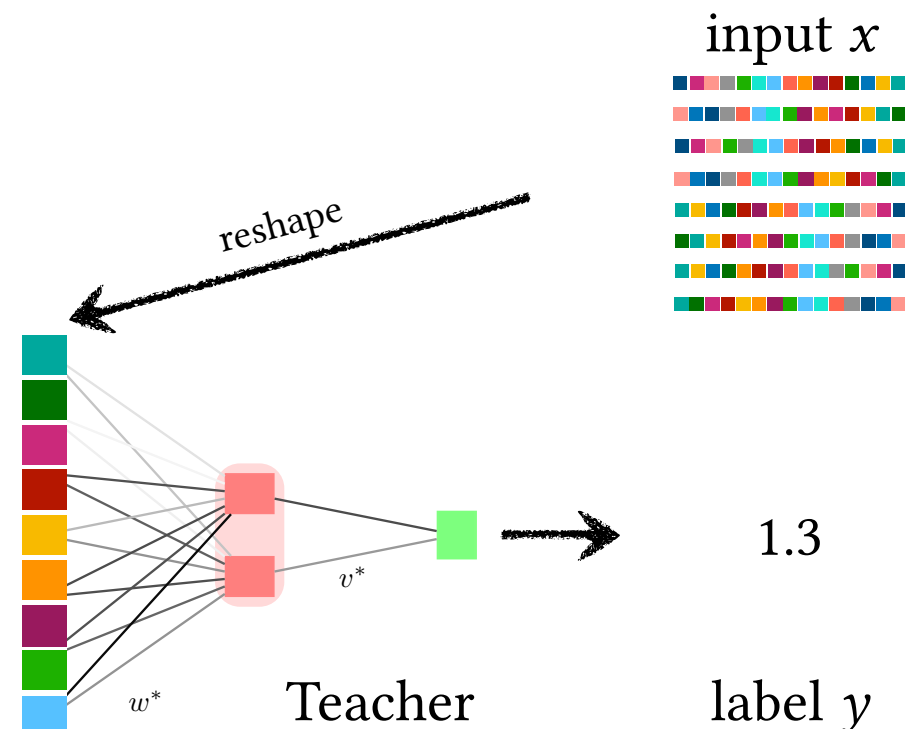


Games of Go

How does data structure impact learning in neural networks?

The vanilla teacher-student setup

Supervised regression task with samples $(x \in \mathbb{R}^N, y \in \mathbb{R})$
drawn from some distribution $q(x, y)$



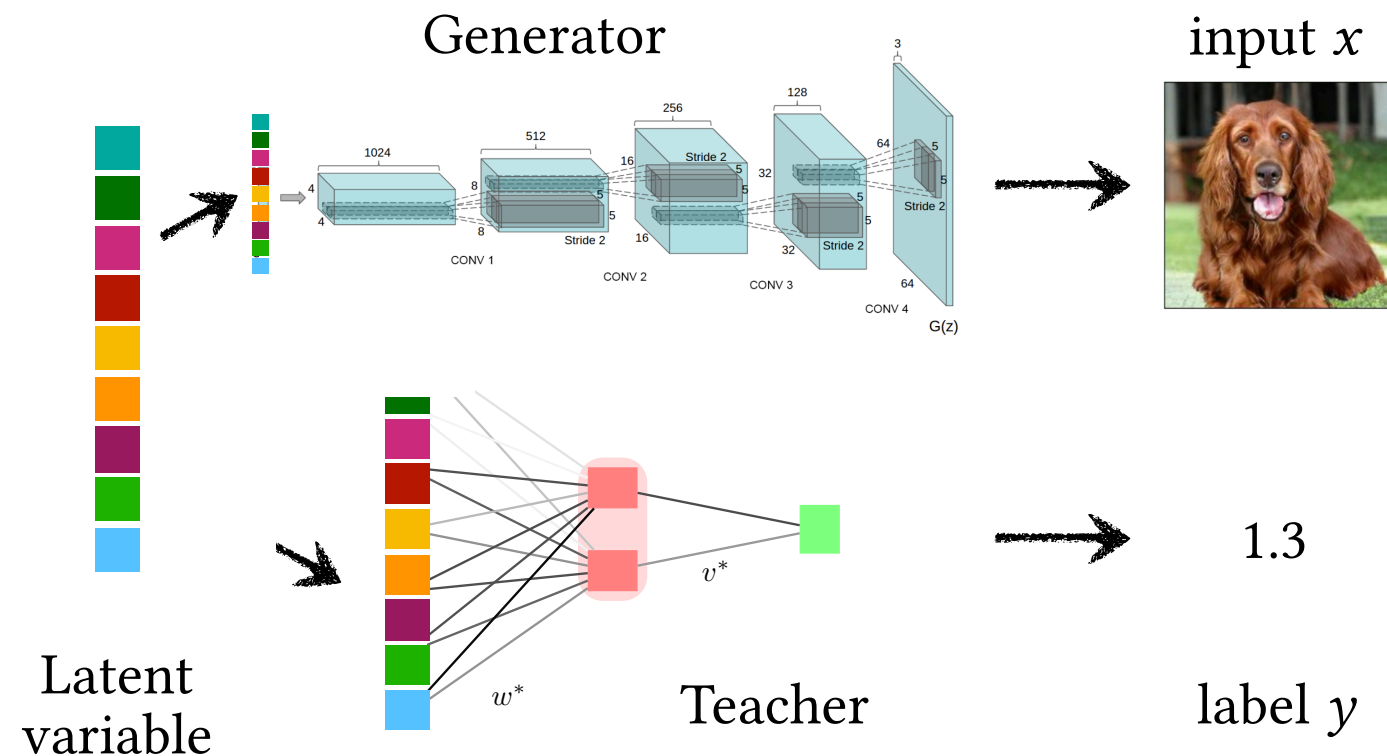
Inputs are i.i.d. draws
from the normal distribution:

$$x = (x_i), \quad x_i \sim \mathcal{N}(0, 1)$$

Generate label by applying a
teacher network directly to the input:

$$y = \phi_{\tilde{\theta}}(x) = \sum_{m=1}^M v^m \tilde{g}(\tilde{w}^m x)$$

The hidden manifold model



Generate an
i.i.d. latent variable

Generate the input
using a generative network

Apply a teacher to the
latent representation

$$c = (c_r), \quad c_r \sim \mathcal{N}(0, 1)$$

$$x = \mathcal{G}(c)$$

$$y = \phi_{\tilde{\theta}}(c)$$

Gaussian Equivalence Theorem

We give rigorous conditions under which we can analyse learning from data coming from single-layer generators.

Dynamical equations for two-layer students

The equations track the test error of two-layer students trained on deep generative models.

The Gaussian Equivalence Property

Goal: compute the prediction mean-squared error at all times.

For the **vanilla-teacher student** with i.i.d. inputs x :

Saad & Solla, (1995)
Biehl & Schwarze (1995)

$$\text{pmse}(\theta, \tilde{\theta}) = \mathbb{E}_x \left(\sum_{k=1}^K v^k g(w^k x) - \sum_{m=1}^M \tilde{v}^m \tilde{g}(\tilde{w}^m x) \right)^2$$

*Student network
(trying to learn)*

*Teacher network
(creates the data)*

The Gaussian Equivalence Property

Goal: compute the prediction mean-squared error at all times.

For the **vanilla-teacher student** with i.i.d. inputs x :

Saad & Solla, (1995)
Biehl & Schwarze (1995)

$$\text{pmse}(\theta, \tilde{\theta}) = \mathbb{E}_x \left(\sum_{k=1}^K v^k g(w^k x) - \sum_{m=1}^M \tilde{v}^m \tilde{g}(\tilde{w}^m x) \right)^2$$

Average over
the inputs x

$\lambda^k \sim w^k x$

$\nu^m \sim \tilde{w}^m x$

The Gaussian Equivalence Property

Goal: compute the prediction mean-squared error at all times.

For the **vanilla-teacher student** with i.i.d. inputs x :

Saad & Solla, (1995)
Biehl & Schwarze (1995)

$$\text{pmse}(\theta, \tilde{\theta}) = \mathbb{E}_{\lambda, \nu} \left(\sum_{k=1}^K v^k g(\lambda^k) - \sum_{m=1}^M \tilde{v}^m \tilde{g}(\nu^m) \right)^2$$

Average over
the *local fields* (λ, ν)

$\lambda^k \sim w^k x$

$\nu^m \sim \tilde{w}^m x$

Key random variables
for online learning
and replicas (batch)

The Gaussian Equivalence Property

Goal: compute the prediction mean-squared error at all times.

For the **vanilla-teacher student** with i.i.d. inputs x :

Saad & Solla, (1995)
Biehl & Schwarze (1995)

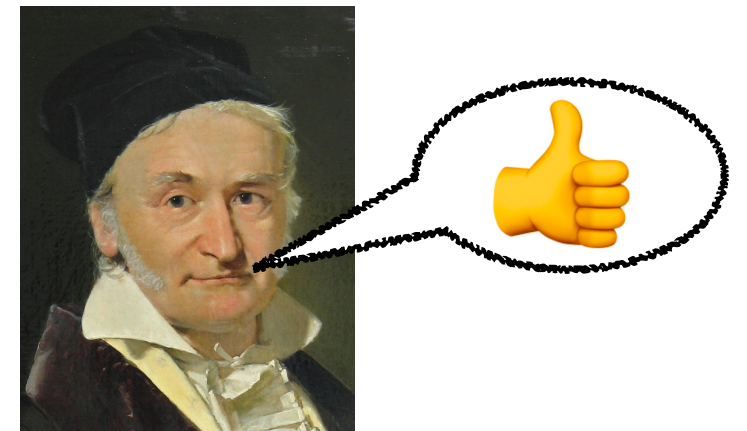
$$\text{pmse}(\theta, \tilde{\theta}) = \mathbb{E}_{\lambda, \nu} \left(\sum_{k=1}^K v^k g(\lambda^k) - \sum_{m=1}^M \tilde{v}^m \tilde{g}(\nu^m) \right)^2$$

$$\mathbb{E} x_i x_j = \delta_{ij} \quad \begin{array}{l} \boxed{\lambda^k} \sim \sum_i w_i^k x_i \\ \boxed{\nu^m} \sim \sum_i \tilde{w}_i^m x_i \end{array}$$

Gaussian Equivalence Property:
 (λ, ν) are jointly Gaussian

Hence, the *pmse* is a function of only the second moments of (λ, ν) :

$$Q^{k\ell} \equiv \mathbb{E} \lambda^k \lambda^\ell, \quad R^{km} \equiv \mathbb{E} \lambda^k \nu^m, \quad T^{mn} \equiv \mathbb{E} \nu^m \nu^n$$



The Gaussian Equivalence Theorem

Setup: Fully connected, single layer generator $\mathcal{G} : \mathbb{R}^D \rightarrow \mathbb{R}^N$

$$x_n = \mathcal{G}_n(c) = \sigma(a_n^\top c)$$

with the teacher acting on the latent variable c : $y = \phi_{\tilde{\theta}}(c)$

$$\mathbb{E} x_i x_j = \Omega_{ij}$$

$$\lambda^k \sim \sum_i w_i^k x_i$$

$$\nu^m \sim \sum_r \tilde{w}_r^m c_r$$



They're still
(sometimes)
Gaussian!

Theorem: Let P be the distribution of the pair (λ, ν) and let \hat{P} be the Gaussian distribution with the same first and second moments. Then...

$$d_{\text{MS}}(P, \hat{P}) = O \left(\left\| \frac{1}{\sqrt{N}} W M_1^{1/2} \right\|^2 + \left\| \frac{1}{\sqrt{N}} W M_2^{1/2} \right\|^2 + \frac{1}{\sqrt{N}} \left\| \frac{1}{\sqrt{D}} \tilde{W} A^\top \right\|^2 + \frac{1}{\sqrt{N}} \right)$$

The Gaussian Equivalence Theorem

Theorem: Let P be the distribution of the pair (λ, ν) and let \hat{P} be the Gaussian distribution with the same first and second moments. Then...

$$\mathcal{G} : \mathbb{R}^D \rightarrow \mathbb{R}^N$$

$$x_n = \mathcal{G}_n(c) = \sigma(a_n^\top c)$$

$$y = \phi_{\tilde{\theta}}(c)$$

$$d_{\text{MS}}(P, \hat{P}) = O \left(\left\| \frac{1}{\sqrt{N}} W M_1^{1/2} \right\|^2 + \left\| \frac{1}{\sqrt{N}} W M_2^{1/2} \right\|^2 + \frac{1}{\sqrt{N}} \left\| \frac{1}{\sqrt{D}} \tilde{W} A^\top \right\|^2 + \frac{1}{\sqrt{N}} \right)$$

Student weights

Teacher weights

Generator weights

Related to input correlations

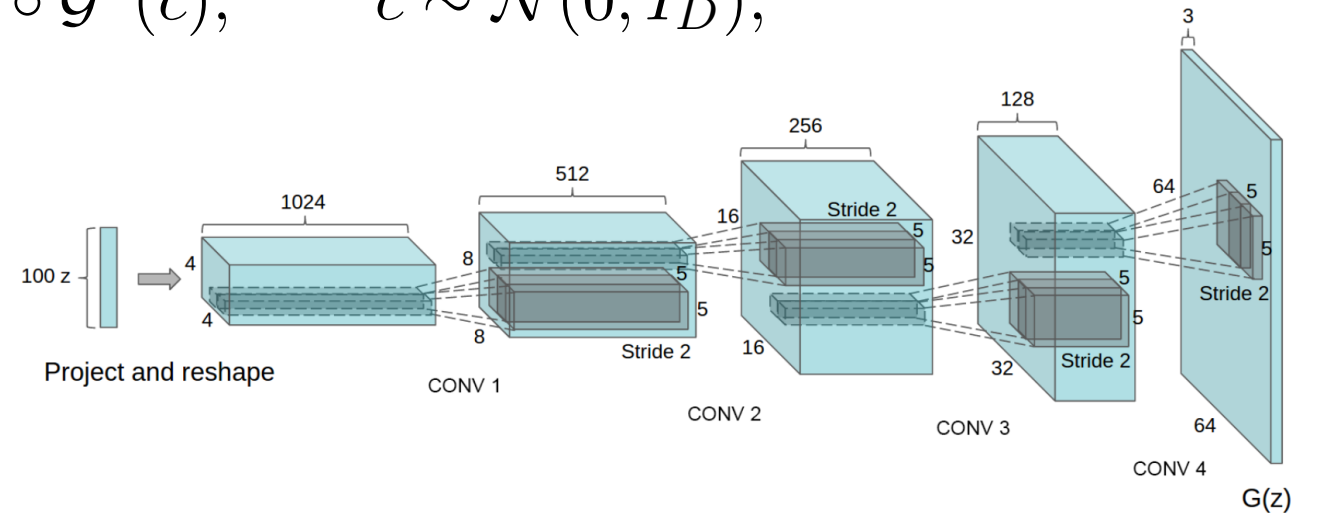
Related work

- Works in wide network limit rely on RMT and thus random weights
- Mei & Montanari; Couillet et al. introduce related equivalent Gaussian models for integrals w.r.t. spectral densities.
- Large body of work on low-dim projections of high-dim data being Gaussian - we quantify how Gaussian they look like.

The deep hidden manifold model

$$x = \mathcal{G}(c) = \mathcal{G}^L \cdots \mathcal{G}^3 \circ \mathcal{G}^2 \circ \mathcal{G}^1(c), \quad c \sim \mathcal{N}(0, I_D),$$

$$y = \phi_{\tilde{\theta}}(c)$$



- Train the student using online SGD:

*Deep Convolutional GAN
Radford et al., ICLR 2016*

$$\theta_{\mu+1} = \theta_{\mu} - \eta \nabla_{\theta} E(\theta) |_{\theta_{\mu}, x_{\mu}, y_{\mu}^*}$$

Goal: Derive a closed set of equations for the order parameters

$$Q^{k\ell} \equiv \mathbb{E} \lambda^k \lambda^{\ell}, \quad R^{km} \equiv \mathbb{E} \lambda^k \nu^m$$

that track the dynamics of a two-layer student trained using online SGD on the deep hidden manifold.

Dynamical equations for two-layer students

Statement:

$$Q^{k\ell} \equiv \mathbb{E} \lambda^k \lambda^\ell, \quad R^{km} \equiv \mathbb{E} \lambda^k \nu^m$$

$$Q^{k\ell} = \int d\mu_\Omega(\rho) \rho q^{k\ell}(\rho)$$

*Spectral density of
input-input covariance*

$$\begin{aligned} \frac{\partial q^{k\ell}(\rho)}{\partial t} = & -\eta \left(\rho \sum_{j \neq k}^K \left[v^k v^j q^{k\ell}(\rho) h_{(1)}^{kj}(Q) + v^k v^j q^{j\ell}(\rho) h_{(2)}^{kj}(Q) \right] + \rho v^k v^k q^{k\ell}(\rho) h_{(3)}^k(Q) \right. \\ & - v^k \sum_n^M \left[\rho \tilde{v}^n q^{k\ell}(\rho) h_{(4)}^{kn}(Q, R, T) + \frac{1}{\sqrt{\delta}} \tilde{v}^n r^{\ell n}(\rho) h_{(5)}^{kn}(Q, R, T) \right] \\ & \left. + \text{all of the above with } \ell \rightarrow k, k \rightarrow \ell \right) + \eta^2 \gamma v^k v^\ell h_{(6)}^{k\ell}(Q, R, T, v, \tilde{v}). \end{aligned}$$

$$R^{km} = \frac{1}{\sqrt{\delta}} \int d\mu_\Omega(\rho) r^{km}(\rho)$$

$$\begin{aligned} \frac{\partial r^{km}(\rho)}{\partial t} = & -\eta v^k \left(\rho \sum_{j \neq k}^K \left[v^j r^{km}(\rho) h_{(1)}^{kj}(Q) + v^j \rho r^{jm}(\rho) h_{(2)}^{kj}(Q) \right] + v^k \rho r^{km}(\rho) h_{(3)}^k(Q) \right. \\ & \left. - \sum_n^M \left[\rho \tilde{v}^n r^{km}(\rho) h_{(4)}^{kn}(Q, R, T) + \frac{1}{\sqrt{\delta}} \tilde{v}^n h_{(5)}^{kn}(Q, R, T) \right] \right). \end{aligned}$$

Dynamical equations for two-layer students

Statement:

$$Q^{k\ell} \equiv \mathbb{E} \lambda^k \lambda^\ell, \quad R^{km} \equiv \mathbb{E} \lambda^k \nu^m$$

$$Q^{k\ell} = \int d\mu_\Omega(\rho) \, \rho \, q^{k\ell}(\rho)$$

$$R^{km} = \frac{1}{\sqrt{\delta}} \int d\mu_\Omega(\rho) \, r^{km}(\rho)$$

Remarkably, the generator only appears via two covariance matrices:

$$\Omega_{ij} = \mathbb{E} x_i x_j$$

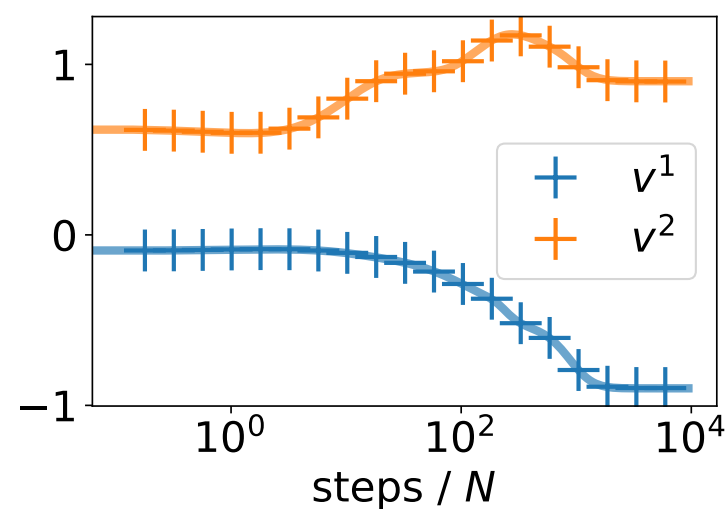
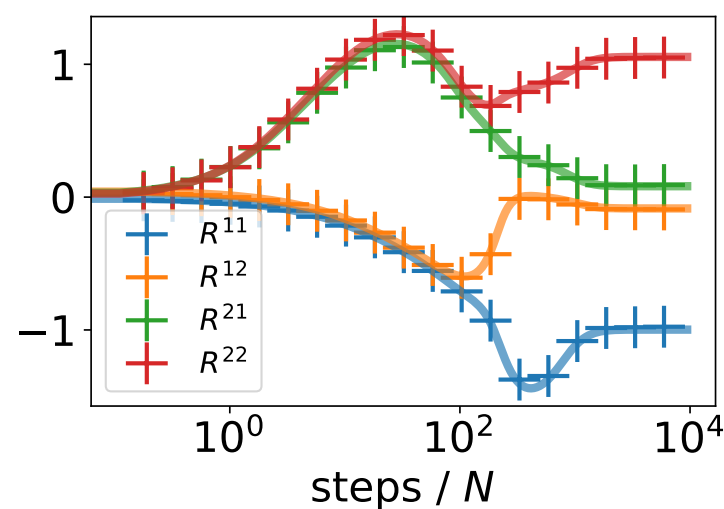
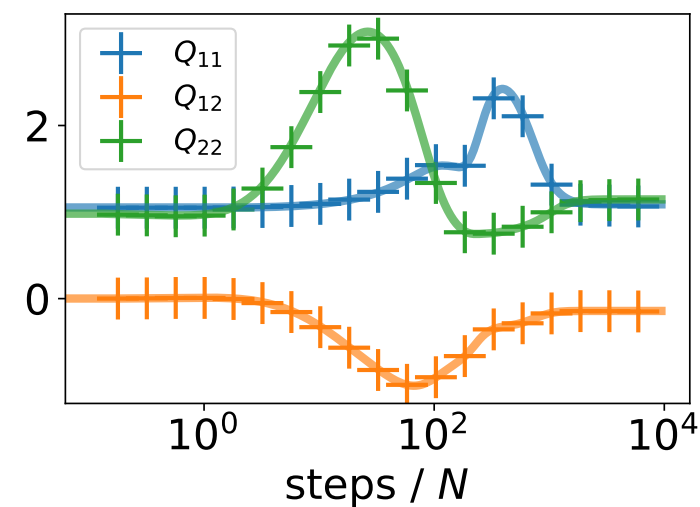
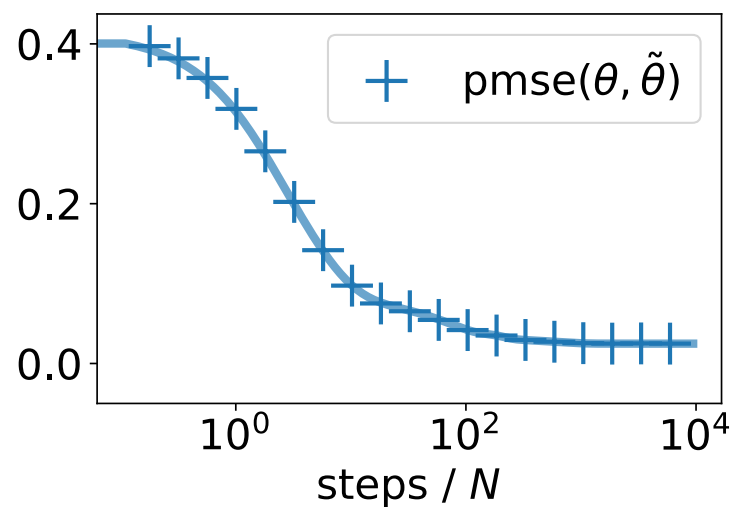
*Input-input
correlations*

$$\Phi_{ir} = \mathbb{E} x_i c_r$$

*Input-latent
correlations*

Single-layer fully connected generator

$$x_n = \mathcal{G}_n(c) = \sigma(a_n^\top c) \quad y = \phi_{\tilde{\theta}}(c)$$

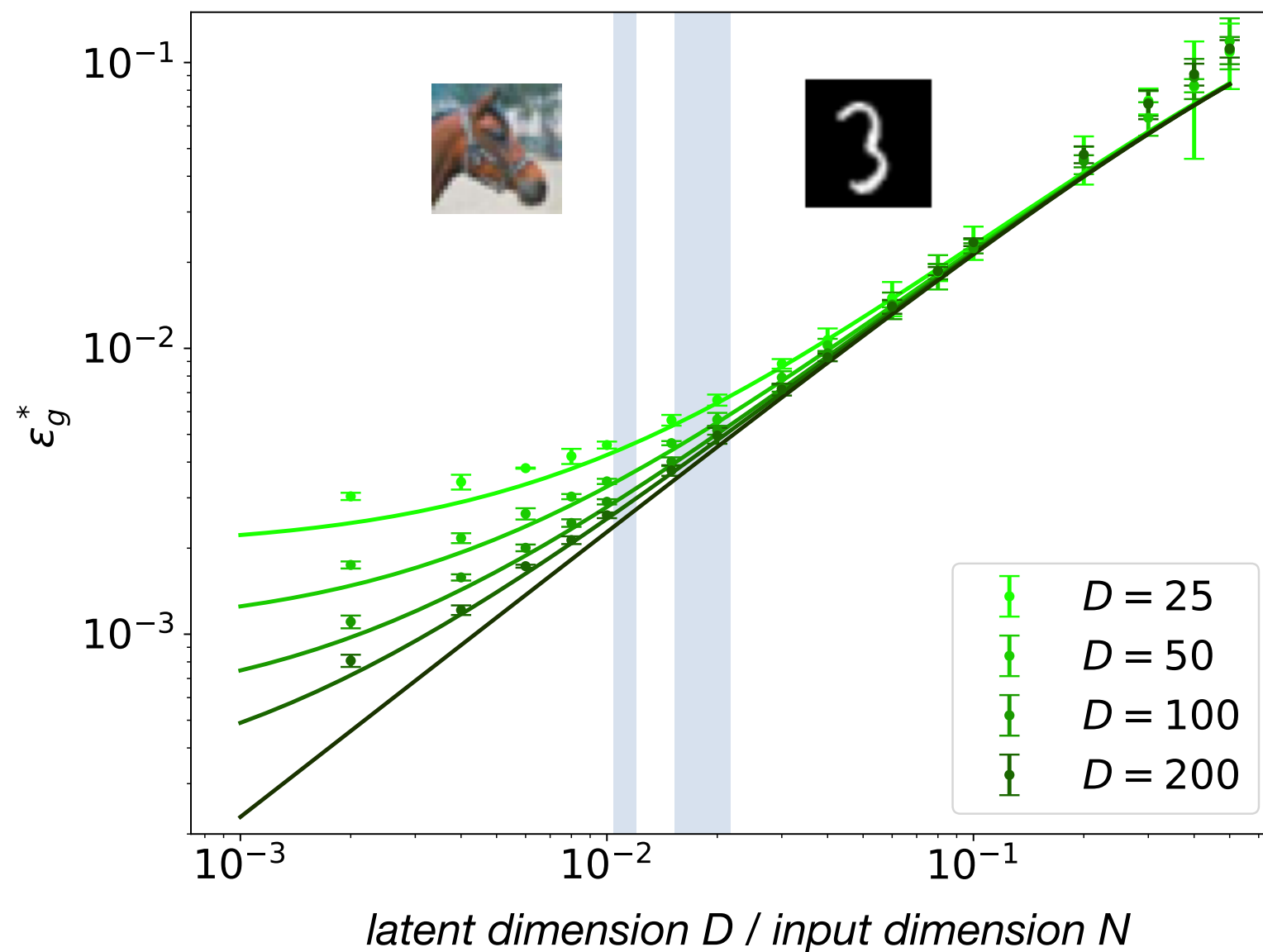


$$\sigma(x) = \text{sgn}(x), \quad g(x) = \text{erf}(x/\sqrt{2})$$

$$M=K=2, \eta = 0.2, D=800, N=8000$$

The ODEs **predict performance** of two-layer NN

$$x_n = \mathcal{G}_n(c) = \sigma(a_n^\top c) \quad y = \phi_{\tilde{\theta}}(c)$$

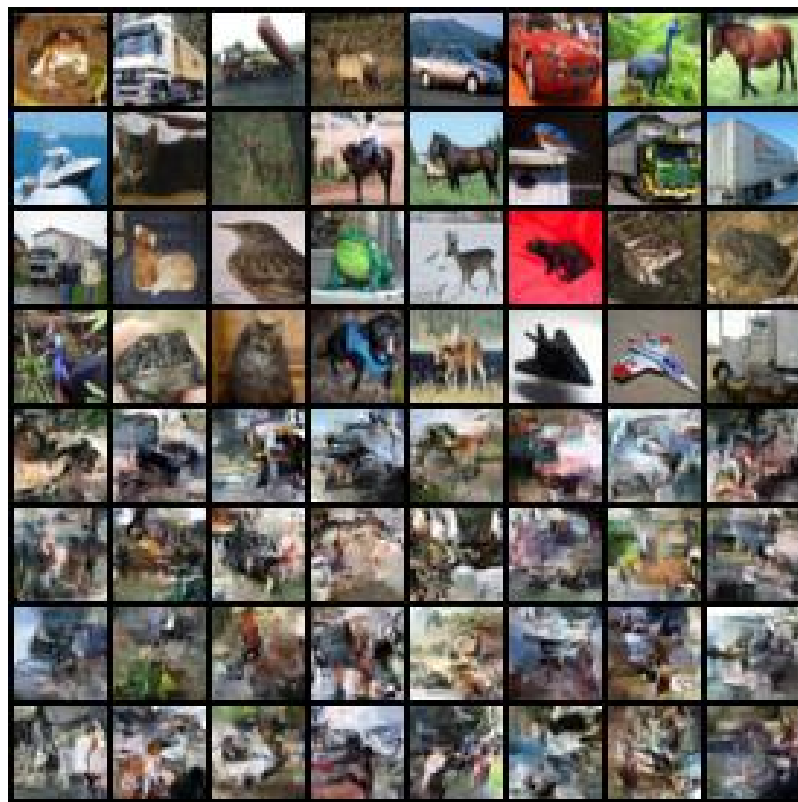


$$\sigma(x) = \text{sgn}(x), \quad g(x) = \text{erf}(x/\sqrt{2})$$
$$M=K=2, \eta = 0.2$$

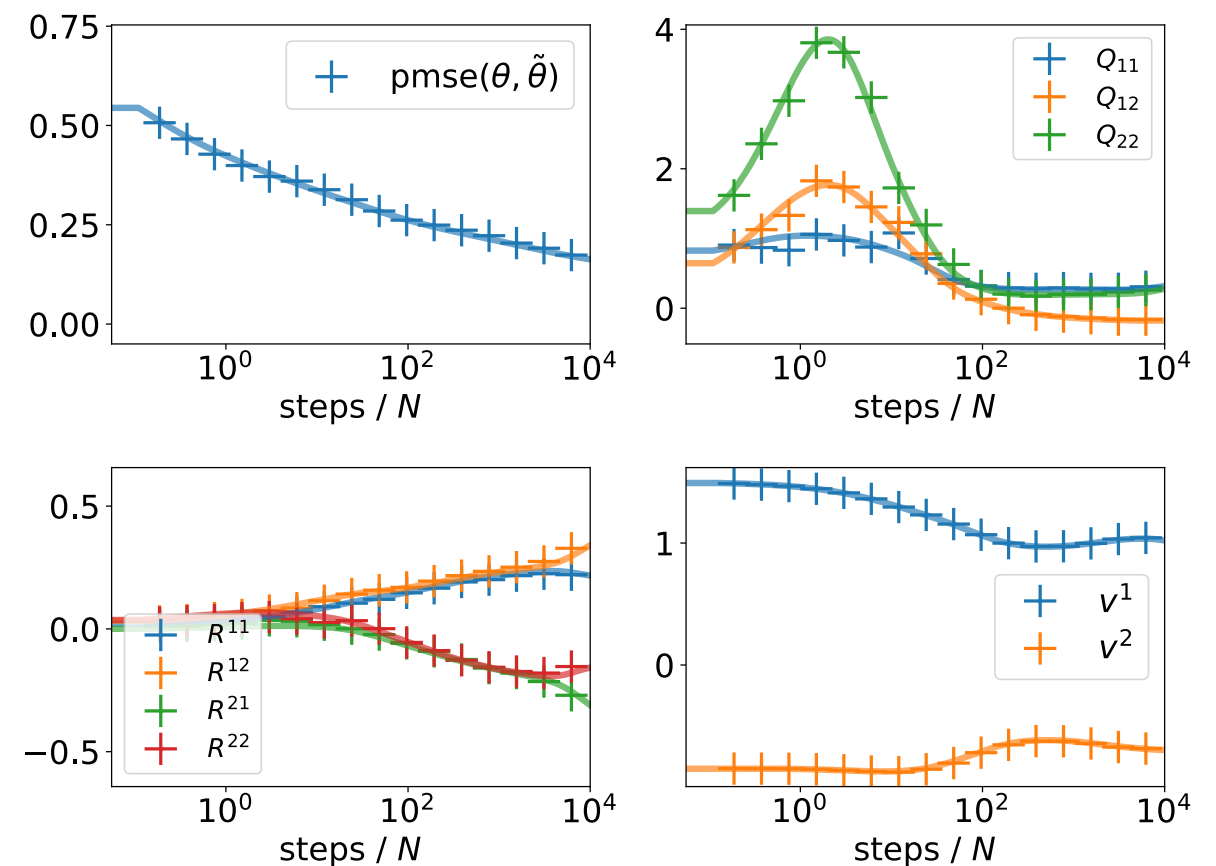
The realNVP: normalising flows

L. Dinh, J. Sohl-Dickstein, S. Bengio (ICLR 2017)

- Generate images using a series of invertible transformations, e.g. convolutions etc.
- Trained **realNVP** (>6M parameters) on CIFAR10



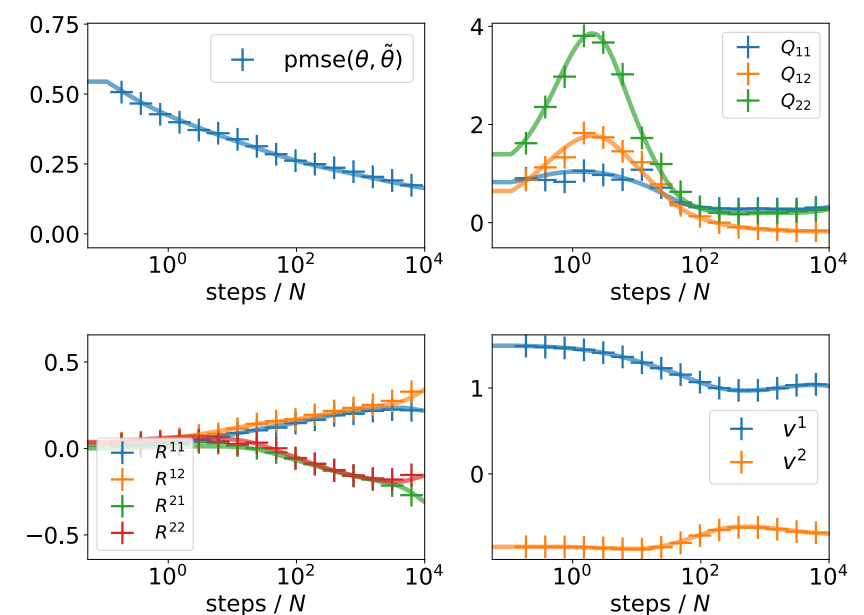
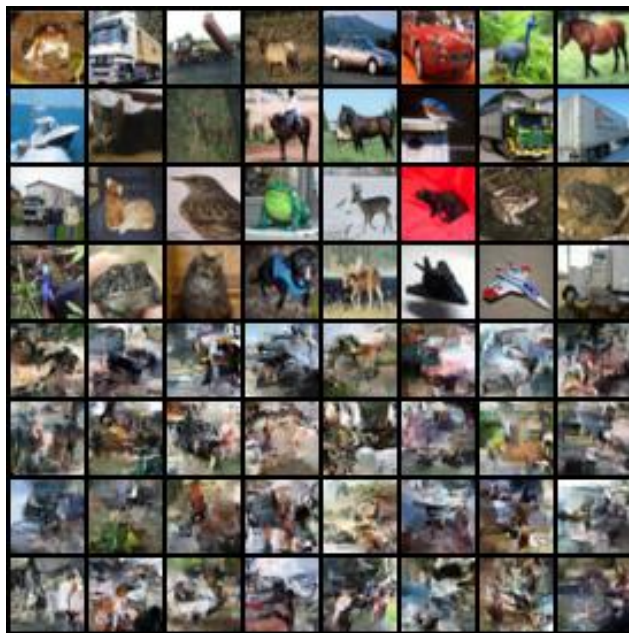
Top half: CIFAR10 images
Bottom half: Samples from realNVP
trained on CIFAR10



$M=K=2, \eta = 0.2, D=3072, N=3072$

Conclusion

- We study can study the **impact of data on learning** by analysing NN trained on data from generative models
- The GEP gives **analytical control** over the learning dynamics, and is **rigorous** for single-layer generators.



References

- **arXiv:1909.11500**
w/ M. Mézard, F. Krzakala, L. Zdeborová
- **arXiv:2006.14709**
w/ G. Reeves, M. Mézard, F. Krzakala, L. Zdeborová



<https://github.com/sgoldt>