

# Loss Landscape and Performance in Deep Learning

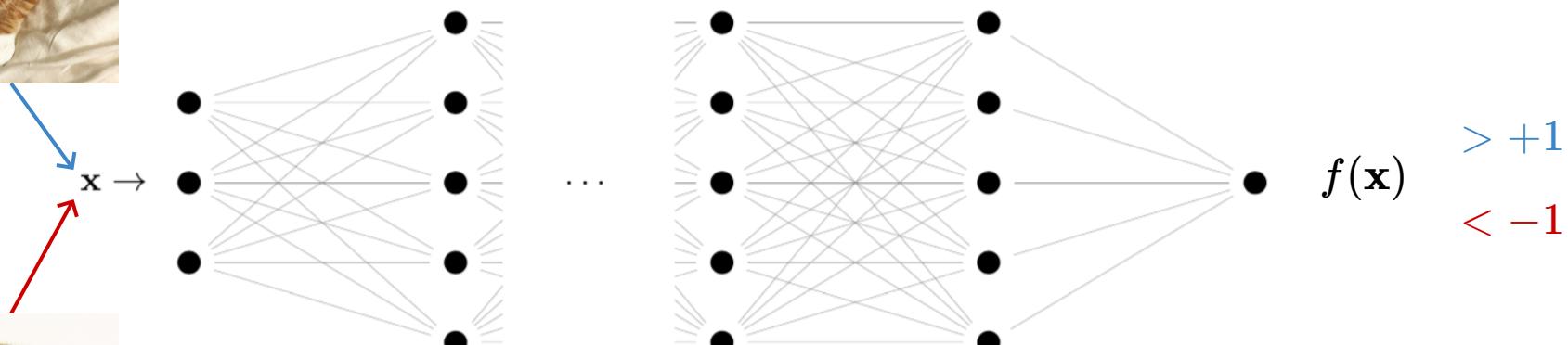


**Stefano Spigler**

M. Geiger, A. Jacot, S. d'Ascoli, M. Baity-Jesi,  
L. Sagun, G. Biroli, C. Hongler, M. Wyart

arXivs: 1901.01608; 1810.09665; 1809.09349

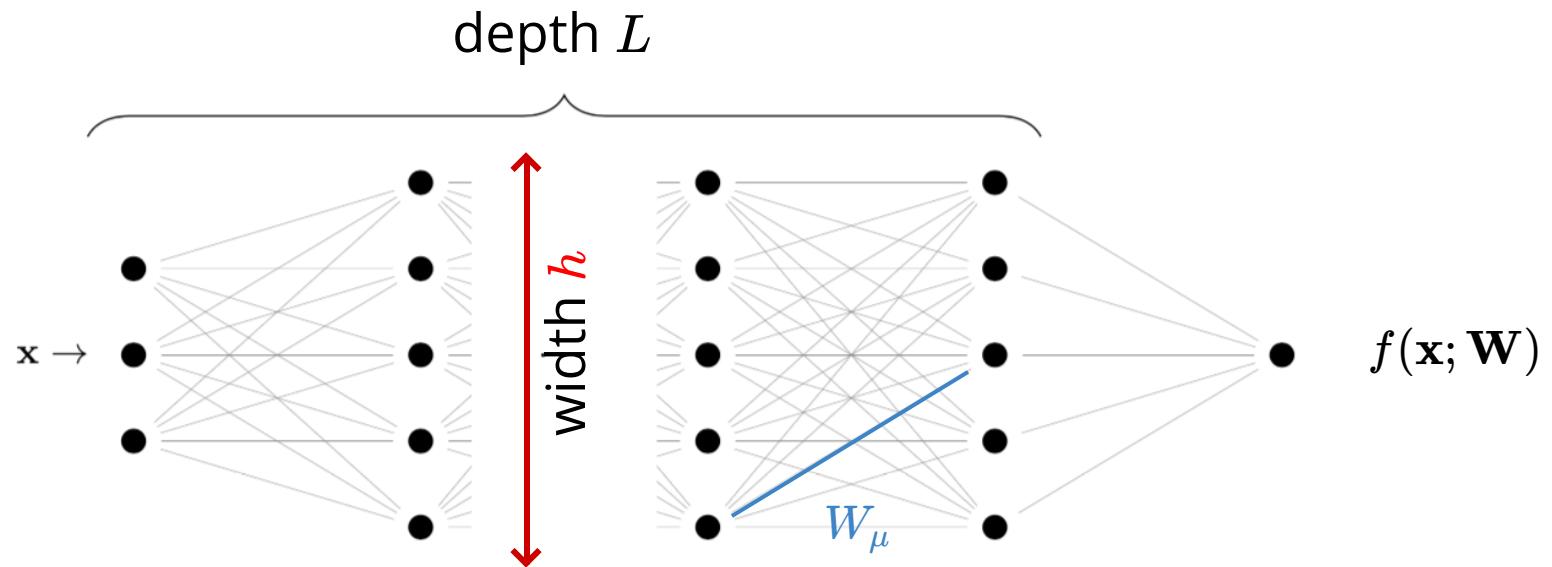
# (Supervised) Deep Learning



- Learning from examples: **train set**
- Is able to predict: **test set**
- Not understood why it works so well!
  - How many data are needed to learn?
  - What network size?

# Set-up: Architecture

- Deep net  $f(\mathbf{x}; \mathbf{W})$  with  $N \sim h^2 L$  parameters



- Alternating *linear* and *nonlinear* operations!

# Set-up: Dataset

- $P$  training data:

$$\mathbf{x}_1, \dots, \mathbf{x}_P$$

- Binary classification:

$$\mathbf{x}_i \rightarrow \text{label } y_i = \pm 1$$

$\pm 1$  = cats/dogs, yes/no, even/odd...

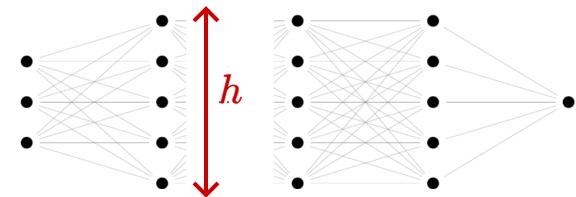
- Independent test set to evaluate performance

Example - MNIST (parity):

70k pictures, digits 0, ..., 9;  
use parity as label

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

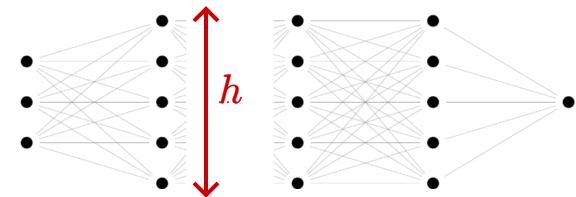
# Outline



Vary **network size**  $N$  ( $\sim h^2$ ):

1. Can networks fit **all** the  $P$  training data?
  
2. Can networks overfit? Can  $N$  be too large?  
→ Long term goal: how to choose  $N$ ?

# Outline



Vary **network size**  $N$  ( $\sim h^2$ ):

1. Can networks fit **all** the  $P$  training data?
2. Can networks overfit? Can  $N$  be too large?  
→ Long term goal: how to choose  $N$ ?

# Learning

- Find parameters  $\mathbf{W}$  such that  $\text{sign}f(\mathbf{x}_i; \mathbf{W}) = y_i$  for  $i \in \text{train set}$

Binary classification:

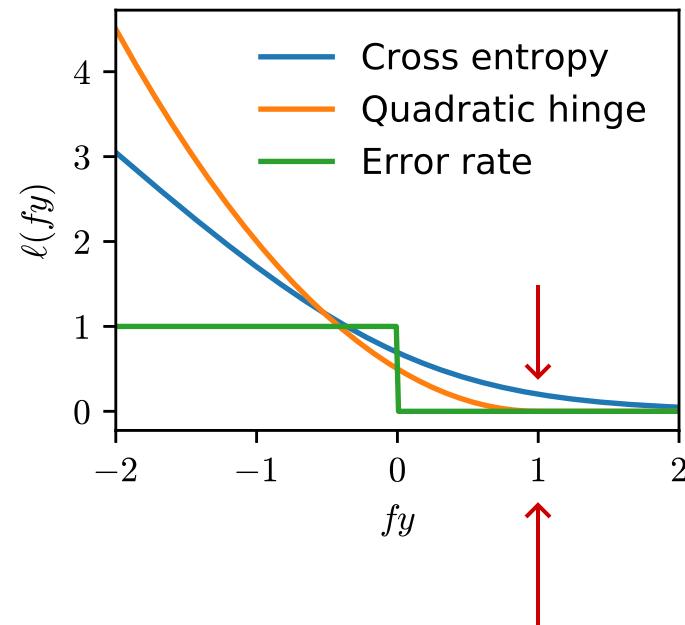
$$y_i = \pm 1$$

- **Minimize some loss!**

$$\mathcal{L}(\mathbf{W}) = \sum_{i=1}^P \ell(y_i f(\mathbf{x}_i; \mathbf{W}))$$

Hinge loss:

- $\mathcal{L}(\mathbf{W}) = 0$  if and only if  $y_i f(\mathbf{x}_i; \mathbf{W}) > 1$  for all patterns

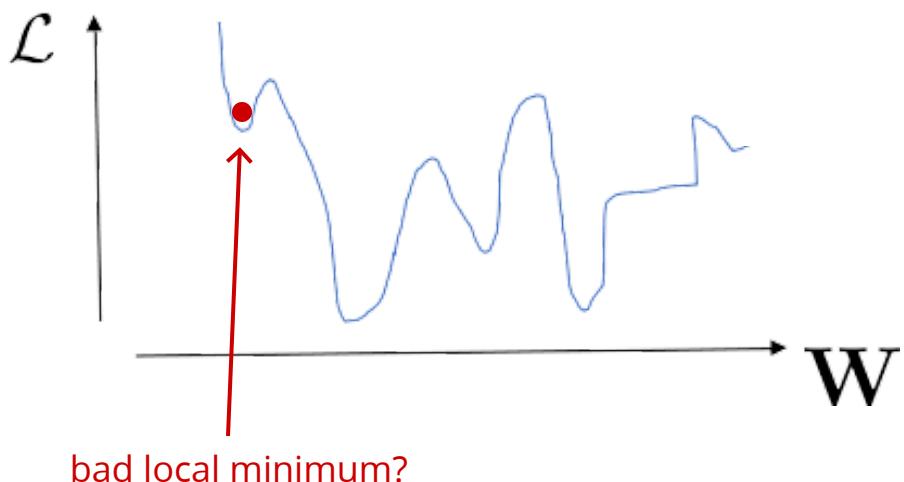


(classified correctly with  
some margin)

# Learning dynamics = descent in loss landscape

- Minimize loss  $\longleftrightarrow$  **gradient descent**
- Start with **random initial conditions!**

Random, high dimensional, not convex landscape!



- Why not stuck in bad local minima?
- What is the landscape geometry?

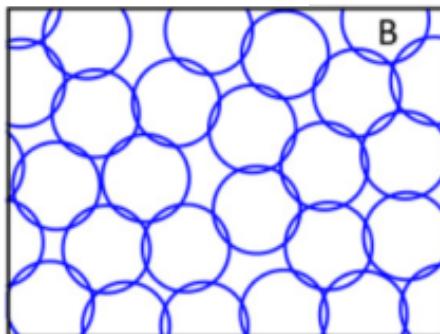
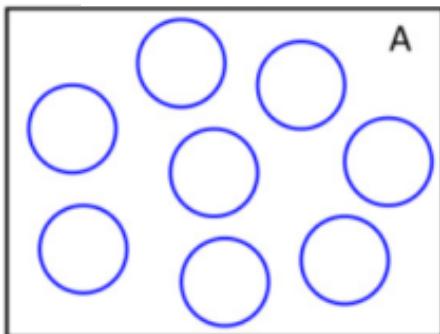
in practical settings:

- Many flat directions are found!

Soudry, Hoffer '17; Sagun et al. '17; Cooper '18;  
Baity-Jesy et al. '18 - arXiv:1803.06969

# Analogy with granular matter: Jamming

## Random packing:

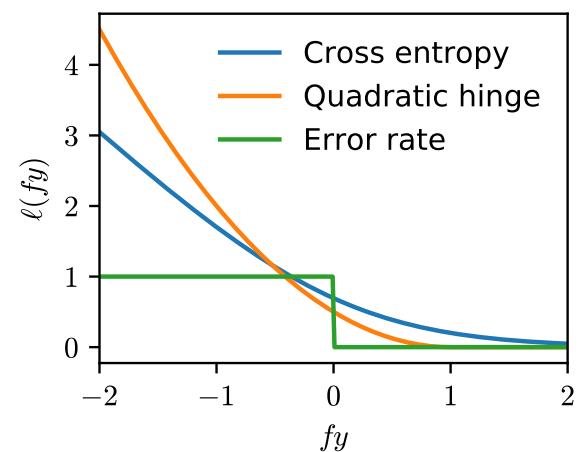


- random initial conditions
- minimize energy  $\mathcal{L}$
- either find  $\mathcal{L} = 0$  or  $\mathcal{L} > 0$

Upon increasing density  $\rightarrow$  transition

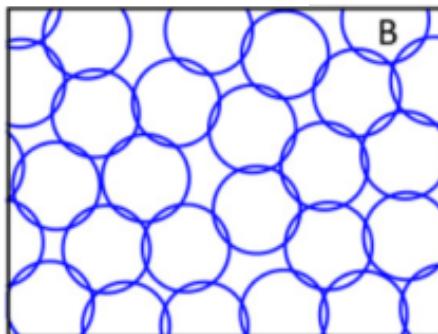
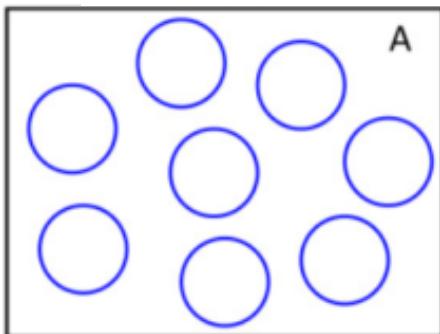
sharp transition with finite-range interactions

this is why we use the **hinge loss!**



# Analogy with granular matter: Jamming

## Random packing:



- random initial conditions
- minimize energy  $\mathcal{L}$
- either find  $\mathcal{L} = 0$  or  $\mathcal{L} > 0$

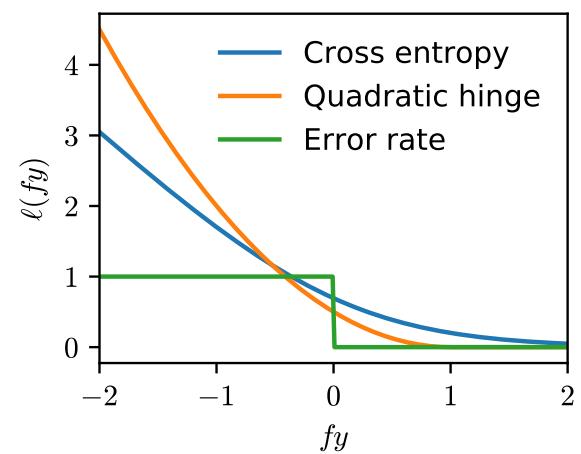
Upon increasing density  $\rightarrow$  transition

sharp transition with finite-range interactions

this is why we use the **hinge loss!**

Shallow networks  $\longleftrightarrow$  packings of **spheres**: Franz and Parisi, '16

Deep nets  $\longleftrightarrow$  packings of **ellipsoids**!



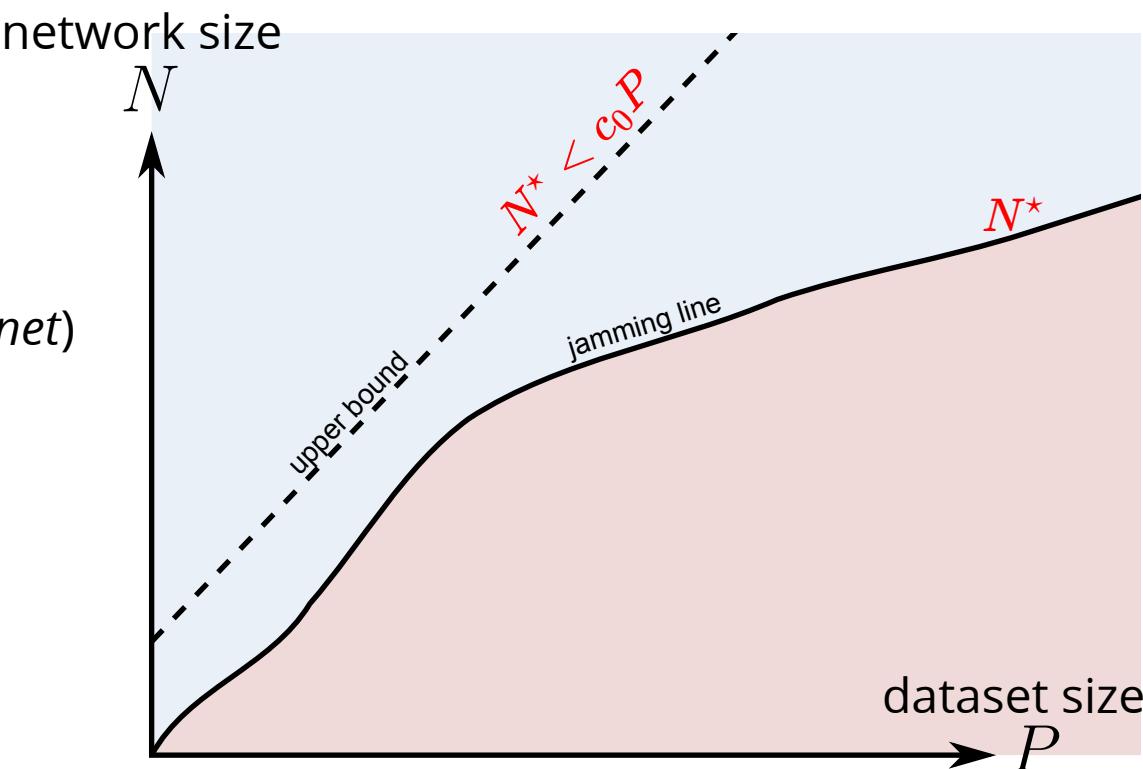
# Theoretical results: Phase diagram

- When  $N$  is large,  $\mathcal{L} = 0$
- Transition at  $N^*$

(if signals propagate through the net)

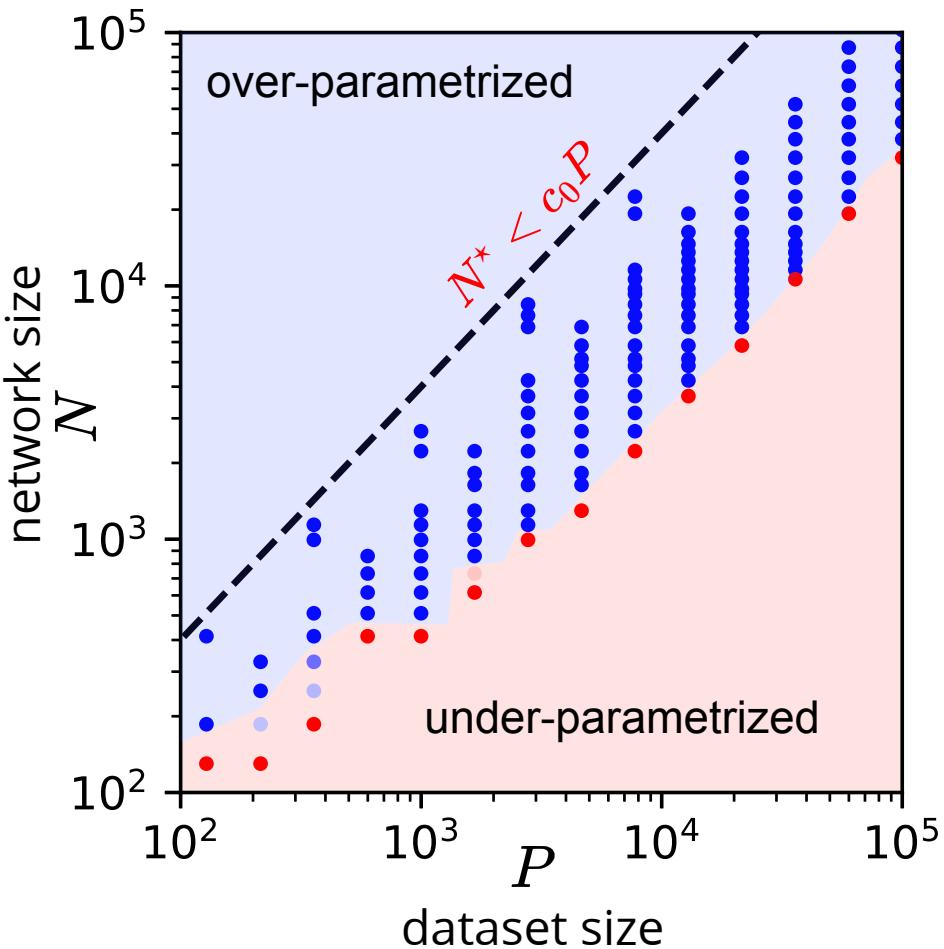
$$N^* < c_0 P$$

typically  $c_0 = \mathcal{O}(1)$



# Empirical tests: MNIST (parity)

Geiger et al. '18 - arXiv:1809.09349;  
Spigler et al. '18 - arXiv:1810.09665



No local minima are found  
when **overparametrized**!

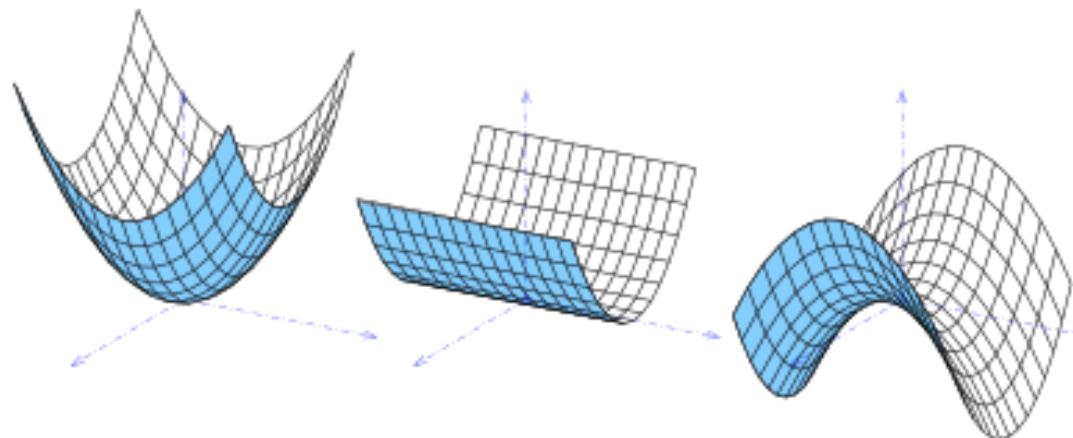
- Above  $N^*$  we have  $\mathcal{L} = 0$
- Solid line is the bound  $N^* < c_0 P$

We don't find local minima when overparametrized...

...shape of the landscape?

Local curvature:  
second order approximation

w.r.t parameters  $\mathbf{W}$



Information captured by **Hessian matrix**:  $\mathcal{H}_{\mu\nu} = \frac{\partial^2}{\partial \mathbf{w}_\mu \partial \mathbf{w}_\nu} \mathcal{L}(\mathbf{W})$

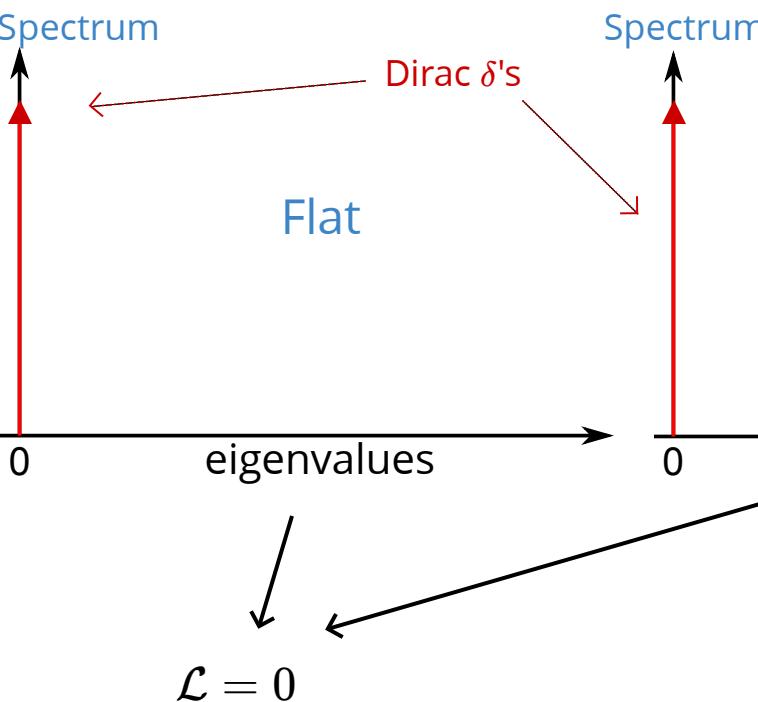


**Spectrum of the Hessian** (eigenvalues)

# Flat directions

Geiger et al. '18 - arXiv:1809.09349

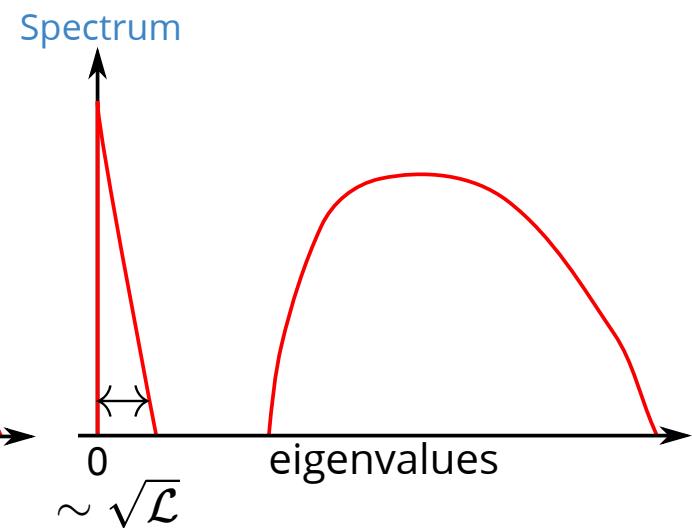
Over-parametrized  
 $N > N^*$



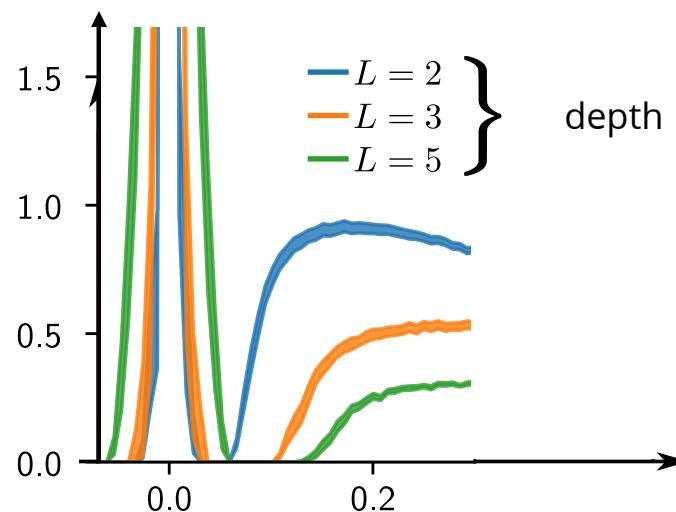
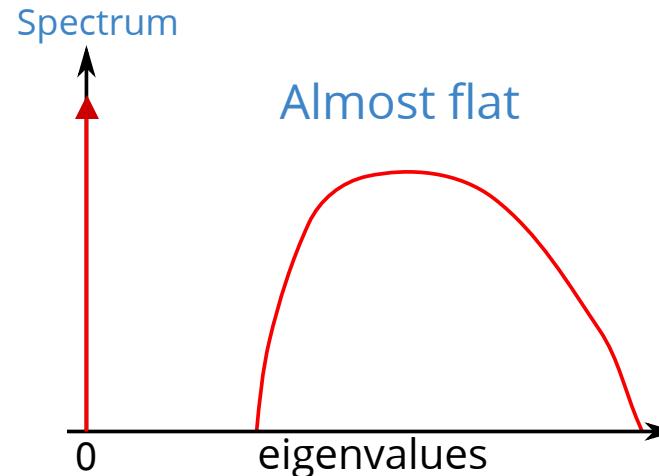
Jamming  
 $N \approx N^*$



Under-parametrized  
 $N < N^*$

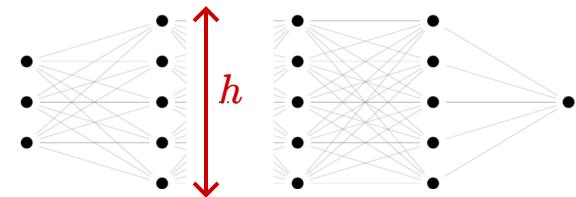


Jamming  
 $N \approx N^*$



From numerical simulations:  
*(at the transition)*

# Outline



Vary **network size**  $N$  ( $\sim h^2$ ):

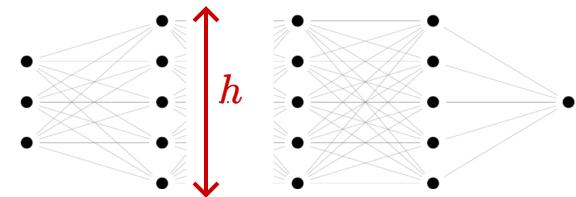
1. Can networks fit **all** the  $P$  training data?

**Yes**, deep networks **fit all data** if  $N > N^*$   $\rightarrow$  *jaming transition*

2. Can networks overfit? Can  $N$  be too large?

$\rightarrow$  Long term goal: how to choose  $N$ ?

# Outline



Vary **network size**  $N$  ( $\sim h^2$ ):

1. Can networks fit **all** the  $P$  training data?

**Yes**, deep networks **fit all data** if  $N > N^*$   $\rightarrow$  *jaming transition*

2. Can networks overfit? Can  $N$  be too large?

$\rightarrow$  Long term goal: how to choose  $N$ ?

Ok, so just crank up  $N$  and fit everything?

Generalization? → Compute **test error**  $\epsilon$

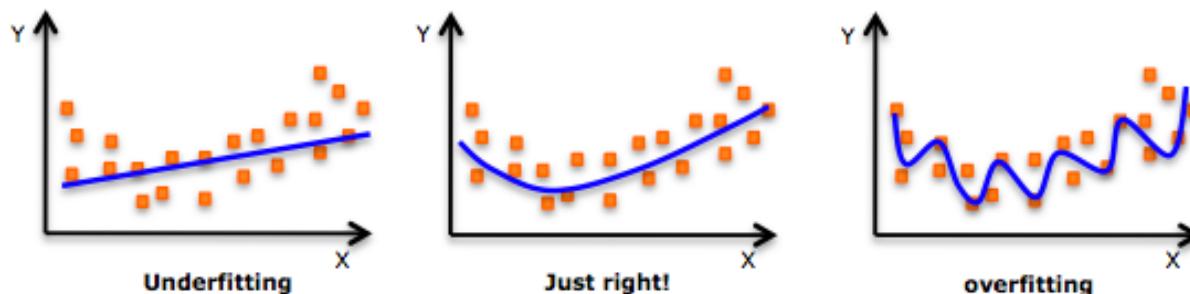
But wait... what about **overfitting**?

Ok, so just crank up  $N$  and fit everything?

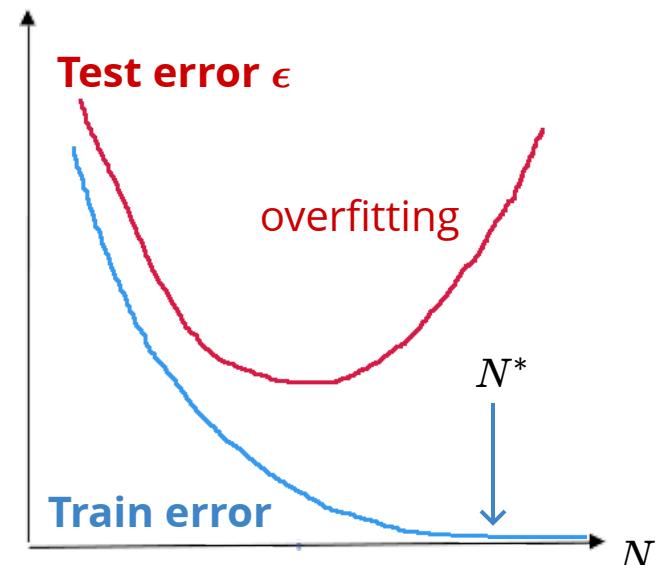
Generalization? → Compute **test error**  $\epsilon$

But wait... what about **overfitting**?

example: polynomial fitting

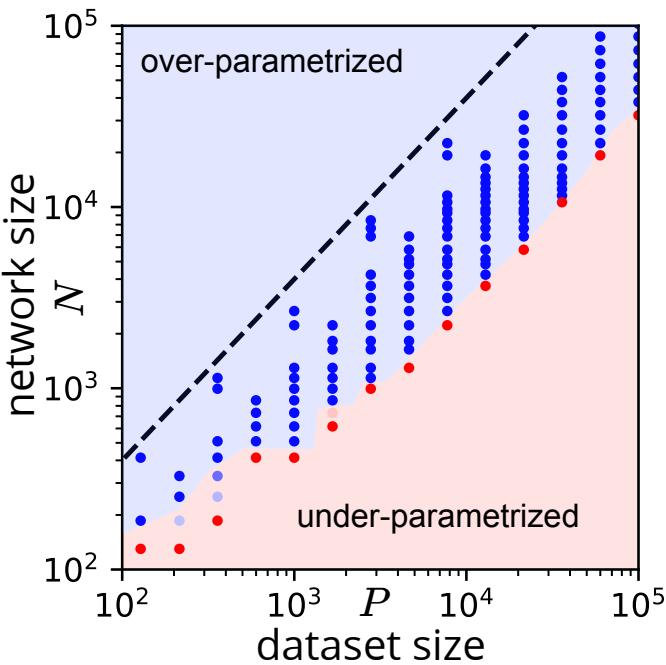


$N \sim$  polynomial degree



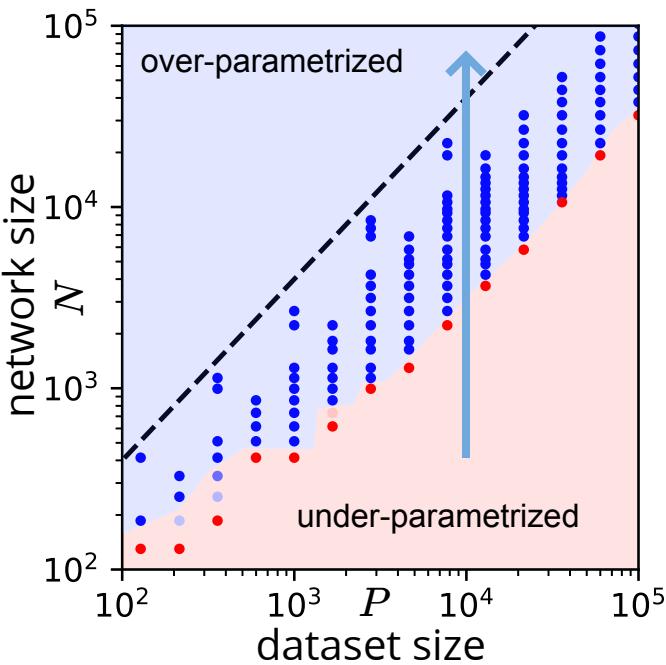
# Overfitting?

Spigler et al. '18 - arXiv:1810.09665



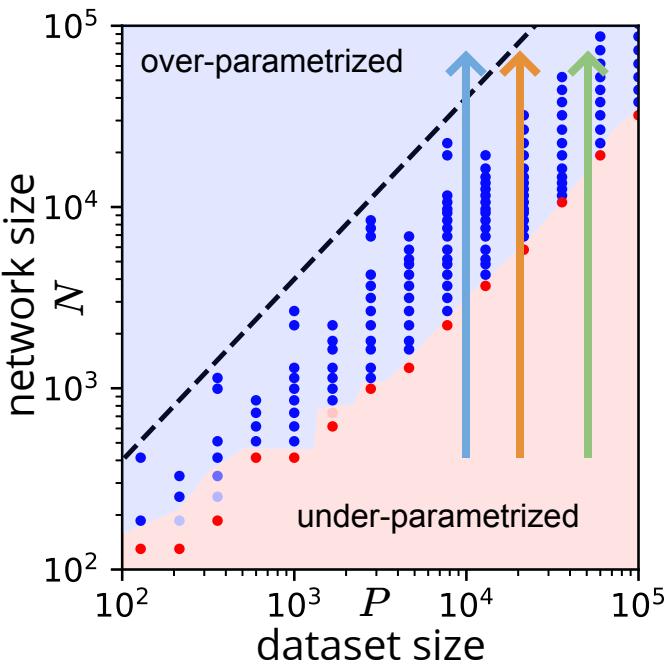
# Overfitting?

Spigler et al. '18 - arXiv:1810.09665



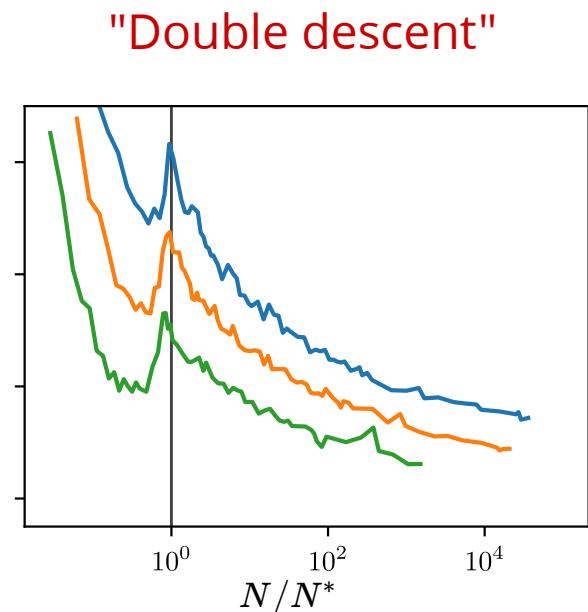
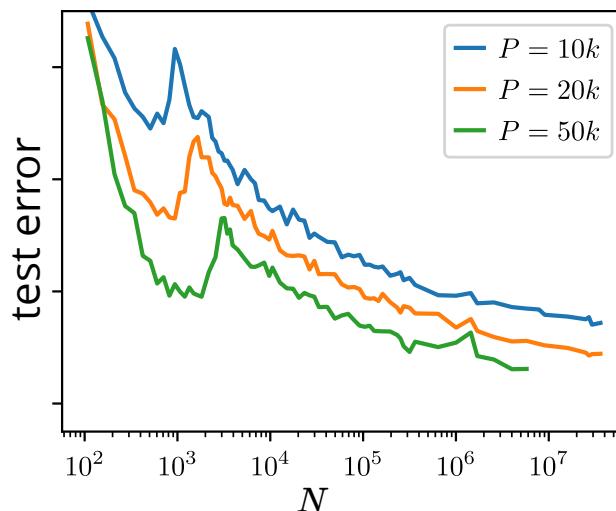
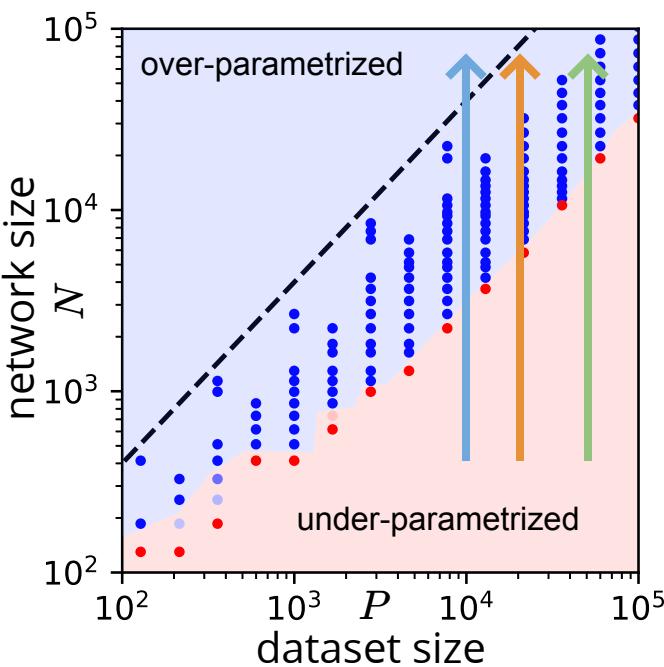
# Overfitting?

Spigler et al. '18 - arXiv:1810.09665



# Overfitting?

Spigler et al. '18 - arXiv:1810.09665



- **Test error decreases monotonically with  $N$ !**

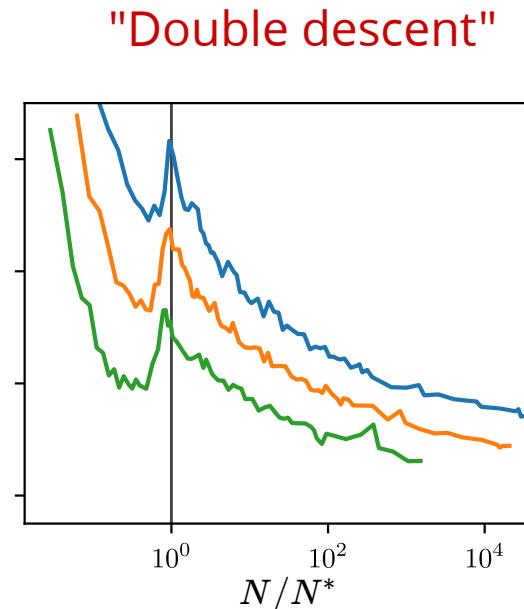
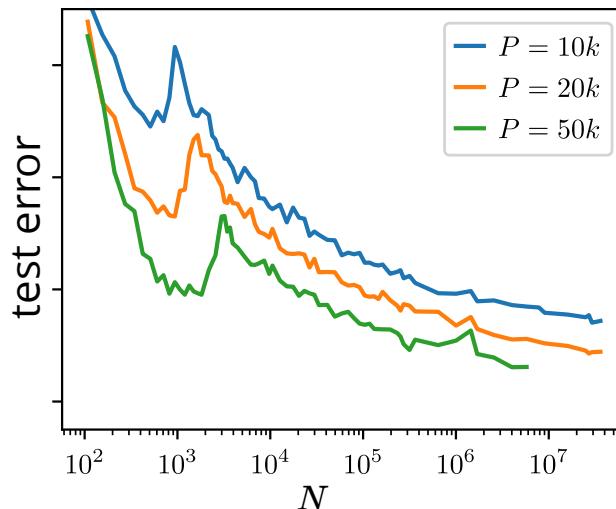
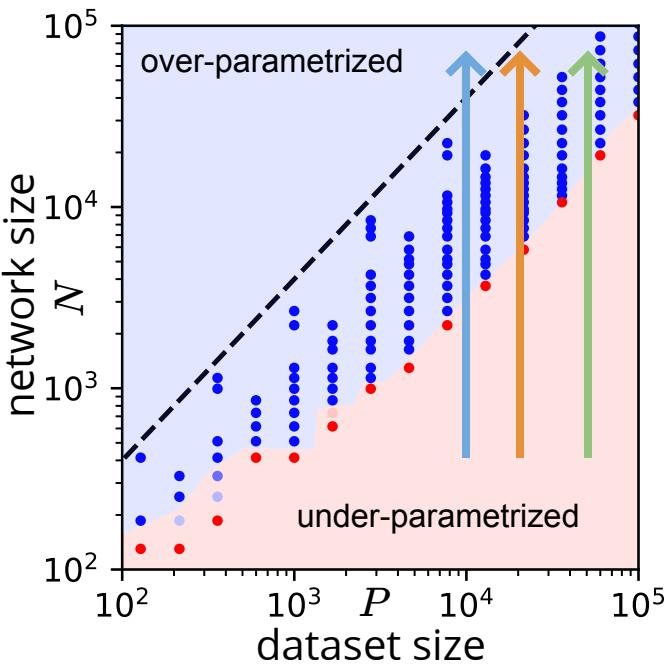
(after the peak)

- **Cusp** at the jamming transition

Advani and Saxe '17;  
Spigler et al. '18 - arXiv:1810.09665;  
Geiger et al. '19 - arXiv:1901.01608

# Overfitting?

Spigler et al. '18 - arXiv:1810.09665



- **Test error decreases monotonically with  $N$ !**

We know why: Fluctuations!

(after the peak)

- **Cusp** at the jamming transition

Advani and Saxe '17;

Spigler et al. '18 - arXiv:1810.09665;  
Geiger et al. '19 - arXiv:1901.01608

## Ensemble average

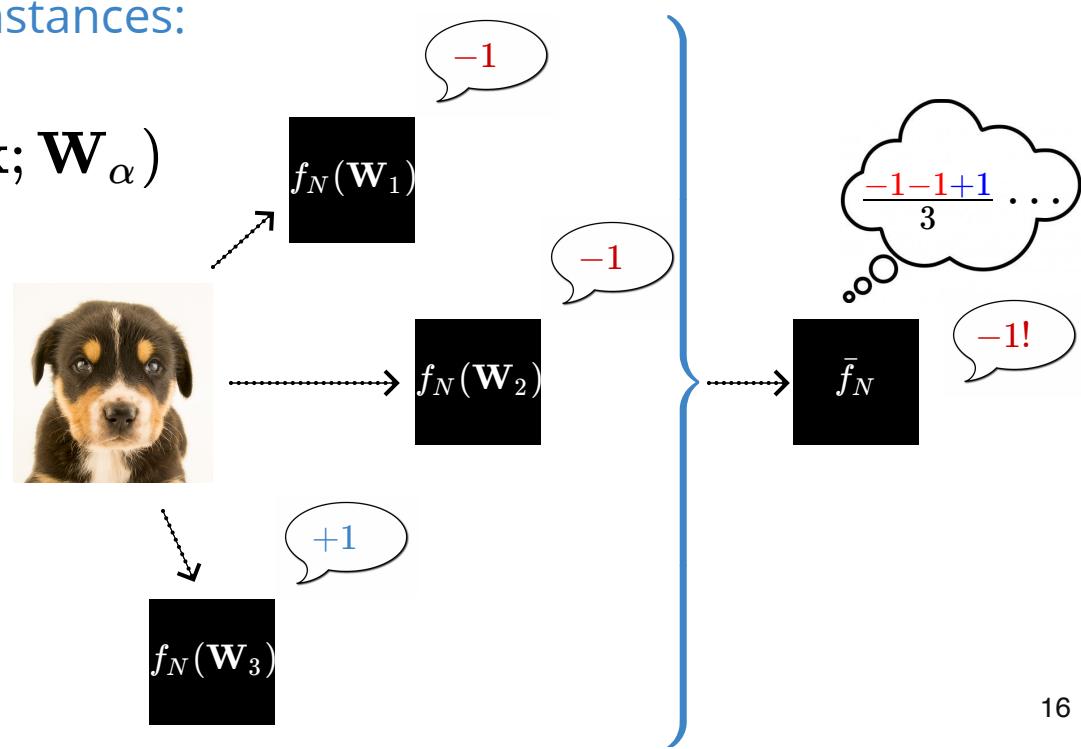
- Random initialization → output function  $f_N$  is **stochastic**
- Fluctuations: quantified by **average** and **variance**

# Ensemble average

- Random initialization → output function  $f_N$  is **stochastic**
- Fluctuations: quantified by **average** and **variance**

**ensemble average** over  $n$  instances:

$$\bar{f}_N^n(\mathbf{x}) \equiv \frac{1}{n} \sum_{\alpha=1}^n f_N(\mathbf{x}; \mathbf{W}_\alpha)$$



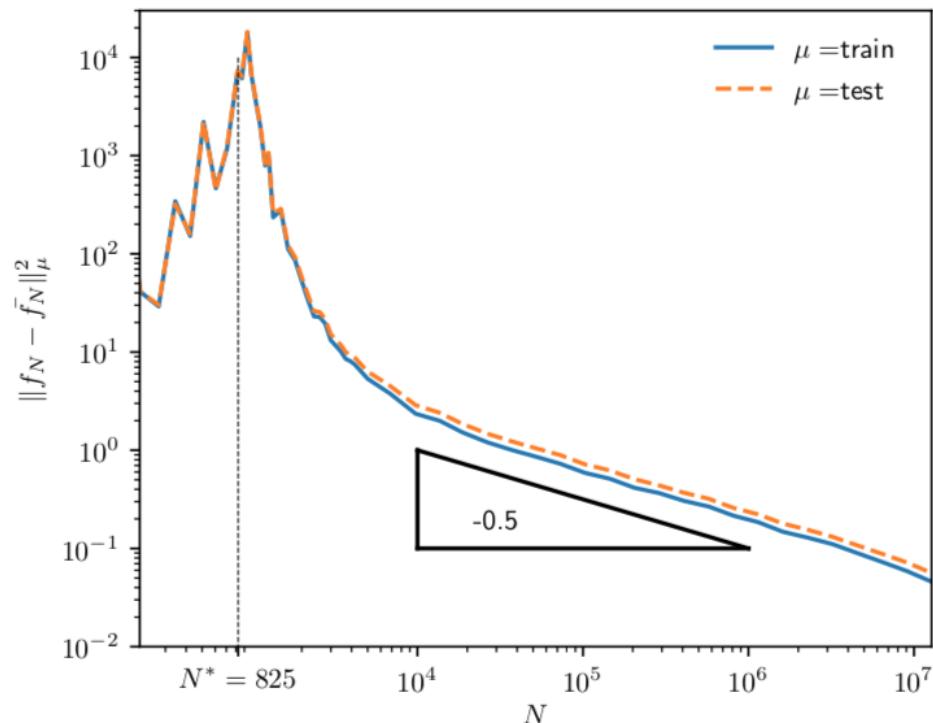
# Ensemble average

- Random initialization → output function  $f_N$  is **stochastic**
- Fluctuations: quantified by **average** and **variance**

Define some norm over the output functions:

**ensemble variance** (fixed  $n$ ):

$$\|f_N - \bar{f}_N^n\|^2 \sim N^{-\frac{1}{2}}$$



# Ensemble average

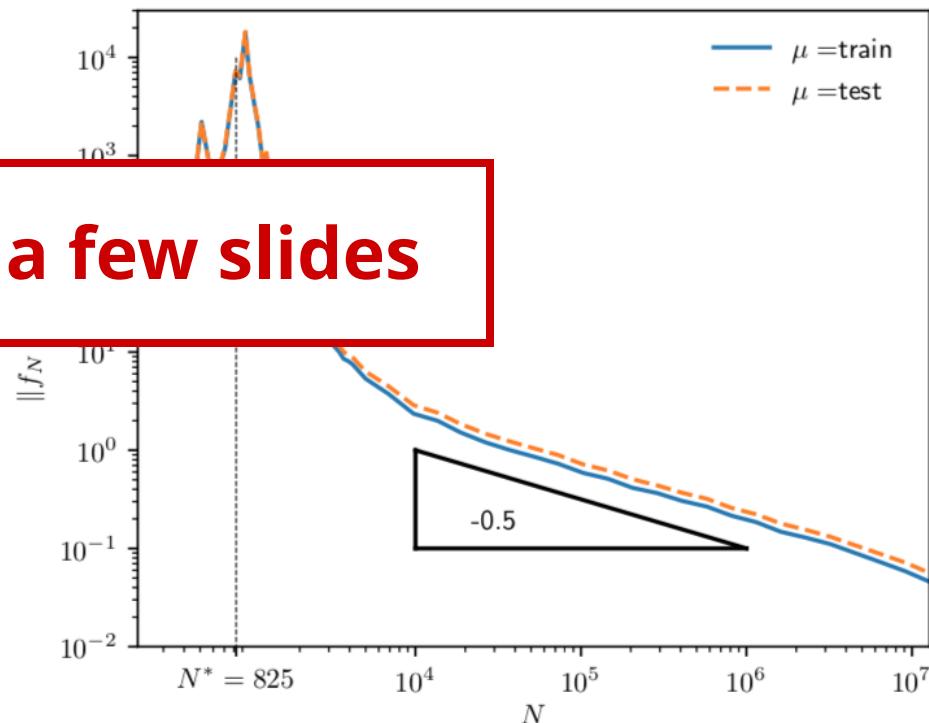
- Random initialization → output function  $f_N$  is **stochastic**
- Fluctuations: quantified by **average** and **variance**

Define some norm over the output functions:

ensemble va

Explained in a few slides

$$\|f_N - \bar{f}_N^n\|^2 \sim N^{-\frac{1}{2}}$$



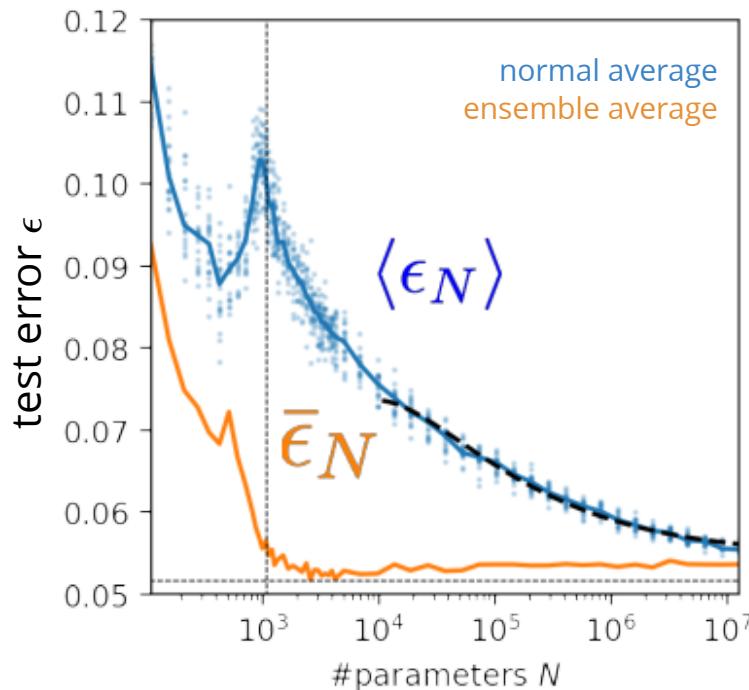
Remark: *test error of ensemble average*  $\neq$  *average test error*

$$\bar{f}_N^n(\mathbf{x}) \rightarrow \bar{\epsilon}_N \quad \{f(\mathbf{x}; \mathbf{W}_\alpha)\} \rightarrow \langle \epsilon_N \rangle$$

Remark: *test error of ensemble average*  $\neq$  *average test error*

$$\bar{f}_N^n(\mathbf{x}) \rightarrow \bar{\epsilon}_N$$

$$\{f(\mathbf{x}; \mathbf{W}_\alpha)\} \rightarrow \langle \epsilon_N \rangle$$



- **Test error** increases with fluctuations
- **Ensemble test error** is nearly flat after  $N^*$ !

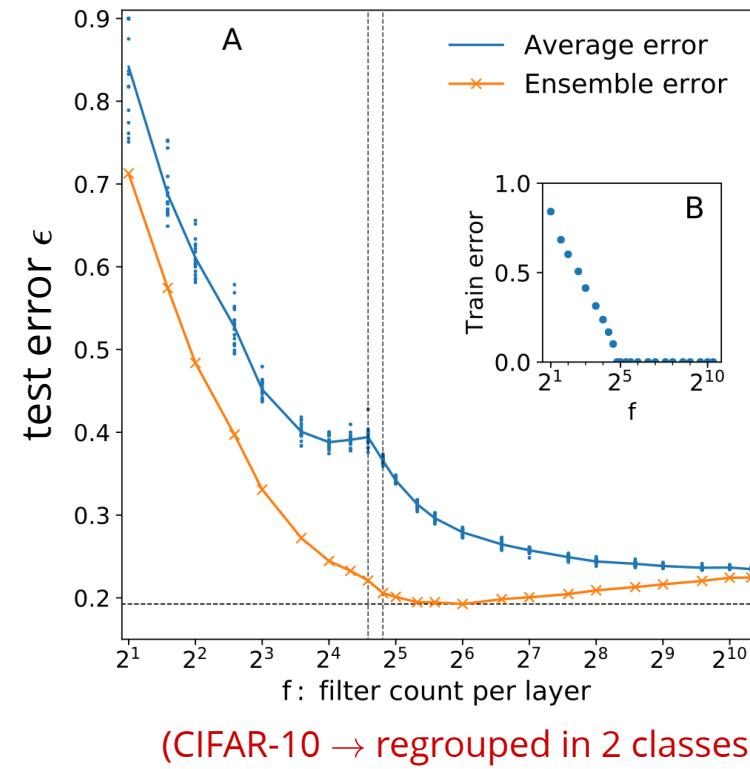
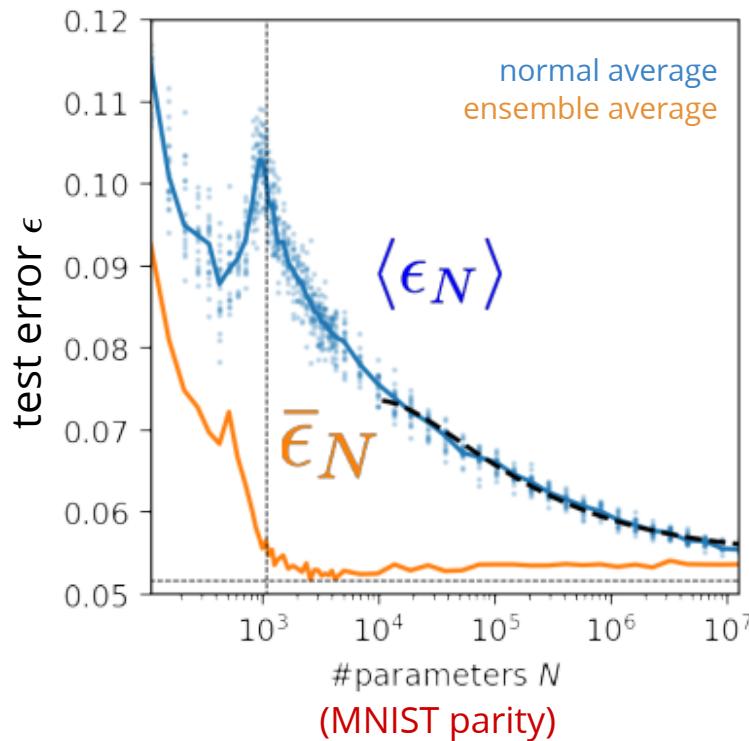
# Fluctuations increase error

Geiger et al. '19 - arXiv:1901.01608

Remark: *test error of ensemble average*  $\neq$  *average test error*

$$\bar{f}_N^n(\mathbf{x}) \rightarrow \bar{\epsilon}_N$$

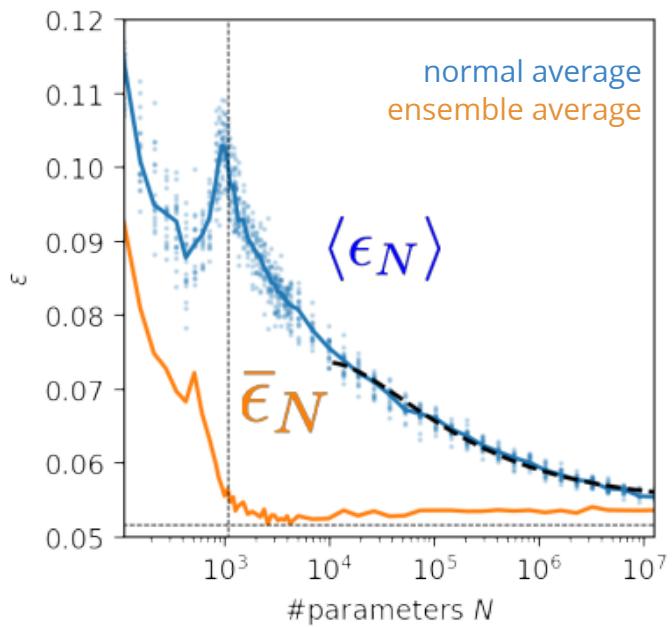
$$\{f(\mathbf{x}; \mathbf{W}_\alpha)\} \rightarrow \langle \epsilon_N \rangle$$



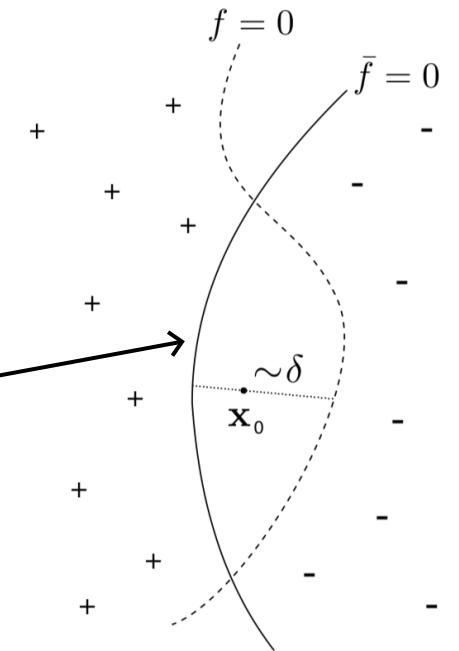
(CIFAR-10  $\rightarrow$  regrouped in 2 classes)

# Scaling argument!

Geiger et al. '19 - arXiv:1901.01608

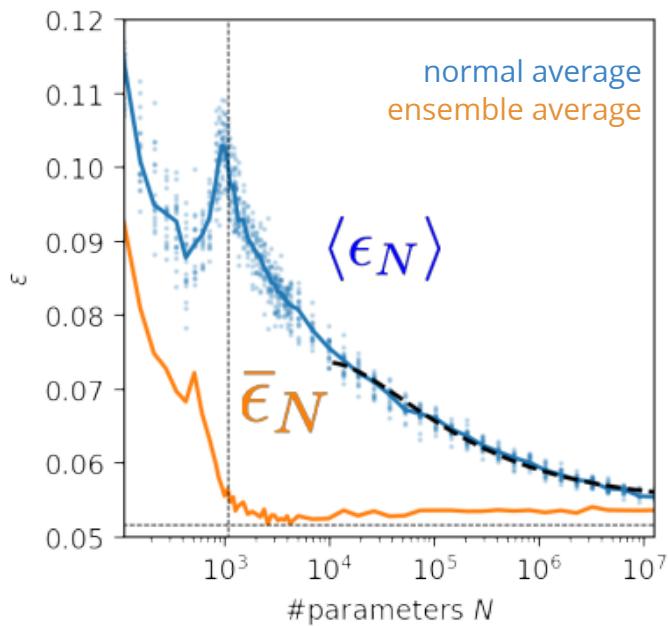


decision boundaries:

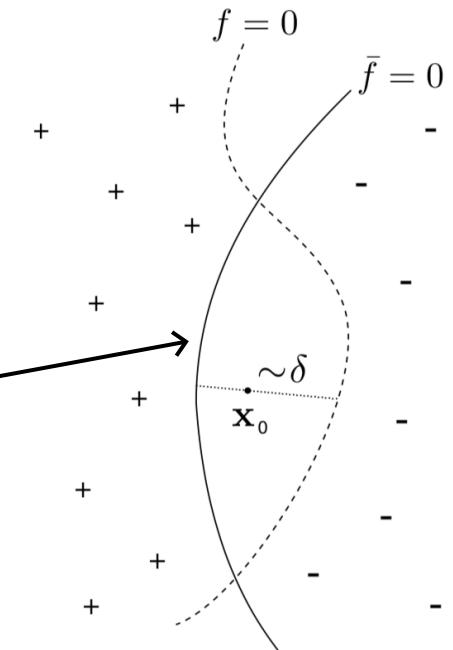


# Scaling argument!

Geiger et al. '19 - arXiv:1901.01608



decision boundaries:

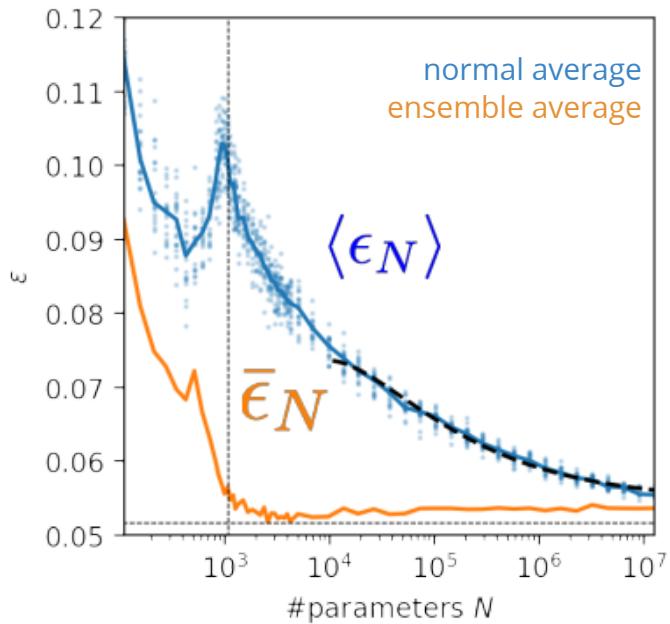


Smoothness of test error as function of  
decision boundary + symmetry:

$$\langle \epsilon_N \rangle - \bar{\epsilon}_N \sim \|f_N - \bar{f}_N\|^2 \sim N^{-\frac{1}{2}}$$

# Scaling argument!

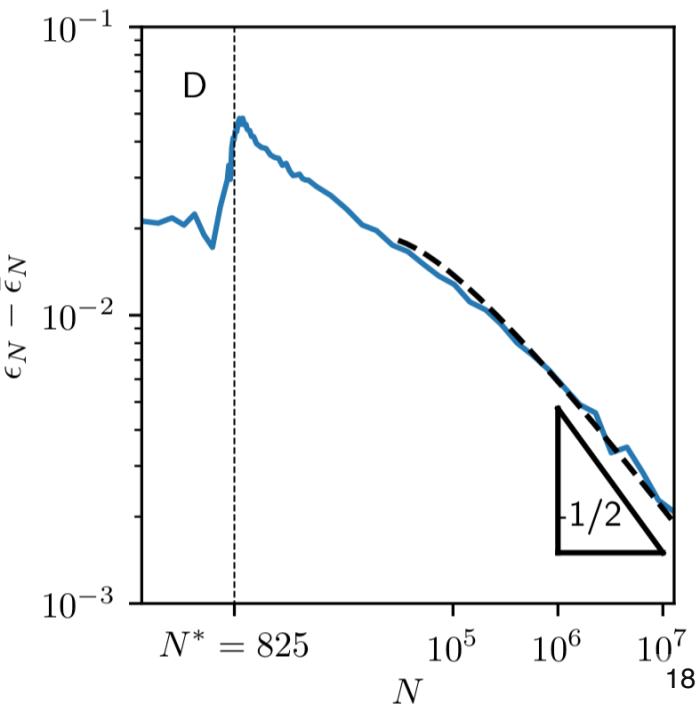
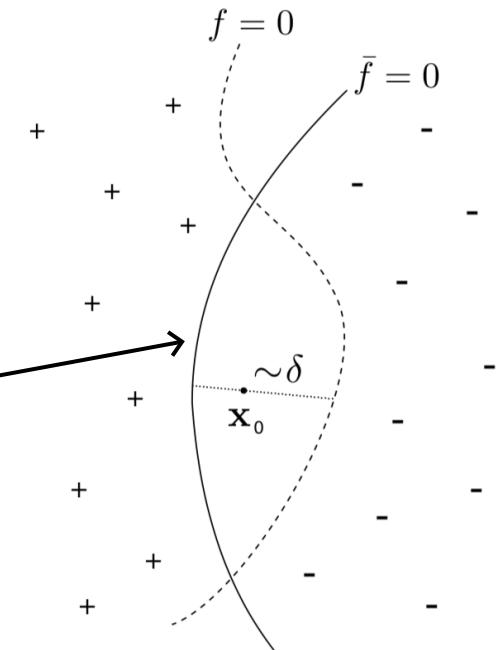
Geiger et al. '19 - arXiv:1901.01608



Smoothness of test error as function of decision boundary + symmetry:

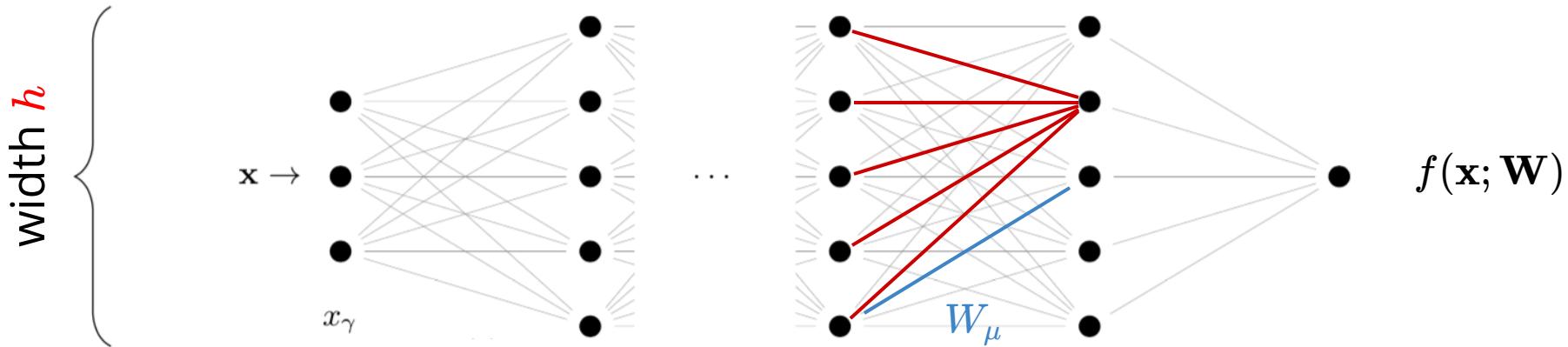
$$\langle \epsilon_N \rangle - \bar{\epsilon}_N \sim \|f_N - \bar{f}_N\|^2 \sim N^{-\frac{1}{2}}$$

decision boundaries:



# Infinitely-wide networks: Initialization

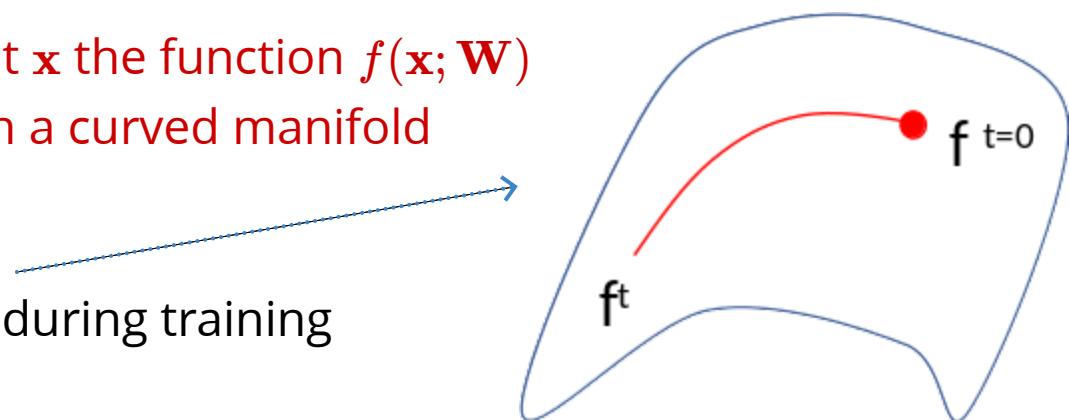
Neal '96; Williams '98; Lee et al '18; Schoenholz et al. '16



- Weights: each initialized as  $W_\mu \sim h^{-\frac{1}{2}} \mathcal{N}(0, 1)$
- Neurons sum  $h$  signals of order  $h^{-\frac{1}{2}}$   $\longrightarrow$  **Central Limit Theorem**
- Output function becomes a **Gaussian Random Field** as  $h \rightarrow \infty$

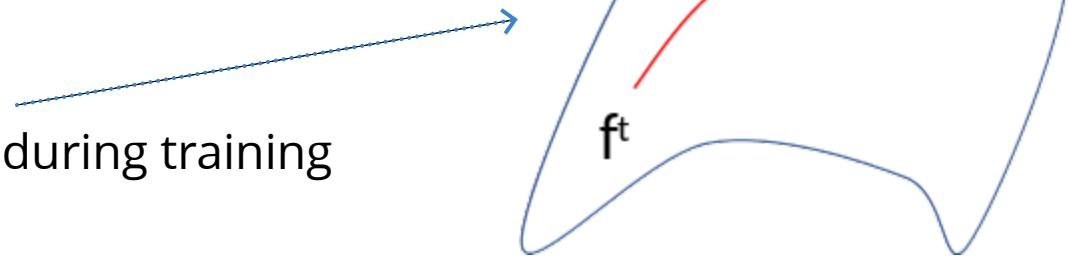
For an input  $\mathbf{x}$  the function  $f(\mathbf{x}; \mathbf{W})$   
lives on a curved manifold

- For **small** width  $h$ :  $\nabla_{\mathbf{W}} f$  evolves during training



For an input  $\mathbf{x}$  the function  $f(\mathbf{x}; \mathbf{W})$   
lives on a curved manifold

- For **small** width  $h$ :  $\nabla_{\mathbf{W}} f$  evolves during training
- For **large** width  $h$ :  $\nabla_{\mathbf{W}} f$  is constant during training

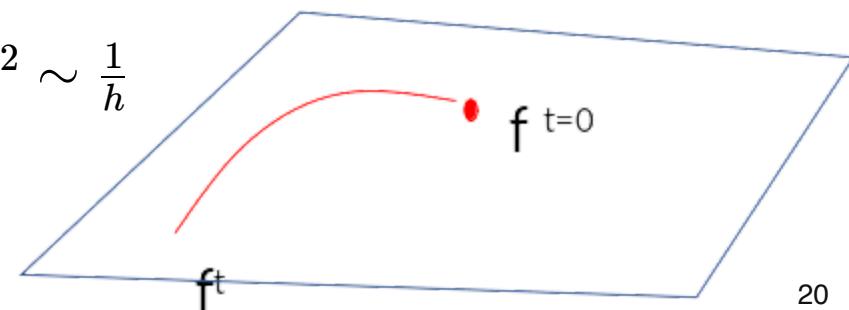


The manifold becomes linear!



Lazy learning:

- weights don't change much:  $\|\mathbf{W}^t - \mathbf{W}^{t=0}\|^2 \sim \frac{1}{h}$
- enough to change the output  $f$  by  $\sim \mathcal{O}(1)$ !



# Neural Tangent Kernel

- Gradient descent implies:

$$\frac{d}{dt} f(\mathbf{x}; \mathbf{W}^t) = \sum_{i=1}^P \Theta^t(\mathbf{x}, \mathbf{x}_i) \ y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$$

# Neural Tangent Kernel

- Gradient descent implies:

convolution with a kernel

$$\frac{d}{dt} f(\mathbf{x}; \mathbf{W}^t) = \sum_{i=1}^P \Theta^t(\mathbf{x}, \mathbf{x}_i) y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$$

$\Theta^t(\mathbf{x}, \mathbf{x}') = \nabla_{\mathbf{W}} f(\mathbf{x}; \mathbf{W}^t) \cdot \nabla_{\mathbf{W}} f(\mathbf{x}'; \mathbf{W}^t)$

The diagram illustrates the gradient of a neural tangent function. It shows the derivative with respect to time  $t$  of the function  $f(\mathbf{x}; \mathbf{W}^t)$  as a sum over  $P$  data points  $i$ . Each term in the sum is the product of the kernel value  $\Theta^t(\mathbf{x}, \mathbf{x}_i)$  (which is highlighted with a red box) and the weighted gradient  $y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$ . Above the equation, a blue brace groups the sum and the kernel term under the heading "convolution with a kernel". Below the equation, a red arrow points from the highlighted kernel term to the formula for the kernel itself, which is defined as the dot product of the gradients of the function at points  $\mathbf{x}$  and  $\mathbf{x}'$ .

The formula for the *kernel*  $\Theta^t$  is useless, unless...

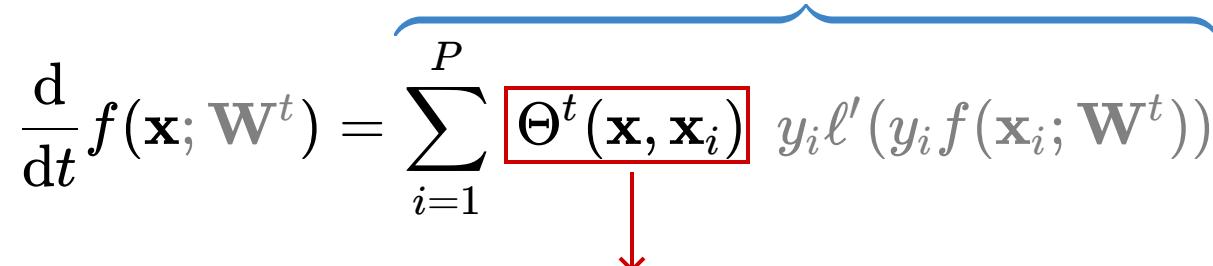
# Neural Tangent Kernel

- Gradient descent implies:

convolution with a kernel

$$\frac{d}{dt} f(\mathbf{x}; \mathbf{W}^t) = \sum_{i=1}^P \Theta^t(\mathbf{x}, \mathbf{x}_i) y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$$

$\Theta^t(\mathbf{x}, \mathbf{x}') = \nabla_{\mathbf{W}} f(\mathbf{x}; \mathbf{W}^t) \cdot \nabla_{\mathbf{W}} f(\mathbf{x}'; \mathbf{W}^t)$



The formula for the *kernel*  $\Theta^t$  is useless, unless...

Theorem. (informal)  $\lim_{\text{width } h \rightarrow \infty} \Theta^t(\mathbf{x}, \mathbf{x}') \equiv \Theta_\infty(\mathbf{x}, \mathbf{x}')$

Jacot et al. '18

Deep learning = learning with a **kernel** as  $h \rightarrow \infty$

# Finite $N$ asymptotics?

Geiger et al. '19 - arXiv:1901.01608;  
Hanin and Nica '19;  
Dyer and Gur-Ari '19

- **Evolution in time** is **small**:  $\|\Theta^t - \Theta^{t=0}\|_F \sim 1/h \sim N^{-\frac{1}{2}}$
- **Fluctuations** are much **larger**:  $\Delta\Theta^{t=0} \sim 1/\sqrt{h} \sim N^{-\frac{1}{4}}$   
*at  $t = 0$*

$$f(\mathbf{x}; \mathbf{W}^t) = \int dt \sum_{i=1}^P \Theta^t(\mathbf{x}, \mathbf{x}_i) y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$$

# Finite $N$ asymptotics?

Geiger et al. '19 - arXiv:1901.01608;  
Hanin and Nica '19;  
Dyer and Gur-Ari '19

- **Evolution in time** is small:  $\|\Theta^t - \Theta^{t=0}\|_F \sim 1/h \sim N^{-\frac{1}{2}}$
- **Fluctuations** are much **larger**:  $\Delta\Theta^{t=0} \sim 1/\sqrt{h} \sim N^{-\frac{1}{4}}$   
at  $t = 0$

$$f(\mathbf{x}; \mathbf{W}^t) = \int dt \sum_{i=1}^P \Theta^t(\mathbf{x}, \mathbf{x}_i) y_i \ell'(y_i f(\mathbf{x}_i; \mathbf{W}^t))$$



Then:  $\|f_N - \bar{f}_N\|^2 \sim (\Delta\Theta^{t=0})^2 \sim N^{-\frac{1}{2}}$

The output function fluctuates similarly to the kernel

# Conclusion

1. Can networks fit **all** the  $P$  training data?

- **Yes**, deep networks **fit all data** if  $N > N^*$   $\rightarrow$  *jamming transition*

2. Can networks overfit? Can  $N$  be too large?

- *Initialization* induces *fluctuations* in output that increase *test error*
- **No overfitting:** error keeps decreasing past  $N^*$  because *fluctuations diminish*  
check Geiger et al. '19 - arXiv:1906.08034 for more!

→ Long term goal: how to choose  $N$ ?

(tentative) **Right after jamming**, and do **ensemble averaging!**

3. How does the test error scale with  $P$ ?

check Spigler et al. '19 - arXiv:1905.10843 !