Applied Machine Learning Days (EPFL 2022) – AI & Physics track

Statistical physics insights on stochastic gradient descent



Francesca Mignacco

Ph.D. student @ Institute of Theoretical Physics, CEA Saclay Supervised by of Lenka Zdeborová & Pierfrancesco Urbani

March 28th, 2022



For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- Why don't heavily parameterized neural networks overfit the data?
- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?



... since almost 30 years !

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

- What is the effective number of parameters?
- Why doesn't backpropagation head for a poor local minima?
- When should one stop the backpropagation and use the current parameters?

Reflections after refereeing papers for NIPS, Leo Breiman, 1995

Why don't heavily parameterized neural networks overfit the data?



... since almost 30 years !

For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

What is the effective number of parameters? Why doesn't backpropagation head for a poor local minima?

1) require the understanding of static properties

Why don't heavily parameterized neural networks overfit the data?

When should one stop the backpropagation and use the current parameters?



For instance, there are many important questions regarding neural networks which are largely unanswered. There seem to be conflicting stories regarding the following issues:

What is the effective number of parameters? Why doesn't backpropagation head for a poor local minima?

1) require the understanding of static properties

(2) require the understanding of dynamical properties

Why don't heavily parameterized neural networks overfit the data? When should one stop the backpropagation and use the current parameters?



The building blocks of deep learning

Data Structure



[Understanding deep learning is also a job for physicists, Zdeborová, L. (2020). Nature Physics, 16(6), 602-604]

Architecture

The building blocks of deep learning

Data Structure



[Understanding deep learning is also a job for physicists, Zdeborová, L. (2020). Nature Physics, 16(6), 602-604]

Architecture

Focus of this talk



The problem:

Given some training data and labels $\left\{ \mathbf{x}_{\mu}, y_{\mu} \right\}_{\mu=1}^{M}$ sampled independently from $p_{x,y}$



The problem:

Given some training data an

 $\hat{\mathbf{w}} = \operatorname{argmin} \mathscr{L}(\mathbf{w})$ W

nd labels
$$\left\{ \mathbf{x}_{\mu}, y_{\mu} \right\}_{\mu=1}^{M}$$
 sampled independently from f

Learn the predictor $\hat{y} = f_{\hat{w}}(\mathbf{x})$ by minimising the empirical risk:

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \left[\sum_{\mu=1}^{M} \ell\left(y_{\mu}, f_{\mathbf{w}}(\mathbf{x}_{\mu}) \right) + r(\mathbf{w}) \right]$$



The problem:

Given some training data an

 $\hat{\mathbf{w}} = \operatorname{argmin} \mathscr{L}(\mathbf{w})$ W

How can we characterize the **predictor's error**? $\varepsilon_{\text{gen}} = \mathbb{E} \left[\left(y_{\text{new}} - f_{\hat{\mathbf{w}}}(\mathbf{x}_{\text{new}}) \right)^2 \right]$

nd labels
$$\left\{ \mathbf{x}_{\mu}, y_{\mu} \right\}_{\mu=1}^{M}$$
 sampled independently from f

Learn the predictor $\hat{y} = f_{\hat{w}}(\mathbf{x})$ by minimising the empirical risk:

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \left[\sum_{\mu=1}^{M} \ell\left(y_{\mu}, f_{\mathbf{w}}(\mathbf{x}_{\mu}) \right) + r(\mathbf{w}) \right]$$

$$\varepsilon_{\text{train}} = \frac{1}{M} \sum_{\mu=1}^{M} \left(y_{\mu} - f_{\hat{\mathbf{w}}}(\mathbf{x}_{\mu}) \right)^2$$



Data model:

 $\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$ $y_{\mu} = \operatorname{sign}\left(\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^{*}\right)$

Data model:

$$\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$$
$$y_{\mu} = \operatorname{sign}\left(\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^{*}\right)$$

Teacher-student model:

The labels are generated from a **teacher vector**.



4



Data model:

$$\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$$
$$y_{\mu} = \operatorname{sign}\left(\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^{*}\right)$$

Teacher-student model:

The labels are generated from a **teacher vector**.

[Abbara, Aubin, Krzakala, Zdeborová, MSML 2020] [Aubin, Krzakala, Lu, Zdeborová, *NeurIPS* 2020]





4



$\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$ Data model: $y_{\mu} = \operatorname{sign}\left(\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^{*}\right)$

Rademacher bound for function class: $f_{\mathbf{w}}(\mathbf{x}) = \operatorname{sign}\left(\mathbf{x}^{\mathsf{T}}\mathbf{w}\right)$

[Abbara, Aubin, Krzakala, Zdeborová, *MSML* 2020] [Aubin, Krzakala, Lu, Zdeborová, *NeurIPS* 2020]



4



$\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$ Data model: $y_{\mu} = \operatorname{sign}\left(\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^{*}\right)$

Rademacher bound for function class: $f_{\mathbf{w}}(\mathbf{x}) = \operatorname{sign}\left(\mathbf{x}^{\mathsf{T}}\mathbf{w}\right)$

Logistic regression

Bayes-optimal performance

> [Abbara, Aubin, Krzakala, Zdeborová, MSML 2020] [Aubin, Krzakala, Lu, Zdeborová, *NeurIPS* 2020]



Data Structure



Architecture



Data Structure





















W

WEIGHTS UPDATE: $\mathbf{w}^{t+1} = \mathbf{w}^t - dt \left[\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C} \left(y_{\mu}, f_{\mathbf{w}(t)}(\mathbf{x}_{\mu}) \right) + \nabla_{\mathbf{w}} r \left(\mathbf{w}(t) \right) \right]$





WEIGHTS UPDATE: $\mathbf{w}^{t+1} = \mathbf{w}^t - dt \left[\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C} \left(y_{\mu}, f_{\mathbf{w}(t)}(\mathbf{x}_{\mu}) \right) + \nabla_{\mathbf{w}} r \left(\mathbf{w}(t) \right) \right]$





WEIGHTS UPDATE: $\mathbf{w}^{t+1} = \mathbf{w}^t - dt \left[\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C} \left(y_{\mu}, f_{\mathbf{w}(t)}(\mathbf{x}_{\mu}) \right) + \nabla_{\mathbf{w}} r \left(\mathbf{w}(t) \right) \right]$



t=2W



WEIGHTS UPDATE: $\mathbf{w}^{t+1} = \mathbf{w}^t - dt \left[\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C} \left(y_{\mu}, f_{\mathbf{w}(t)}(\mathbf{x}_{\mu}) \right) + \nabla_{\mathbf{w}} r \left(\mathbf{w}(t) \right) \right]$



Bad Global Minima Exist and SGD Can Reach Them

Shengchao Liu, Dimitris Papailiopoulos University of Wisconsin–Madison

t=2

Dimitris Achlioptas University of California, Santa Cruz

[Liu, Papailiopoulos, Achlioptas, D. (2019), arXiv:1906.02613]



- WHICH BASINS ARE ATTRACTIVE ?
- $\bullet~When/Why$ is SGD better than GD ?

• What is the role of SGD effective noise ?

Related works (non-exhaustive list)

$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\tilde{\nabla}_{R}\mathscr{L}(\mathbf{w})$



Related works (non-exhaustive list)

$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\tilde{\nabla}_{R}\mathscr{L}(\mathbf{w})$





Related works (non-exhaustive list)

$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\tilde{\nabla}_B \mathscr{L}(\mathbf{w})$





 $\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\tilde{\nabla}_B \mathscr{L}(\mathbf{w})$ **APPROXIMATE** WEIGHTS GRADIENT

Related works (non-exhaustive list)



= computed only on the training samples in the current mini-batch

Related works (non-exhaustive list)

• Modelling the <u>noise</u> of SGD:

$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t \,\nabla \mathcal{L}(\mathbf{w}) + \mathrm{d}t \left(\nabla \mathcal{L}(\mathbf{w}) - \tilde{\nabla}_{\mathcal{B}} \mathcal{L}(\mathbf{w}) \right)$

Related works (non-exhaustive list)

• Modelling the <u>noise</u> of SGD:

 $\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t \,\nabla \mathscr{L}(\mathbf{w}) + \mathrm{d}t \left(\nabla \mathscr{L}(\mathbf{w}) - \tilde{\nabla}_{\mathscr{B}} \mathscr{L}(\mathbf{w}) \right)$

Gradient descent

SGD noise

Related works (non-exhaustive list)

• Modelling the <u>noise</u> of SGD:

$$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\nabla\mathcal{L}$$

Gradient descent

• Gaussian (Central limit theorem) [Mandt et al., '16; Jastrzebski et al., '17; Li et al., '17; Hu et al., '17; Zhu et al., '18; Chaudhari, Soatto, '18, ...]

 $\mathscr{C}(\mathbf{w}) + \mathrm{d}t \left(\nabla \mathscr{L}(\mathbf{w}) - \widetilde{\nabla}_{\mathscr{B}} \mathscr{L}(\mathbf{w}) \right)$ SGD noise

Related works (non-exhaustive list)

• Modelling the <u>noise</u> of SGD:

$$\mathbf{w}(t + \mathrm{d}t) = \mathbf{w}(t) - \mathrm{d}t\,\nabla\mathcal{L}$$

Gradient descent

- Gaussian (Central limit theorem)
- Fat-tailed (Generalized CLT) [Simsekli et al., '19; Gurbuzbalaban et al., '20; Simsekli et al., '20]

 $\mathscr{C}(\mathbf{w}) + \mathrm{d}t \left(\nabla \mathscr{L}(\mathbf{w}) - \tilde{\nabla}_{\mathscr{B}} \mathscr{L}(\mathbf{w}) \right)$ SGD noise

[Mandt et al., '16; Jastrzebski et al., '17; Li et al., '17; Hu et al., '17; Zhu et al., '18; Chaudhari, Soatto, '18, ...]

Related works (non-exhaustive list)

- Tracking the <u>whole learning trajectory</u>:
 - Gradient descent in (deep) linear networks [Bös, Opper, '97; Saxe, McClelland, Ganguli, '13; Advani, Saxe, Sompolinksy, '20]

stic Modelling the dynamics of **mini-batch** gradient descent (SGD) stochastic

- Tracking the <u>whole learning trajectory</u>:
 - Gradient descent in (deep) linear networks
 - Online SGD

i.e., a <u>new sample</u> is used to update the weights at each time step no distinction between train & test





stochastic Modelling the dynamics of <u>mini-batch-g</u>radient descent (SGD)

Related works (non-exhaustive list)

- Tracking the whole learning trajectory:
 - Gradient descent in (deep) linear networks
 - Online SGD for two-layer neural networks



stochastic Modelling the dynamics of <u>mini-batch-g</u>radient descent (SGD)

Related works (non-exhaustive list)

• Tracking the whole learning trajectory:

- Gradient descent in (deep) linear networks
- Online SGD for two-layer neural networks

-with finite hidden layer

[in physics: Saad, Solla, '95; Saad, '09; Goldt, Advani, Saxe, Krzakala, Zdeborová, '19; Goldt, Mézard, Krzakala, Zdeborová, '19; Refinetti, et al., '20, '21]



stochastic Modelling the dynamics of <u>mini-batch-g</u>radient descent (SGD)

Related works (non-exhaustive list)

• Tracking the whole learning trajectory:

- Gradient descent in (deep) linear networks
- Online SGD for two-layer neural networks
- -with finite hidden layer

with infinitely wide hidden layer
[Rotskoff, Vanden-Eijnden, '18; Mei, Montanari, Nguyen, '18; Chizat, Bach, '18; ...]



- Tracking the <u>whole learning trajectory</u>:
 - Multi pass SGD in shallow neural networks
 - Bordelon, Pehlevan, '21]

[FM, Krzakala, Urbani, Zdeborová, '20; FM, Urbani, Zdeborová, '21; FM, Urbani, '21;

• Tracking the <u>whole learning trajectory</u>:

• Multi - pass SGD in shallow neural networks

[FM, Krzakala, Urbani, Zdeborová, '20; FM, Urbani, Zdeborová, '21; FM, Urbani, '21; Bordelon, Pehlevan, '21]

KEY INGREDIENTS: Generative model for the data + Typ

Generative model for the data + Typical-case scenario + Thermodynamic limit (#samples, #dimensions $\rightarrow \infty$)

Based on:FM, Krzakala, Urbani, Zdeborová, NeurIPS 2020 ;FM, Urbani, Zdeborová, Machine Learning: Science & Technology (2020) ;FM, Urbani, arXiv preprint: arXiv:2112.10852

Joint works with:



Pierfrancesco Urbani @IPhT



Lenka Zdeborová @EPFL



Florent Krzakala @EPFL



- Sampling protocols for the mini batch $\mathcal{B}(t)$
 - $\mu \in \mathscr{B}(t)$ with probability $b \in (0,1]$ i.i.d. at each time
 - $\frac{1}{M} \langle |\mathscr{B}(t)| \rangle = b \quad \text{"BATCH SIZE"}$



Persistent SGD: two-state Markov jump process $\frac{1}{M} \langle |\mathscr{B}(t)| \rangle = b, \quad \mathcal{T} \text{ "PERSISTENCE TIME"}$

- Sampling protocols for the mini batch $\mathscr{B}(t)$
 - $\mu \in \mathscr{B}(t)$ with probability $b \in (0,1]$ i.i.d. at each time
 - $\frac{1}{M} \langle |\mathscr{B}(t)| \rangle = b \quad \text{"BATCH SIZE"}$



#samples Thermodynamic limit: #dimensions, #samples $\longrightarrow \infty$, at fixed sample complexity $\alpha = \frac{1}{\text{#dimensions}}$

Markovian dynamics of strongly coupled degrees of freedom $\rightarrow \infty$

[Sompolinsky, Zippelius, '81; Mézard, Parisi, Virasoro, '87; Sompolinsky, Crisanti, Sommers, '88 ... Agoritsas, Birol, Urbani, Zamponi (2018)]

Stochastic gradient flow equations

$$\dot{\mathbf{w}}(t) = -\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C}\left(\mathbf{w}(t), y_{\mu}, \mathbf{x}_{\mu}\right) - \nabla_{\mathbf{w}} r(\mathbf{w}(t)) \in \mathbb{I}$$







Markovian dynamics of strongly coupled degrees of freedom $\rightarrow \infty$

dimensional reduction

Non-Markovian dynamics of one effective degree of freedom

ĥ

W

[Sompolinsky, Zippelius, '81; Mézard, Parisi, Virasoro, '87; Sompolinsky, Crisanti, Sommers, '88 ... Agoritsas, Birol, Urbani, Zamponi (2018)]

#samples Thermodynamic limit: #dimensions, #samples $\longrightarrow \infty$, at fixed sample complexity $\alpha = \frac{1}{\text{#dimensions}}$

Stochastic gradient flow equations

$$\dot{\mathbf{w}}(t) = -\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C}\left(\mathbf{w}(t), y_{\mu}, \mathbf{x}_{\mu}\right) - \nabla_{\mathbf{w}} r(\mathbf{w}(t)) \in \mathbb{R}$$

One effective stochastic process + one ODE

$$(t) = -\tilde{\lambda}(t)h(t) - s(t)\Lambda'(y,h(t),m(t)) + \int_0^t dt' M_R(t,t')h(t') - h(t) = -\lambda m(t) - \mu(t)$$



ĥ

Markovian dynamics of strongly coupled degrees of freedom $\rightarrow \infty$

dimensional reduction

Non-Markovian dynamics of one effective degree of freedom

#samples Thermodynamic limit: #dimensions, #samples $\longrightarrow \infty$, at fixed sample complexity $\alpha = \frac{1}{\text{#dimensions}}$

Stochastic gradient flow equations

$$\dot{\mathbf{w}}(t) = -\sum_{\mu \in \mathscr{B}(t)} \nabla_{\mathbf{w}} \mathscr{C}\left(\mathbf{w}(t), y_{\mu}, \mathbf{x}_{\mu}\right) - \nabla_{\mathbf{w}} r(\mathbf{w}(t)) \in \mathbb{I}$$

One effective stochastic process + one ODE

$$(t) = -\tilde{\lambda}(t) h(t) - s(t) \Lambda'(y, h(t), m(t)) + \int_0^t dt' M_R(t, t') h(t') - \delta(t') h(t') + \delta(t') h(t') + \delta(t') h(t') h(t') + \delta(t') h(t') h(t') + \delta(t') h(t') h(t') + \delta(t') h(t') h(t') h(t') + \delta(t') h(t') h(t') h(t') + \delta(t') h(t') h(t'$$

 $\dot{m}(t) = -\lambda m(t) - \mu(t)$

To be solved self-consistently

[Eissfeller, Opper, '97; Roy et al., '19; Manacorda et al., '20]



The problem: BINARY GAUSSIAN MIXTURES CLASSIFICATION

Data distribution:

Goal: learn the <u>best separating</u> hyperplane



BINARY GAUSSIAN MIXTURES CLASSIFICATION The problem:

Vanilla SGD



[FM, F. Krzakala, P. Urbani, L. Zdeborová, NeurIPS (2020)] $\alpha = 2, 1/\tau = 0.6, dt = 0.2, simulations at N = 500$







Characterizing SGD noise

The problem: BINARY GAUSSIAN MIXTURES CLASSIFICATION

If the dynamics was at equilibrium — Fluctuation-Dissipation Theorem

$$R(t,t') = -\frac{1}{T}\partial_t C(t,t') \Theta(t-t')$$

Characterizing SGD noise

The problem: BINARY GAUSSIAN MIXTURES CLASSIFICATION

Linear response

$$R(t,t') = \lim_{\{H_j \to 0\}} \frac{1}{N} \sum_{j=1}^N \frac{\delta w_j(t)}{\delta H_j(t')} \leq \frac{1}{N} \sum_{j=1}^N \frac{\delta w_j(t)}{\delta H_j(t')} \leq$$

Integrating on both sides and rescaling:

If the dynamics was at equilibrium \longrightarrow Fluctuation-Dissipation Theorem

$$R(t,t') = -\frac{1}{T}\partial_t C(t,t') \Theta(t-t')$$

Correlation

$$C(t,t') = \frac{1}{N} \mathbf{w}(t)^{\top} \mathbf{w}(t)^{\top$$

$$\bar{\chi}(t,t') = \frac{1}{T} \left(1 - \bar{C}(t,t') \right)$$



Characterizing SGD noise

The problem: BINARY GAUSSIAN MIXTURES CLASSIFICATION

Linear response

$$R(t,t') = \lim_{\{H_j \to 0\}} \frac{1}{N} \sum_{j=1}^N \frac{\delta w_j(t)}{\delta H_j(t')} \leq \frac{1}{N} \sum_{j=1}^N \frac{\delta w_j(t)}{\delta H_j(t')} \leq$$

Integrating on both sides and rescaling:

$$\bar{\chi}(t,t') = \frac{1}{T} \left(1 - \bar{C}(t,t') \right)$$

If the dynamics was at equilibrium — Fluctuation-Dissipation Theorem

$$R(t,t') = -\frac{1}{T}\partial_t C(t,t') \Theta(t-t')$$

Correlation

$$C(t,t') = \frac{1}{N} \mathbf{w}(t)^{\top} \mathbf{w}(t)^{\top$$

Directly available from DMFT



The effective temperature of vanilla SGD

Changing the time step:



An effective fluctuation-dissipation theorem holds for SGD at stationarity in the under-parametrized regime.



The effective temperature of persistent SGD



THE SIGN-RETRIEVAL PROBLEM :

Data model: $\mathbf{x}_{\mu} \sim \mathcal{N}\left(\mathbf{0}, \mathbf{I}_{N}\right)$ $y_{\mu} = |\mathbf{x}_{\mu}^{\mathsf{T}}\mathbf{w}^*|$

Loss function:

$$\mathscr{E}\left(y_{\mu}, \mathbf{w}^{\mathsf{T}}\mathbf{x}_{\mu}\right) = \left(y_{\mu}^{2} - (\mathbf{w}^{\mathsf{T}}\mathbf{x}_{\mu})^{2}\right)^{2}$$



 $\mathbf{W}^{\top}\mathbf{W}^{*}$



FM, P. Urbani, L. Zdeborová, Machine Learning: Science and Technology (2021)

 $\mathbf{w}^{\top}\mathbf{w}^{*}$



FM, P. Urbani, L. Zdeborová, Machine Learning: Science and Technology (2021)

Vanilla SGD and GD get stuck in local minima

 $\mathbf{W}^{\top}\mathbf{W}^{*}$



FM, P. Urbani, L. Zdeborová, Machine Learning: Science and Technology (2021)

Persistent SGD achieves perfect recovery

Vanilla SGD and GD get stuck in local minima





The algorithmic noise produces an effective self-annealing.

FM, P. Urbani, L. Zdeborová, Machine Learning: Science and Technology (2021)



The theory keeps finite-size effects under control.

FM, P. Urbani, L. Zdeborová, Machine Learning: Science and Technology (2021)



Based on:

B. Loureiro, C. Gerbelot, H. Cui, S. Goldt, F. Krzakala, M. Mezard, L. Zdeborová, "Learning curves of generic features maps for realistic datasets with a teacher-student model". Advances in Neural Information Processing Systems 34 (2021).

Based on:

Advances in Neural Information Processing Systems 34 (2021).

The model:

 $egin{array}{c|c} oldsymbol{u}^{\mu} & \in \mathbb{R}^{p+d} \ oldsymbol{v}^{\mu} \end{array}$

$$^{l} \sim \mathcal{N}\left(0, \begin{bmatrix}\Psi & \Phi\\ \Phi^{\top} & \Omega\end{bmatrix}\right)$$

i.i.d.

Based on:

Advances in Neural Information Processing Systems 34 (2021).



$$^{l} \sim \mathcal{N}\left(0, \begin{bmatrix}\Psi & \Phi\\ \Phi^{\top} & \Omega\end{bmatrix}\right)$$

i.i.d.

Based on:

Advances in Neural Information Processing Systems 34 (2021).



$$^{l} \sim \mathcal{N}\left(0, \begin{bmatrix}\Psi & \Phi\\ \Phi^{\top} & \Omega\end{bmatrix}\right)$$

i.i.d.

Based on:

Advances in Neural Information Processing Systems 34 (2021).

The model:

 $egin{array}{c|c} oldsymbol{u}^{\mu} & \in \mathbb{R}^{p+d} \ oldsymbol{v}^{
u} & \in \mathbb{R}^{p+d} \end{array}$

 $y^{\mu} = f_0 \left(\frac{1}{\sqrt{p}} \boldsymbol{\theta}_0^{\top} \boldsymbol{u}^{\mu} \right)$

$$\hat{w} \sim \mathcal{N}\left(0, \begin{bmatrix} \Psi & \Phi \\ \Phi^{ op} & \Omega \end{bmatrix}
ight)$$

LEARNING
 $\hat{w} = \operatorname*{arg\,min}_{w \in \mathbb{R}^d} \left[\sum_{\mu=1}^n g\left(rac{w^{ op}v^{\mu}}{\sqrt{d}}, y^{\mu}
ight) + r(w)$



Modelling realistic data structure: the <u>asymptotic-time</u> learning curves

$\mathbf{x} \sim \text{DCGAN}$ TRAINED ON CIFAR10



Logistic regression, $f = \text{sign}, \mathbf{v} = \text{relu}(W_1 \text{relu}(W_2 \mathbf{x}))$

[Loureiro, et al. Advances in Neural Information Processing Systems 34 (2021)]



Ridge regression, f = id, $\mathbf{v} = relu\left(W_1 relu\left(W_2 \mathbf{x}\right)\right)$



MORE OPEN QUESTIONS:

- How to characterise the noise when the dynamics stops $(T_{\text{eff}} \rightarrow 0)$?
- Can we link the noise magnitude to generalisation properties?
- Exploring SGD with more realistic data models ([Goldt, et al., 2020]), optimisation strategies (e.g., momentum [Mannelli, Urbani, NeurIPS 2021]) and architectures (e.g., one-hidden layer networks)

. . .

Thank you for your attention!

