# Lecture VI – Wide multilayer networks and the Neural Tangent Kernel (NTK)

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

October, 2023

The Neural Tangent Kernel (NTK)

Wide networks and Gaussian Processes

The NTK is constant during training Example – regression and  $\mathcal{L}_{LS}$ 

Wide and deep networks and classification

# Notation

- Neural network predictor  $f(x; \theta)$ , where  $x \in \mathbb{R}^d$
- For each layer l = 1 : L of dimension  $m_l$ , with  $x^0 \equiv x$ , and  $z^L \equiv f(x)$

$$z^{l+1} = W^{l+1}x^{l} + b^{l+1} \qquad x^{l+1} = \phi(z^{l+1})$$
(1)

Here  $x^{l,l+1}, z^{l+1}, b^{l+1}$  are column vectors  $W^{l+1}$  is a  $m_{l+1} \times m_l$  matrix,  $\phi()$  is the non-linearity/activation function.

The weights

$$W_{ij}^{l} = \sigma_w w_{ij}^{l} / \sqrt{m_l}, \qquad b_j^{l} = \sigma_b \beta_j^{l}, \quad ext{Known as NTK parametrization}$$
 (2)

- ▶ Parameter vector  $\theta$  = vector { $w^{1:L}$ ,  $\beta^{1:L}$ }  $\in \mathbb{R}^p$  initialized i.i.d. ~ N(0, 1)
- $\sigma_{w,b}$  are fixed hyper-parameters,  $1/\sqrt{m_l}$  normalizes the expected norm of  $W^l$  columns • Loss  $\mathcal{L}(y, f)$
- We want to analize the behavior of this network f() at initialization and during training, when m<sub>1:L</sub> very large
- Three approximations help analysis
  - (A1) continuous time training, called gradient flow
  - (A2)  $m_{1:L} \to \infty$  in the wide limit, we can apply the Central Limit Theorem (CLT), and Gaussian Processes
  - (A3) parameters  $\theta$  do not change much during training, i.e.  $\theta_t \theta_0$  is small

#### The Gradient Flow

• Assume training by gradient descent on  $\hat{\mathcal{L}} = \sum_{i} \mathcal{L}(y^{i}, f(x^{i}))$ 

• The gradient of  $\hat{\mathcal{L}}$ 

$$\nabla_{\theta} \hat{\mathcal{L}} = \sum_{i} \frac{\partial \mathcal{L}}{\partial f} (y^{i}, f(x^{i}; \theta)) \nabla_{\theta} f(x^{i}, \theta) = \nabla_{\theta} f_{\mathcal{D}} \nabla_{f} \mathcal{L}_{\mathcal{D}} \qquad \in \mathbb{R}^{p}$$
(3)

where  $\nabla_f \mathcal{L}_{\mathcal{D}} = [\frac{\partial \mathcal{L}}{\partial f}(y^i, f(x^i; \theta))]_{i=1:n} \in \mathbb{R}^n$ ,  $\nabla_\theta f_{\mathcal{D}} = [\nabla_\theta f(x^i, \theta)]_{i=1:n} \in \mathbb{R}^{p \times n}$ Assume (A1) gradient descent with infinitezimal time steps. In other words, the parameters evolve by an ordinary differential equation

$$\dot{\theta} = -\eta \nabla_{\theta} f_{\mathcal{D}} \nabla_{f} \mathcal{L}_{\mathcal{D}} \quad \in \mathbb{R}^{p}$$

$$(4)$$

$$\dot{f} = \sum_{j=1}^{p} \frac{\partial f}{\partial \theta_j} \frac{\partial \theta_j}{\partial t} = (\nabla_{\theta} f)^T \dot{\theta} \in \mathbb{R}$$
(5)

$$\dot{f}_{\mathcal{D}} = -\eta \underbrace{(\nabla_{\theta} f_{\mathcal{D}})^{T} \nabla_{\theta} f_{\mathcal{D}}}_{G} \nabla_{f} \mathcal{L}_{\mathcal{D}} \in \mathbb{R}^{p}$$
(6)

•  $G \equiv \nabla_{\theta} f_{\mathcal{D}}^T \nabla_{\theta} f_{\mathcal{D}} \equiv \kappa(X, X)$  is a Gram matrix! • Therefore, we define the Neural Tangent Kernel (NTK) by

$$\kappa(\mathbf{x},\mathbf{x}') = \nabla_{\theta} f(\mathbf{x};\theta)^{\mathsf{T}} \nabla_{\theta} f(\mathbf{x}';\theta)$$
(7)

### Gradient flow and NTK - summary

$$\begin{split} \dot{\theta} &= -\eta \nabla_{\theta} f_{\mathcal{D}} \nabla_{f} \mathcal{L}_{\mathcal{D}} \quad \in \mathbb{R}^{p} \\ \dot{f}_{\mathcal{D}} &= -\eta G \nabla_{f} \mathcal{L} \quad \in \mathbb{R}^{p} \\ \kappa(x, x') &= \nabla_{\theta} f(x)^{T} \nabla_{\theta} f(x') \end{split}$$

•  $f_X$ ,  $\nabla_{\theta} f_X$ , G depend only on the inputs X,  $\theta$ 

- ▶  $\nabla_f \mathcal{L}$  depends only on the correct outputs Y, and predicted outputs, i.e. on Y and  $\theta$
- This holds for any predictor! So what is special about neural networks?
- First, we will analyze κ for very wide neural networks with random parameters (e.g. at initialization)
- Then, we will analyze what happens during training under assumption (A3)

# Wide NN's Gaussian Process (GP)

- **\triangleright** This is about  $f_0$ , a NN initialized with Gaussian independent parameters. For simplicity, we denote it as f.
- ► Assume  $\theta^{1:L-1}$  fixed, only  $W^L$ ,  $b^L$  random as in (2) ► Recall  $f(x) = W^L x^{L-1}(x) + b^L$  for any x with  $x^{L-1} \in \mathbb{R}^{m_L}$
- ▶  $f(x) = \text{sum of } m_{l-1} \text{ i.i.d. random variables, hence } f(x) \sim Normal \text{ by CLT, for } m_{l-1} \text{ large}$ Randomness is over weights W<sup>L</sup>, b<sup>L</sup>!!!
- We have E[f(x)] = 0 and

$$Cov(f(x), f(x')) = E[(W^{L}x^{L-1} + b^{L})(W^{L}(x')^{L-1} + b^{L})] = \frac{\sigma_{w}^{2}}{m_{L-1}}(x^{L-1})^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^{T}(x')^{L-1} + \sigma_{b}^{2} \equiv \kappa^{L}(x^{L-1})^{T}(x')^$$

where  $x^{L-1}, (x')^{L-1} \in \mathbb{R}^{m_{L-1}}$  are the outputs of the (L-1)-th layer for inputs x, x' $\blacktriangleright \kappa^{L}$  is a positive definite kernel Exercise Prove this.

- $\blacktriangleright$  f(x) is a random function of x
- The distribution of f(x) defined as above, is called a Gaussian Pocess
- More generally, it can be shown [Jacot, Gabriel, Hongler, NeurIPS 2018] that, when all  $\theta$ parameters are sampled as in (??),  $f_0(x) \sim GP(0, \kappa^L)$
- Q1 What is the kernel  $\kappa^{L}$  of this GP
- Q2 This is all nice, but  $\theta$  changes during training. What can we say about  $\theta_t$ ,  $f_t$  after training?

Q1: Idea.

From (8), for layer I = 1 : L we have

$$\kappa^{l}(\mathbf{x}, \mathbf{x}') = E[z_{j}^{l}(\mathbf{x})z_{j}^{l}(\mathbf{x}')] = \frac{\sigma_{w}^{2}}{m_{l-1}} (\mathbf{x}^{l-1})^{T} (\mathbf{x}')^{l-1} + \sigma_{b}^{2}$$
(9)

with  $x^{l-1} = \phi(z^{l-1})$ . Note also that  $z_j^l$  are i.i.d. so it does not matter which j we choose.

- ▶ In particular,  $\kappa^1(x, x') = \frac{\sigma_w^2}{m_1} x^T x' + \sigma_b^2$  is deterministic
- ... and  $\kappa^{l}$  is random for l > 1.
- However, when  $m_l \to \infty$ ,  $\frac{1}{m_{l-1}} (x^{l-1})^T (x')^{l-1} \to E[*]$
- More specifically, this expectation can be written as

$$E[*] = \int \int \phi(z)\phi(z') Normal(\begin{bmatrix} z\\ z' \end{bmatrix}; 0, \kappa_{x,x'}^{l-1}) dz dz'.$$
(10)

In the above z, z' represent the  $z^{l-1}(x), z^{l-1}(x')$  variables, sampled from the level lNormal distribution, which has covariance given by  $\kappa^{l-1}$ , namely

$$\kappa_{x,x'}^{l-1} = \begin{bmatrix} \kappa^{l-1}(x,x) & \kappa^{l-1}(x,x') \\ \kappa^{l-1}(x',x) & \kappa^{l-1}(x',x') \end{bmatrix}.$$
 (11)

• Hence, the limit of  $\kappa^{l}(x, x')$  when  $m_{1:l} \to \infty$ , is a deterministic kernel for all *l*. [Jacot, Gabriel, Hongler, NeurIPS 2018] derived this recursion (next page).

#### Q1: A recursive expression for the Neural Tangent Kernel

#### [Jacot, Gabriel, Hongler, NeurIPS 2018]

- *L* fixed,  $m \to \infty$
- Simplified expression for  $m_{0:L} = m$ ,  $\sigma_w = \sigma_b = 1$
- Then the NTK  $\kappa \equiv \kappa^L$  is defined recursively by layer

$$\kappa^{1}(x, x') = \Sigma^{1}(x, x'), \quad \Sigma^{1}(x, x') = \frac{1}{m}x^{T}x' + 1$$
 (12)

$$\kappa^{l+1}(x,x') = \kappa^{l}(x,x')\dot{\Sigma}^{l+1}(x,x') + \Sigma^{l+1}(x,x'), \tag{13}$$

with

$$\Sigma^{l+1}(x, x') = L^{\phi}_{\Sigma^{l}(x, x')},$$
(14)

$$\dot{\Sigma}^{l+1}(x, x') = \mathsf{L}^{\phi'}_{\Sigma^{l}(x, x')},\tag{15}$$

and

- $\mathsf{L}^{\phi}_{\Sigma} = E[\phi(X)\phi(X')] \operatorname{with}(X,X') \sim N(0, \begin{bmatrix} \Sigma(X,X) & \Sigma(X,X') \\ \Sigma(X,X') & \Sigma(X',X') \end{bmatrix} (16)$
- ▶ In other words, at level l + 1,  $X \equiv x^l$ ,  $X' \equiv (x')^l$  are sampled from a GP with kernel  $\Sigma^l$ , and  $\Sigma^{l+1}(x, x')$ ,  $\dot{\Sigma}^{l+1}(x, x')$  represent their (scalar) covariance after passing through the non-linearities  $\phi$ ,  $\phi'$  (where  $\phi'$  is the derivative of  $\phi$ )

# Summary so far

Now, we understand the random intialization of wide networks, with *L* layers.

$$f_0 \sim GP(0, \kappa^L) \tag{17}$$

where  $\kappa^{L}$  is a kernel that depends only on  $\phi$  (and  $\sigma^{2}_{b,w}$ )

What next?

- Analysis of training by linearization
- ▶ Then, the NTK limit for  $L \rightarrow \infty$  and its relevance for classification and regression

# The Linearized Network $f^{\text{lin}}$

Notation:  $\theta_{0,t}$ ,  $f_{0,t}$  = parameters, predictor at times 0, t

Here we use (A3), the assumption that the parameters θ change little during training. Extensive evidence supports this assumption.

First order Taylor expansion of  $f_t$  around  $f_0$ 

$$f_t^{\rm lin}(x) = f_0(x) + \nabla_\theta f_0(x)^T (\theta_t - \theta_0)$$
(18)

non-linear in x, linear in  $\theta$ 

$$\nabla_{\theta} f_t^{\text{lin}} = \nabla_{\theta} f_0 \tag{19}$$

$$\kappa(x, x') = \nabla_{\theta} f_0(x)^T \nabla_{\theta} f_0(x') \quad \text{constant during training} \qquad (20)$$
$$G_0 \equiv \kappa_{X,X} \qquad (21)$$

$$\dot{\theta}_t = -\eta \nabla_{\theta} f_0(\mathsf{X})^T \nabla_f \mathcal{L}(\mathsf{Y}, f_t^{\mathrm{lin}}(\mathsf{x}))$$
(22)

$$\dot{f}_t^{\text{lin}}(x) = -\eta \underbrace{\kappa(x, X) G_0}_{t} \nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(x))$$
(23)

depends on  $\theta_0$ 

## NTK during training - empirical evidence



Figure 1: Convergence of the NTK to a fixed limit Figure 2: Networks function  $f_{\theta}$  near convergence for two widths n and two times t. for two widths n and two times t.

#### Linearized Network dynamics for $\mathcal{L}_{LS}$

For example, for L<sub>LS</sub>(y, f) = ½(f − y)<sup>2</sup>, ∇<sub>f</sub>L<sub>LS</sub>(f, y) = f − y. In this case, equations (22),(23) are a linear system and have an analytic solution.

$$\theta_t - \theta_0 = -\nabla_{\theta} f_0(\mathsf{X})^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) \left( f_0(\mathsf{X}) - \mathsf{Y} \right)$$
(24)

$$f_t^{\rm lin}(X) = (I - e^{-\eta G_0 t}) Y + e^{-\eta G_0 t} f_0(X)$$
(25)

$$f_t^{\text{lin}}(x) = \underbrace{\kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) Y}_{\mu(x)} + \underbrace{f_0(x) - \kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) f_0(\underline{k26})}_{\gamma(x)}$$

#### Notes:

• if  $G_0 \succ 0$  then  $e^{-\eta G_0 t} \rightarrow 0$  for  $t \rightarrow \infty$ 

▶ in discrete time t = 0, 1, 3, ... replace  $e^{at}$  with  $(1 - a)^t$ . Sketch of proof:  $\ln(1 - a)^t = t \ln(1 - a) \approx t(-a)$  for a small; therefore  $e^{-at} \approx (1 - a)^t$ .

• 
$$f_t^{\text{lin}}(x) = f_0(x) + \kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) (Y - f_0(X))$$

Exercise Prove (24),(25),(26) from (22),(23)

# Wide and deep neural networks for classification – Basic quantities and assumptions

#### [Radhakrishnan, Belkin, Ulher, 2022]

- ▶ This paper studies the limits of wide neural networks  $m_l \rightarrow \infty$  for all l = 1 : L when the depth  $L \rightarrow \infty$
- It is already known that for regression  $L \rightarrow \infty$  is NOT OPTIMAL
- Since the NTK depends only of the activation function  $\phi$ , the limit shall only depend on  $\phi$  as well.
- In particular, the limit depends on  $\phi$  only through the following

$A = E[\phi(Z)]$	when $z \sim N(0,1)$
$A' = E[\phi'(Z)]$	when $z \sim N(0,1)$
$B = E[(\phi'(Z))^2]$	when $z \sim N(0,1)$

- Classifier  $f(x) = \lim_{L \to \infty} \operatorname{sgn} Y G^{-1} \kappa^{L}(X, x)$  with  $G = [\kappa^{L}(x^{i}, x^{j})]_{i,j=1:n}$ .
- Additional assumptions
  - Data  $X \subseteq S^d_+$ , vectors of norm 1 with all entries  $\geq 0$ .
  - Simplifying assumptions on NTK parameters (e.g.  $\sigma_w = \sigma_b = 1$ )

## Case $A \neq 0$ : Networks implement majority vote

Theorem (Proposition 1 in [Radhakrishnan, Belkin, Ulher, 2022]) If there is a function  $0 < c(L) < \infty$  so that

$$\lim_{L \to \infty} \frac{\kappa^{L}(x, x')}{c(L)} = c_1 > 0 \text{ for any } x \neq x', \text{ and } \lim_{L \to \infty} \frac{\kappa^{L}(x, x)}{c(L)} \neq c_1,$$
(27)

then

$$\lim_{L \to \infty} f(x) = \operatorname{sgn} \sum_{i=1}^{n} y^{i} \qquad MAJORITY \ CLASSIFIER$$
(28)

• What  $\phi$ 's satisfy theorem? ReLU, all  $\phi$  with  $B \neq 1$ .

Case A = A' = 0: Networks implement 1-nearest neighbor

Theorem (Theorem 3 in [Radhakrishnan, Belkin, Ulher, 2022]) Given x, assume w.l.o.g. that  $x^T x^1 = \max_{i=1:n} x^T x^i$ .

$$\lim_{L \to \infty} \frac{\kappa^L(x, x^i)}{\kappa^L(x, x^1)} = 0.$$
 (29)

and

$$\lim_{x \to \infty} f(x) = \operatorname{sgn} y^1 \qquad 1 - nn \tag{30}$$

Case A = 0,  $A' \neq 0$ : Networks implement singular kernel classifier

Theorem (Theorem 1 in [Radhakrishnan, Belkin, Ulher, 2022])

$$\lim_{L \to \infty} \frac{\kappa^{L}(x, x')}{(A')^{2L}(L+1)} = \frac{R(\|x - x'\|)}{\|x - x'\|^{\alpha}},$$
(31)

with  $\alpha = -4 \frac{\log A'}{\log B'}$  and  $R() \ge 0$ , bounded, and  $R(u) > \delta$  around 0.

L

- if  $\alpha > 0$ ,  $\frac{R(||x-x'||)}{||x-x'||^{\alpha}}$  is singular kernel
- Computationally not a problem: if data  $x^{1:n}$  distinct,  $G_0$  is well defined
- If  $x = x^i$ , set  $f(x) = y^i$ .

### Optimality of singular kernel classifier

Theorem (Theorem 2 in [Radhakrishnan, Belkin, Ulher, 2022]) If A = 0,  $A' \neq 0$  and  $\alpha = -d$  then  $\lim_{L\to\infty} f(x)$  is Bayes-optimal.

What activations \u03c6 satisfy this theorem?

$$\phi^{\text{opt}}(z) = \frac{1}{2^{d/4}} \frac{z^3 - 3z}{\sqrt{6}} + \sqrt{1 - 2^{1 - d/2}} \frac{z^2 - 1}{\sqrt{2}} + \frac{1}{2^{d/4}} z \quad \text{for } d \ge 2.$$
 (32)



#### Optimal singular kernels for d = 4, 8, 16, 32



phi

phi'

4 5

phi

0

----- phi

3

# Summary



b when A ≠ 0, lim<sub>L→∞</sub> κ<sup>L</sup>(x, x') = 0 for x ≠ x', and f(x) = 0 is vanishingly small (useless for regression), but sgnf(x) can be optimal for classification
 c Singular kernel α > d, α < d, majority vote kernel, and 1-nn kernel</li>

# Limits of some activation functions

 $\begin{array}{l} \phi^{\rm opt} \quad {\rm Bayes\ classifier} \\ \hline {\rm ReLU\ majority\ vote} \\ {\rm sigmoid\ } \frac{1}{1+e^{-z}}-\frac{1}{2}\ 1{\rm -nearest\ neighbor} \end{array}$