

Discussion

Neural Tangent Kernel: Convergence and Generalization in Neural Networks

Arthur Jacot

École Polytechnique Fédérale de Lausanne
arthur.jacot@netopera.net

Franck Gabriel

Imperial College London and École Polytechnique Fédérale de Lausanne
franckrgabriel@gmail.com

Clément Hongler

École Polytechnique Fédérale de Lausanne
clement.hongler@gmail.com

NIPS 18

Scaling description of generalization with number of parameters in deep learning

Mario Geiger * ¹, Arthur Jacot * ¹, Stefano Spigler ¹, Franck Gabriel ¹, Levent Sagun ¹, Stéphane d'Ascoli ², Giulio Biroli ², Clément Hongler ¹, and Matthieu Wyart ¹

arxiv 19

Why does deep learning work?

- when can one fit the data (not stuck bad minimum)?

crank up the number of parameters

- Why does it generalize well, even when the number of parameters is large?

Generalization keeps improving with number of parameters...

MENU:

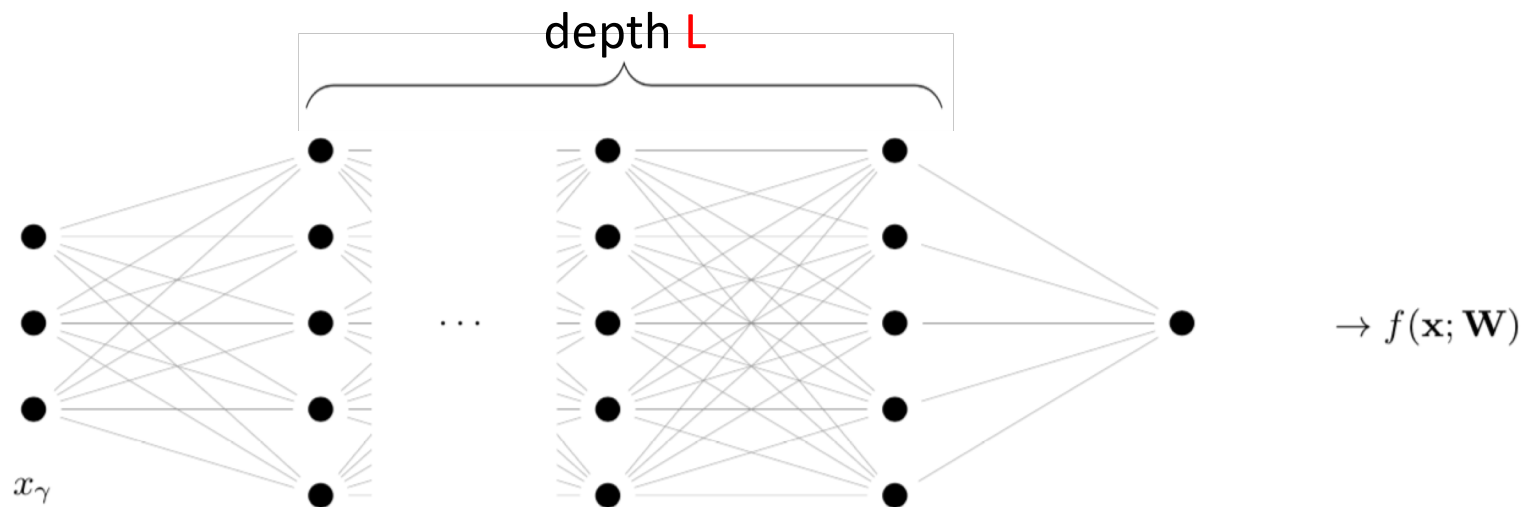
1/ Quantification of evolution of generalization with number of parameters

2/ Neural Tangent Kernel (NTK)

3/ NTK and generalization as number of parameters becomes asymptotically large

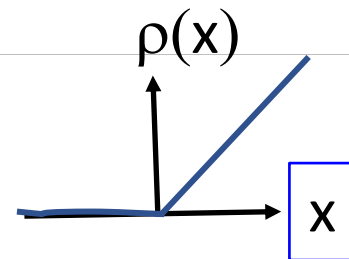
Set-up

- binary classification task, P training data $\{\mathbf{x}_i, y_i = \pm 1\}$
- Deep net $f_{\mathbf{W}}(\mathbf{x}_i)$ with N parameters, width h ($N \sim h^2$)



$$a_\beta = \rho\left(\sum_{\alpha \in \text{previous layer}} W_{\alpha, \beta} a_\alpha - B_\beta\right)$$

ρ : non-linear function



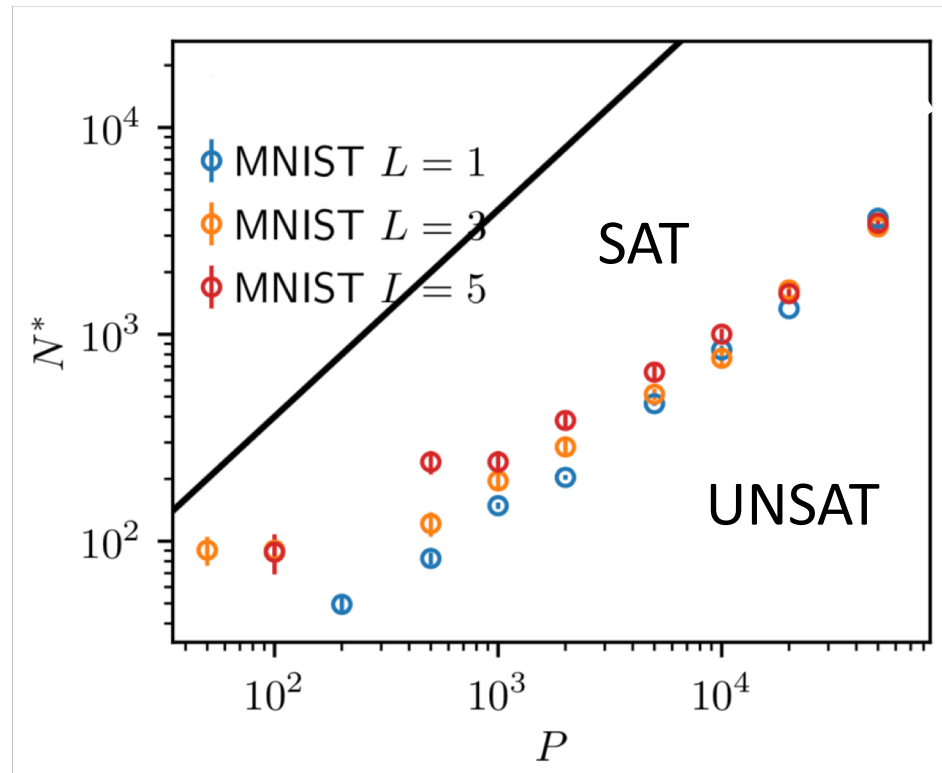
Learning

- Learning: gradient descent in loss function $\mathcal{L} = \frac{1}{P} \sum_i^P l_i(f_{\mathbf{W}}(x_i))$
- Hinge Loss: $l_i(f_{\mathbf{W}}(x_i)) = (f_{\mathbf{W}}(x_i)y_i - 1)^2$ if $f_{\mathbf{W}}(x_i)y_i < 1$
otherwise 0
- $\mathcal{L} = 0 \Leftrightarrow f_{\mathbf{W}}(x_i)y_i > 1 \forall i$ satisfiability problem
- Dynamics stops in the SAT phase.
- Expect transition at some $N^*(P)$

Empirical tests: MNIST (parity)

Geiger et al., arxiv 180909349,

- $6 \cdot 10^4$ images of digits

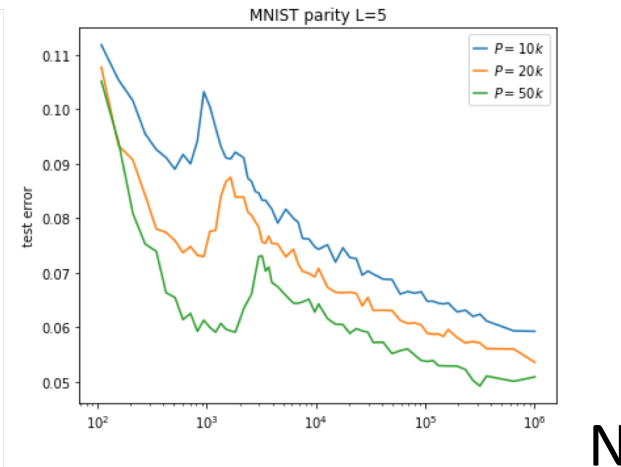
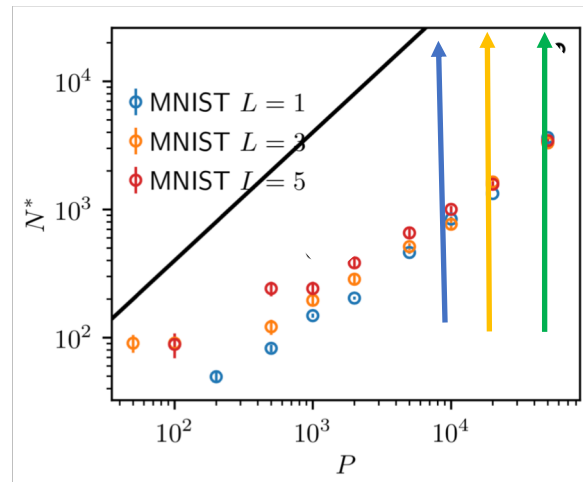


- position of transition depends on dynamics (GD, adams, fire...)

Generalization

Spigler et al. arxiv 1810.09665,

test error



*see also Advani and Saxe 17,
Neal et al. 18, Neyshabur et al., 15, 17.*

2 interesting asymptotic regimes:

- peak at the SAT-UNSAT transition
- performance improves with N in the SAT phase???

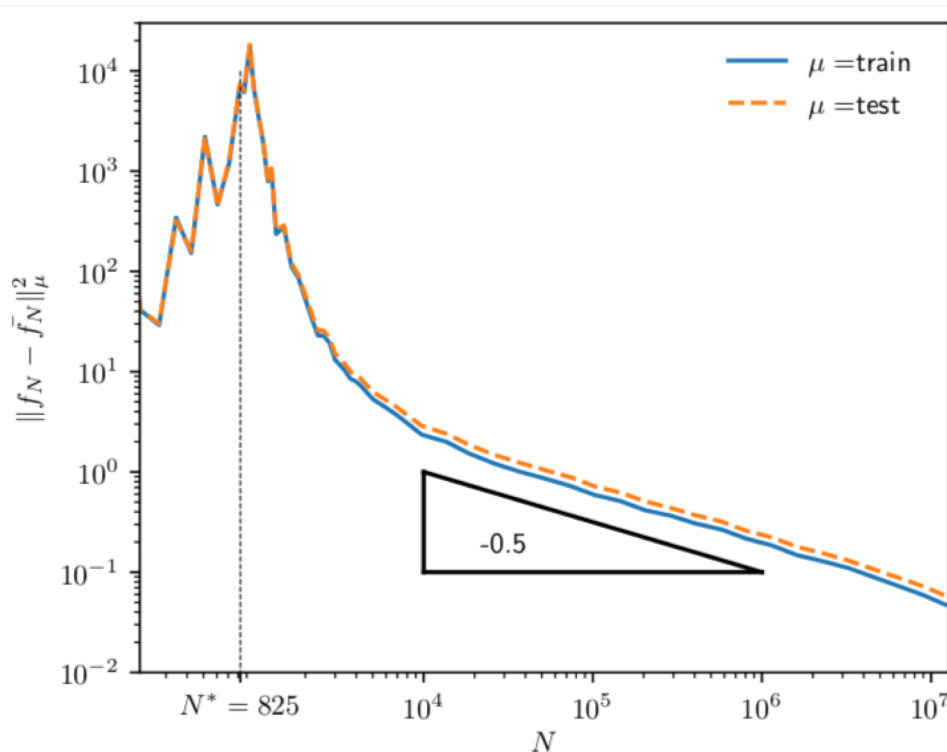
talks Rakhlin, Srebro: increased regularization with N

Quantitative description? importance of $N=\infty$

Quantifying fluctuations induced by initialization

- fixed data set, output function f stochastic due to initialization
- This stochasticity is reduced as N grows [Neal et al. arxiv 1810591](#)

\bar{f}_N : ensemble average of f_N on (20) initial conditions



$$\|f\|_\mu^2 = \int d\mu(x) f(x)^2$$

$$\|f_N - \bar{f}_N\|_\mu \sim N^{-1/4}$$

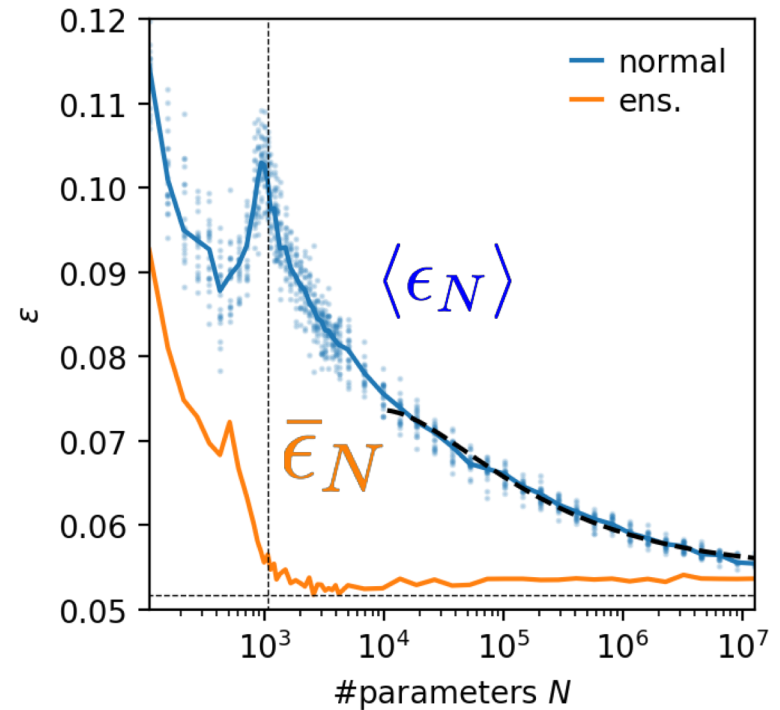
(to be explained by NTK)

Test and practical consequences

Geiger et al., arxiv 1901.01608

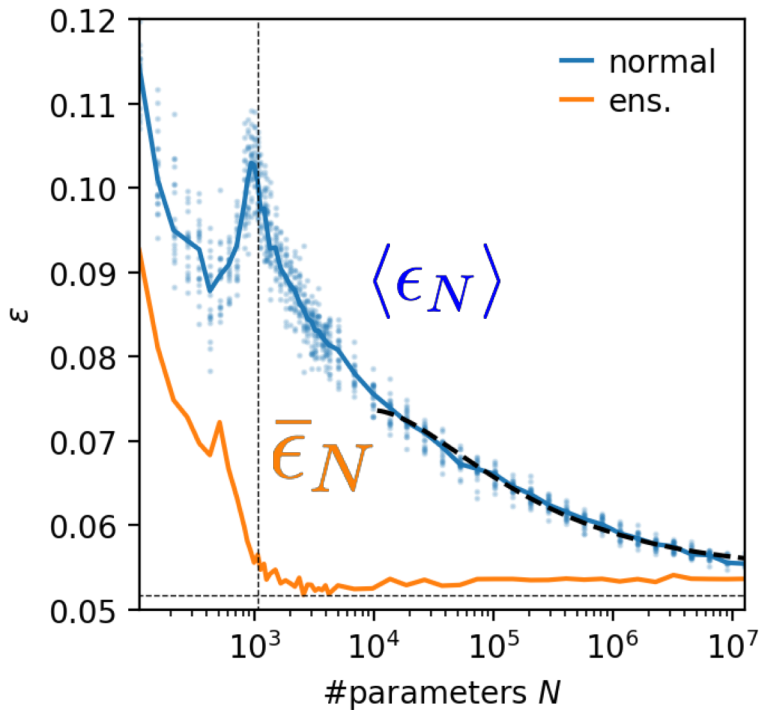
- Reduce fluctuations by averaging

$\bar{\epsilon}_N$: test error of \bar{f}_N

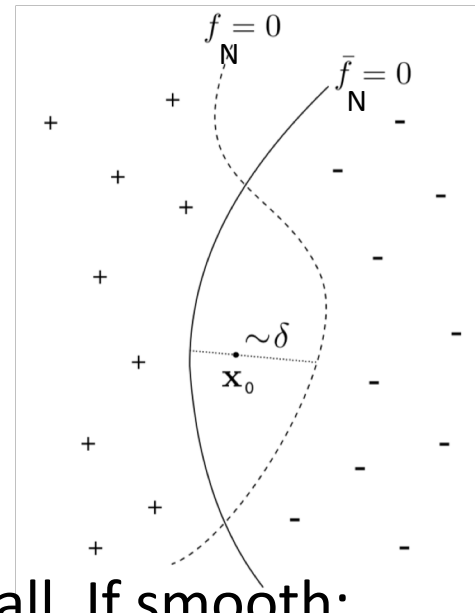


- test error becomes nearly flat for $N > N^*$, optimal near N^*
- *Best procedure: ensemble average near SAT-UNSAT transition!!!*

Scaling argument for generalization error



- seek to compute $\langle \epsilon_N \rangle - \bar{\epsilon}_N$
using $\delta f_N = f_N - \bar{f}_N$ very small



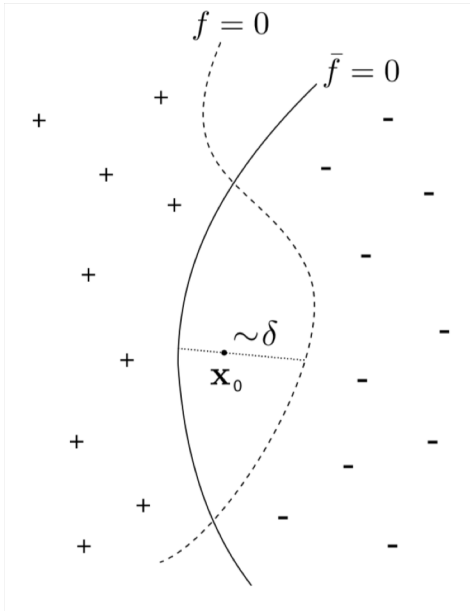
decision
boundary

- signed distances $\delta(x)$ becomes small. If smooth:

$$\delta(x) = \frac{\delta f_N(x)}{\|\nabla \bar{f}_N(x)\|} + \mathcal{O}(\delta f_N^2)$$

$$\begin{cases} \delta(x) \sim \|\delta f_N\|_\mu & (\text{NTK}) \\ \langle \delta(x) \rangle \sim \|\delta f_N\|_\mu^2 \end{cases}$$

Scaling argument for generalization error

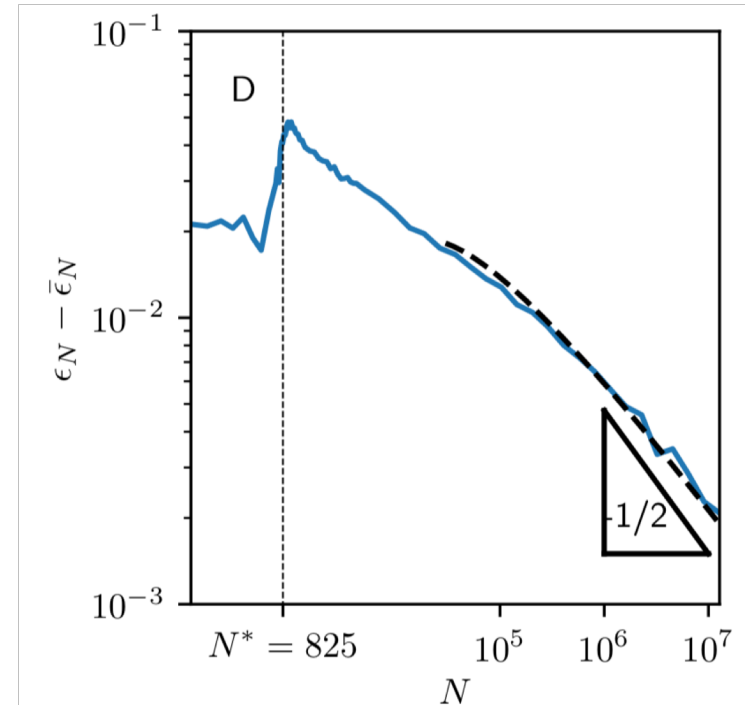


$$\Delta\epsilon = \int_B dx^{d-1} \left[\frac{\partial\epsilon}{\partial\delta(x)}\delta(x) + \frac{1}{2} \frac{\partial^2\epsilon}{\partial^2\delta(x)}\delta^2(x) + \mathcal{O}(\delta^3(x)) \right]$$

$$\langle\Delta\epsilon\rangle = c_0\|\delta f\|^2 + \mathcal{O}(\|\delta f\|^3)$$

expect $c_0 > 0$ if $\bar{\epsilon}$ small

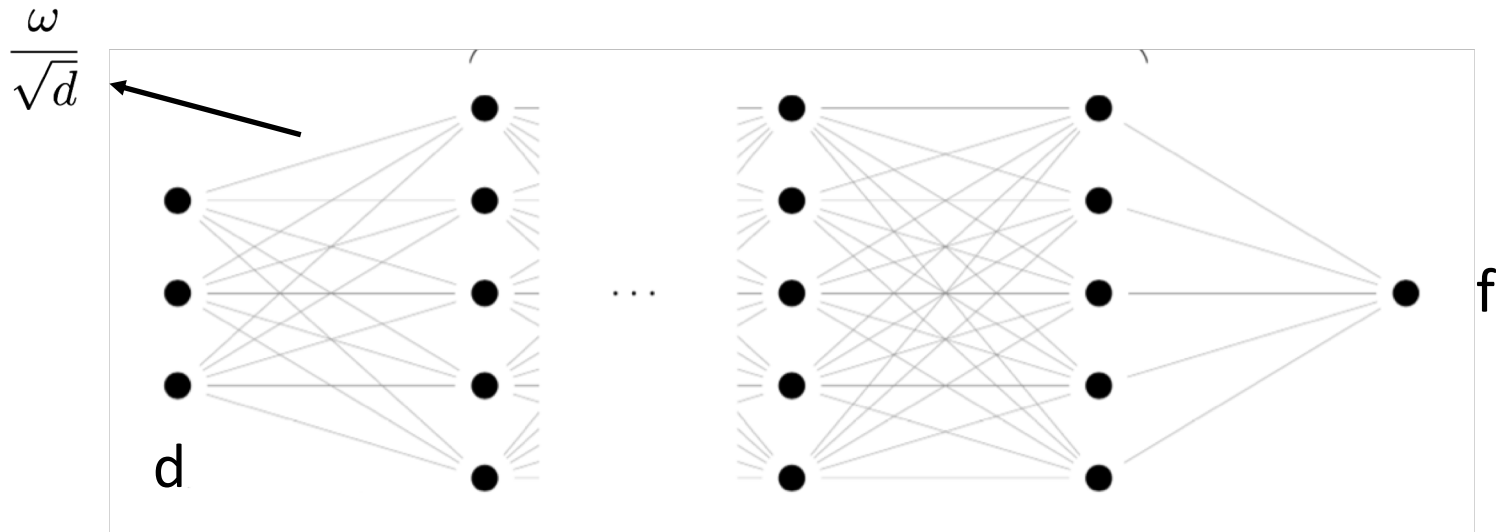
$$\langle\epsilon_N\rangle - \bar{\epsilon}_N \sim \|\bar{f}_N - f_N\|^2 \sim 1/\sqrt{N}$$



Propagation in infinitely wide nets at $t=0$

Neal 96, williams 98, Lee et al 18, Ganguli et al.

set-up: initialization iid weights = $\frac{\omega}{h^{1/2}}$ where $\omega \sim \mathcal{N}(0, 1)$



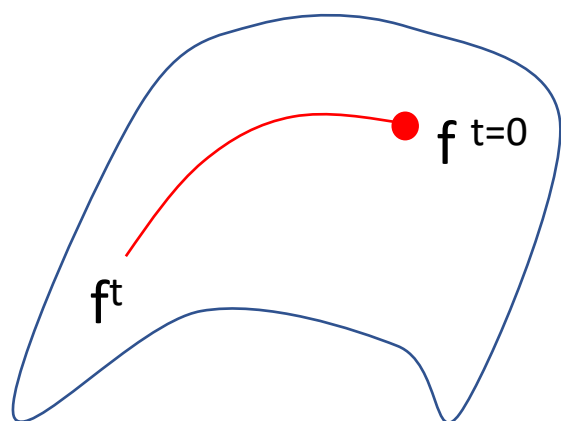
- Non-trivial limit for propagation, pre-activation $\alpha \sim 1$ and $f \sim 1$
- pre-activation and output are iid gaussian processes as $h \rightarrow \infty$

$$\langle \alpha_i^\ell(x) \alpha_j^\ell(x') \rangle = \delta_{i,j} \Sigma_\ell(x, x')$$

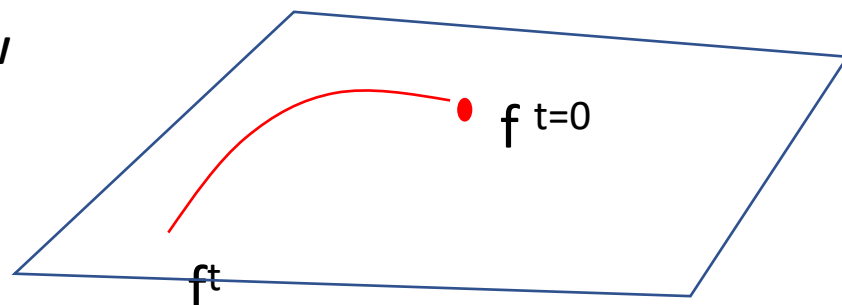
- recursive relation for Σ_ℓ
- results: $\partial f / \partial \alpha \sim 1 / \sqrt{h}$, between hidden neurons $\partial f / \partial \omega \sim 1 / h$

Learning: Neural Tangent Kernel

Jacot, Gabriel, Hongler NIPS 18



manifold f_w



small h : $\partial f / \partial w$ evolves

large h : $\partial f / \partial w$ fixed



“lazy learning”:

- weights change a little bit $\omega^t - \omega^0 \sim 1/h$
- sufficient to change f (positive interference)
- does not change $\partial f / \partial w$

Results

$$\mathcal{L} = \frac{1}{P} \sum_i^P l_i(f_{\mathbf{w}}(x_i)) \quad \text{gradient descent}$$

$$\frac{df(x)}{dt} = -\frac{1}{P} \sum_{i=1}^P \Theta_N^t(x_i, x) l'_i(f(x_i))$$

$$\Theta_N^t(x_i, x) = \sum_{\omega} \frac{\partial f^t(x_i)}{\partial \omega} \frac{\partial f^t(x)}{\partial \omega} \quad \text{useless in general...}$$

Theorem 1: at $t=0$, kernel does not depend on initialization at large N

$$\boxed{\lim_{N \rightarrow \infty} \Theta_N^{t=0}(x_i, x) = \Theta_{\infty}(x_i, x)} \quad (\text{recursive proof})$$

scaling makes sense. All layers contribute to Kernel.

Results

- Theorem 2: kernel does not depend on time

$$\lim_{N \rightarrow \infty} \Theta_N^t(x_i, x) = \Theta_\infty(x_i, x)$$

$$\frac{df(x)}{dt} = -\frac{1}{P} \sum_{i=1}^P \Theta_\infty(x_i, x) l'_i(f(x_i))$$

deep learning equivalent to kernel learning as $N \rightarrow \infty$

Self-consistency:

- learning occurs on time $t \sim \mathcal{O}(1)$
- on that time scale, weights change very little

$$\omega^t - \omega^{t=0} \sim d\mathcal{L}/d\omega \sim \partial f / \partial \omega \sim 1/h$$

Results

- Theorem 3: Dynamics find global minimum of the loss if loss l_i convex and activation function non-polynomial

Gram matrix $\Theta_{\infty}(x_i, x_j)$ positive definite

$$\frac{df(x)}{dt} = -\frac{1}{P} \sum_{i=1}^P \Theta_{\infty}(x_i, x) l'_i(f(x_i))$$

- Result 4: smoothness of $f^t(x)$ can be deduced

$$f^t(x) = f^{t=0}(x) + \sum_{i=1}^P c_i(t) \Theta_{\infty}(x, x_i)$$

Finite N Geiger et al. 19, Jacot et al 19

- Fluctuations of $\Theta_N^{t=0}$ go as $1/\sqrt{h} \sim N^{-1/4}$
- evolution in time much smaller $\theta_N^t - \theta_N^{t=0} \sim 1/\sqrt{N}$

$$\frac{df(x)}{dt} = -\frac{1}{P} \sum_{i=1}^P \Theta_N^t(x_i, x) l'_i(f(x_i))$$

- leads to fluctuations of similar magnitude for output function (proof mean square loss)

$$\|f_N - \bar{f}_N\|_{\mu} \sim N^{-1/4}$$

Is learning features useful?

- neurons pattern of activity barely changes as $N \rightarrow \infty$

$$\alpha^t(x) - \alpha^{t=0}(x) \sim 1/\sqrt{h}$$

- success of deep learning believes to be associated with the emergence of good features....
- Small effect at best on MNIST...

More data needed.

