

①

Lecture 8 : Infinite-width limits

Previous lectures: as a first approximation of highly overparametrized neural networks, we can consider infinite width limits

→ The hope is that:

① There will be some simplifications that result in a compact description of the dynamics which is amenable to analysis.

② It's a good approximation of neural networks used in practice.

③ Conversely, it can inform us how to scale NNs
(Common practice: grid-search over hyperparameters
→ if we had a more principled approach, we could save time and money + now models are too big to do grid-search directly)

Right now, we saw two different limits:

* NTK scaling: As $M \rightarrow \infty$ with multi-layer NNs
 SGD concentrates on kernel method

$$\text{NTK}(\alpha_1, \alpha_2) = h_L(\langle \alpha_1, \alpha_2 \rangle) \rightarrow \text{closed form}$$

easy to compute

* MF scaling: As $M \rightarrow \infty$, 2-layer NNs
 SGD concentrates on Wasserstein gradient flow

$$f = \int \sigma(x; \theta) p(\theta) \quad \partial_t p_t = \nabla_{\theta} \cdot [p_t \nabla_{\theta} U(\theta; p_t)]$$

→ Very different behavior between the 2 limits

What about limits with other scalings?

Many papers have explored how to scale different hyperparameters with width, and how to "transfer" these hyperparameters between models

Today's goal: briefly present a very simple framework
 "ABC"-parametrization [Yang, Hu, 2020]

(3)

Fully-connected multi-layer NNs + SGD
and knock the scaling of 3 sets of hyperparameters as
width $\rightarrow \infty$.

Model: L-hidden layer NNs

all layers are width M

* input $x \in \mathbb{R}^d$

* weight matrices: $W^1 \in \mathbb{R}^{M \times d}$

$W^2, W^3, \dots, W^L \in \mathbb{R}^{M \times M}$

$W^{L+1} \in \mathbb{R}^{1 \times M}$

* non-linearity $\sigma: \mathbb{R} \rightarrow \mathbb{R}$

Compute the output iteratively: $h^1(x) = W^1 x$

$l = 1, \dots, L$

"reactivation"

$\tilde{z}^l(x) = \sigma(h^l(x))$

$\in \mathbb{R}^M$

$h^{l+1}(x) = W^{l+1} \tilde{z}^l(x)$

$\in \mathbb{R}^M$

$$f(x) = W^{L+1} \tilde{z}^L(x) = W^{L+1} \sigma(h^L(x))$$

$$= W^{L+1} \sigma(W^L \tilde{z}^{L-1}(x))$$

$$= W^{L+1} \sigma(W^L \phi(h^{L-1}(x)))$$

$$= W^{L+1} \circ \sigma \circ W^L \circ \dots \circ \sigma \circ W^1 x$$

(4)

We train this NN with batch-1 SGD

Some loss: $l(y, f(x; \Theta))$

$$\begin{aligned} \Theta_{k+1} &= \Theta_k - \eta \nabla_{\Theta} l(y_k, f(x_k; \Theta_k)) \\ &= \Theta_k - \eta l'(y_k, f(x_k; \Theta_k)) \nabla_{\Theta} f(x_k; \Theta_k) \end{aligned}$$

$$f_k(x) := f(x; \Theta_k)$$

$$\rightarrow \text{some with } z_k^l(x) \quad h_k^l(x)$$

$$X_k = l'(y_k, f(x_k; \Theta_k))$$

We want to consider $M \rightarrow \infty$ limit in this dynamic and how we should scale 3 sets of hyperparameters

$$(a) \underline{\text{Layer normalization}}: W^l = \frac{1}{M^{al}} \bar{W}^l$$

$\hookrightarrow \bar{W}^l$ are the trainable parameters

(5)

(b) Initialization: \bar{W}_0^l has iid $N(0, \frac{1}{M^{2\alpha_l}})$

(c) Learning rate: $\eta = \frac{1}{M^c}$

$$\bar{W}_{k+1}^l = \bar{W}_k^l - \frac{1}{M^c} \sum_k X_k \nabla_{\bar{W}^l} f_k(x_k)$$

Set of hyperparameters $\{\alpha_l, b_l\}_{l=1}^{L+1} \cup \{c\}$

Which choice of abc leads to meaningful limits?

What is the behavior of these different limits?

Rank: abc parametrization is degenerate

- At initialization: $W_0^l = \frac{\bar{W}^l}{M^{\alpha_l}} \sim N(0, \frac{1}{M^{2(\alpha_l + b_l)}})$

↳ For $\alpha_l + b_l = \text{cte}$ NN at initialization is the same

6

- When updating:

$$h_{k+1}^l(x) = W_{k+1}^l \tilde{z}_{k+1}^l(x)$$

$$= \frac{1}{M^{al}} \left(W_k^l - \frac{\eta}{M^c} X_k \nabla_{W^l} f_k(x_k) \right)$$

$$= h_k^l(x) - \frac{1}{M^{c+2al}} X_k \nabla_{W^l} f_k(x_k)$$

So NN during SGD is invariant under the following reparametrization

$\forall \theta \in \mathbb{R}$,

$$a_l \leftarrow a_l + \theta$$

$$b_l \leftarrow b_l - \theta$$

$$c \leftarrow c - 2\theta$$

} invariance at initialization
} invariance updates

What can go wrong when taking $M \rightarrow \infty$?

- $h_k^l(x) \rightarrow \infty$ (e.g. c too small)
learning rate too large
 \hookrightarrow blow-up of preactivation
 "parametrization is unstable"
- $f_k(x) = f_0(x)$ (e.g., c too large)
learning rate too small
 \hookrightarrow network does not change in finite time
 "parametrization is trivial"

Example: 1-hidden layer ($L=1$)

(8)

For simplicity $x \in \mathbb{R}$ $W \in \mathbb{R}^{M \times 1}$ $V \in \mathbb{R}^{1 \times M}$

$$f(x) = V \tilde{z}(x) \quad \tilde{z}(x) = \sigma(h(x)) \quad h(x) = Wx$$

$$V = \frac{V}{M^{\alpha_2}} \quad W = \frac{W}{M^{\alpha_1}}$$

* NT parametrization: $\alpha_1 = 0 \quad b_1 = 0$
 $\alpha_2 = \frac{1}{2} \quad b_2 = 0 \quad) \frac{1}{M} \sum v_j \sigma(w_j^T x)$
 $c = 0$

* MF parametrization: $\alpha_1 = b_1 = b_2 = 0$
 $\alpha_2 = 1 \quad) \frac{1}{M} \sum v_j \sigma(w_j^T x)$
 $c = -1 \quad \rightarrow \eta M$

Derivatives: $\nabla_{\tilde{z}} f(x) = V$ $\nabla_h f(x) = V \circ \sigma'(h(x))$

$$\nabla_V f(x) = \frac{1}{M^{\alpha_2}} \tilde{z}(x)$$

$$\nabla_W f(x) = \frac{1}{M^{\alpha_1}} V \circ \sigma'(h(x)) \propto$$

(9)

Update after one step:

$$\Delta h_1(x) = h_1(x) - h_0(x) = -\frac{1}{M^{2a_1+c}} x_0 \otimes [V_0 \odot \sigma(h_0(x))]$$

$$\begin{aligned}\Delta f_1(x) &= V_0 \Delta \beta_1(x) + \Delta V_1 \beta_1(x) \\ &= V_0 \Delta \beta_1(x) - \frac{1}{M^{2a_2+c}} \beta_1(x_0)^T \beta_1(x)\end{aligned}$$

NTP & MFP: $\beta_0(x)$ has coordinates $\Theta(1)$

Feature evolution?

$$\underline{\text{NTP: }} \Delta h_1(x) = \Theta(1) \cdot \underbrace{V_0 \odot \Theta(1)}_{\Theta\left(\frac{1}{M}\right)}$$

each coordinate of preactivation move $\Theta\left(\frac{1}{M}\right)$

$$\underline{\text{MFP: }} \Delta h_1(x) = \Theta(M) \cdot \Theta\left(\frac{1}{M}\right) \odot \Theta(1) \\ = \Theta(1)$$

each coordinate of preactivation move $\Theta(1)$

Feature "kernel" evolution:

$$F_k(x, x') = \frac{1}{M} \mathbf{z}_k(x)^\top \mathbf{z}_k(x')$$

$$\lim_{M \rightarrow \infty} F_k(x, x') = \lim_{M \rightarrow \infty} F_0(x, x') \quad \text{NTP}$$

the feature kernel does not evolve

$$\lim_{M \rightarrow \infty} F_k(x, x') \neq \lim_{M \rightarrow \infty} F_0(x, x') \quad \text{MFP}$$

the kernel changes!

General picture

Symmetries of abc parametrized:

$$\forall \theta \in \mathbb{R}, \quad a_l \leftarrow a_l + \theta \quad b_l \leftarrow b_l - \theta \quad c \leftarrow c - 2\theta$$

result in the same limiting model.

Feature movements scaling: $\Delta \tilde{\gamma}_k^L(x) = \odot(M^{-n})$

$$n := \min(a_{L+1} + b_{L+1}, 2a_{L+1} + c) + c - 1$$

$$+ \min_{l=1}^L [2a_l + \mathbf{1}(l=1)]$$

E.X. NTP $n = \frac{1}{2}$ MFP $n = 0$

\hookrightarrow to avoid blow-up $n \geq 0$ feature learning $n = 0$

(Proof: look at one step, easy)

Stable abc parametrization:

It is stable IFF

- (1) prectives[°] at initialization are $O(1)$
and $f_0 = O(1)$

$l=1$	$2 \leq l \leq L$	$l=L+1$
$a_l + b_l = 0$	$a_l + b_l = \frac{1}{2}$	$a_{L+1} + b_{L+1} \geq \frac{1}{2}$

- (2) features don't blowup $\Delta a_i^l = O(1)$

$$n \geq 0$$

- (3) f_k doesn't blowup

$$2a_{L+1} + c \geq 1 \quad a_{L+1} + b_{L+1} + n \geq 1$$

Proof: easy (do one step)

(13)

Nontivial abc parametrizations

A stable parametrization is non-trivial

$$\text{IFF } a_{L+1} + b_{L+1} + r = 1 \quad \text{or} \quad 2a_{L+1} + c = 1$$

Proof: $\Delta f_1(x) = \Delta W_1^{L+1} \mathbf{z}_1^L(x) + W_0^{L+1} \Delta \mathbf{z}_1^L(x)$

$$\Theta\left(\frac{1}{M^{2a_{L+1}+c}}\right) \cdot \Theta(1) \quad \Theta\left(\frac{1}{M^{a_{L+1}+b_{L+1}}}\right) \Theta(m^{-r})$$

$$\lim_{M \rightarrow \infty} \Delta f_1(x) \neq 0 \quad \text{IFF} \quad \text{one of the 2 terms} \neq 0$$

Dynamical dichotomy

Def: [Feature learning] abc parametrization admits feature learning if as $M \rightarrow \infty$,

$$\boxed{\frac{1}{M} (\mathbf{z}_k^L(x)^T \mathbf{z}_k^L(x') - \mathbf{z}_0^L(x)^T \mathbf{z}_0^L(x')) = \mathcal{S}(1)}$$

for some time $k \geq 1$ and input x, x'

(That is $\Delta \mathcal{Z}_n^L(n)$ has $\Sigma(1)$ coordinates)

Def: [Kernel regime] abc parametrization is in the kernel regime iff for any $k \geq 0$,

$$\lim_{M \rightarrow \infty} f_{k+1}(x) = \lim_{M \rightarrow \infty} f_k(x) - \gamma K(x, x_k) l'(y_k, f_k(x_k))$$

where $K(x, x') = \lim_{M \rightarrow \infty} \langle \nabla_{\mathbb{G}} f(x), \nabla_{\mathbb{H}} f(x') \rangle$

Thm: A non-trivial stable abc parametrization

* admits feature learning iff $r=0$

* is in kernel regime iff $r>0$

Remark: For $r=0$ and $r>0$ plenty of dynamics possible

(e.g. some layers don't move)

• Higher-order Taylor correct^o to NTK are not valid ∞^k width limits

Examples

NTP: $a_l = \begin{cases} 0 & l=1 \\ \frac{1}{2} & l \geq 2 \end{cases}$ $b_l = 0$ $c = 0$

Then $n = \frac{1}{2}$ kernel regime

RFP: $c = 1$ $a_{L+1} = 0$ $b_{L+1} = \frac{1}{2}$

$a_l = 0$ $b_1 = 0$

$b_l = \frac{1}{2}$ $2 \leq l \leq L$

Then $n = \frac{1}{2}$
 \hookrightarrow only train last layer ("random feature" model)

MFP: $L = 1$ $a_1 = 0$ $b_1 = 0$

$a_2 = 1$ $c = -1$

Then $n = 0$ feature learning regime

SP: "standard parametrization" available as default setting in PyTorch

It has $a_l = 0$ for all l

(16)

$b_1 = 0$ and $b_l = \frac{1}{2}$ for $l \geq 2$

$$c = 0$$

↳ learning rate is too large, not a stable abc parametrization. To make it stable, take $c=1$

Then $\eta = \frac{1}{2}$ kernel regime

" μP " : "maximal update parametrization"

↳ parametrization that allows each parameters to be updated maximally without blowing up

i.e. $\Delta h_k^l(x) = \textcircled{n}(1)$ entries for all $l = 1, \dots, L$

$$\textcircled{n}(n^{-\eta_l}) \quad \eta_l = \min(\alpha_{L+1} + b_{L+1}, 2\alpha_{L+1} + c)$$

$$+ c - 1 + 2\alpha_l + \textcircled{1}(l=1)$$

$$\eta = \min_{1 \leq l \leq L} \eta_l \Rightarrow \mu P \quad \eta_l = 0 \text{ for all } l$$

This gives: $a_l = \begin{cases} -\frac{1}{2} & l=1 \\ 0 & 2 \leq l \leq L \\ \frac{1}{2} & l=L+1 \end{cases}$

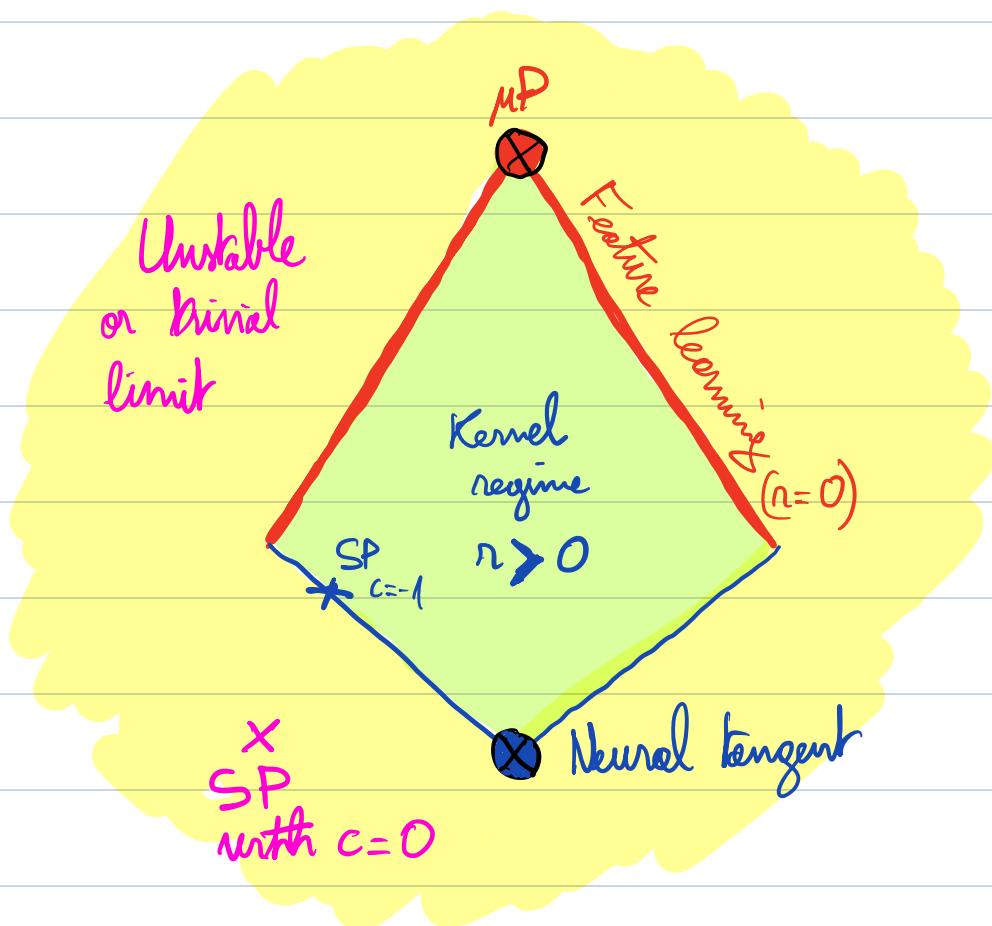
$$b_l = \frac{1}{2}$$

$$c = 0$$

This coincides with MFP for $L=1$

Summary

In space $\{\alpha_2, b_2\} \cup \{c\}$ non binar stable parametrization form a polytop



Application: • μP seem to work well

- "one-shot transfer learning": Let's we want to train a model with 10^6 parameters ($M = 10^6$) how to choose learning rate / layer normalization / variance of initial

Take a smaller model with width M_0

(e.g. $M_0 = 10^5$, i.e. $\approx 100M$ parameters)

do grid search and find optimal η_0 , $N_{\ell,0}$, $\nu_{\ell,0}$

then for model with width M_1 , use directly

$$\eta_1 = \eta_0 \left(\frac{M_0}{M_1} \right)^c$$

$$N_{\ell,1} = N_{\ell,0} \left(\frac{M_1}{M_0} \right)^{\alpha_\ell}$$

$$\nu_{\ell,1} = \nu_{\ell,0} \left(\frac{M_0}{M_1} \right)^{2b_\ell}$$

e.g. take μP

[Greg Yang et al., 22] Using this idea with μP , took GP-3

40M parameters \rightarrow 6.7B parameters

\hookrightarrow the model they get 1) show better performance
 2) only 7% of total pretraining cost

[OpenAI, Anthropic, xAI: have empirical heuristics to scale hyperparameters ... but it's not μP ...]