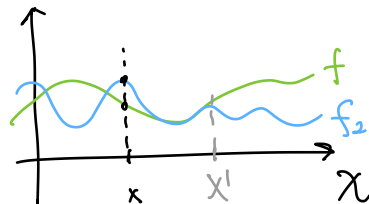# Lecture 17

Wide Neural Networks
     as GPs

No Q3
No HW7

**Gaussian Process** → *distribution over functions*

$\mathcal{X}$ sample space
(continuous)

$\{ f : \mathcal{X} \longrightarrow \mathbb{R} \} = \mathcal{F}$ ···· distribution on $\mathcal{F}$ = process

$\left[ f(x) \right]_{x \in \mathcal{X}}$

distribution of $f(x)$, $f \in \mathcal{F}$

$$\boxed{f(x) \sim N\left( \mu_x, \sigma_x^2 \right)}$$ for all $x \in \mathcal{X}$

$\operatorname{Cov}\left( f(x), f(x') \right) = K(x, x') \in \mathbb{R}$ $\quad G$    Gram matrix for SVM

Gaussian Process G.P.

$\mathcal{D} = \{ x^{1:n} \} \subset \mathcal{X}$

$$\boxed{\left[ K(x^i, x^j) \right]_{i,j=1:n}} = \Sigma \geqslant 0 \text{ for all } \mathcal{D}$$

Required
$\Updownarrow$
$K(,)$ satisfies Mercer condition

$f(x) \perp\!\!\!\perp f(x')$
"$y$"      "$y'$"

$\operatorname{Cov}(y, y') \lesssim 1 \Rightarrow y \approx y' \Rightarrow f$ **smooth**

GP is prior

Non-linear

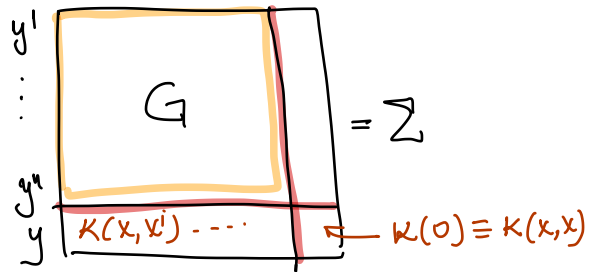Non-parametric → Bayesian regression

$\mathcal{D} = \{(x^i, y^i), i = 1:n\}$

↳ kernel

Prior $\quad GP(0, K)$

Posterior $\quad GP(\mu_{y|x, \mathcal{D}}, G_{x|\mathcal{D}})$

$\mu_{x|\mathcal{D}}$ $\quad$ "covariance"

Prediction given $x \in X$

wanted $\quad \Pr[f(x) | \mathcal{D}] = N(\mu_{x|\mathcal{D}}, \sigma^2_{x|\mathcal{D}})$

$\qquad y$

$\qquad\qquad$ = conditional distribution

$\mu_{x|\mathcal{D}} \equiv E[y] = [k(x, x^i) \cdots] G^{-1} \begin{bmatrix} y^1 \\ \vdots \\ y^n \end{bmatrix}$

$\qquad$ at $x$

$\sigma^2_{x|\mathcal{D}} = K(x, x) - [k(x, x^i) \cdots] G^{-1} \begin{bmatrix} k(x, x^i) \end{bmatrix}$



$\begin{matrix} y^1 \\ \vdots \\ y^n \\ y \end{matrix}$ $\quad \boxed{G} = \Sigma$

$\qquad K(x, x^i) \cdots \qquad\qquad K(0) \equiv K(x, x)$
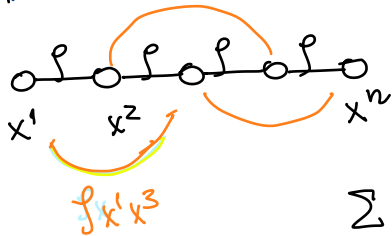
$G = [k(x^i, x^j)]$

$\begin{bmatrix} y \\ y^1 \\ \vdots \\ y^n \end{bmatrix} \sim N(0, \Sigma)$
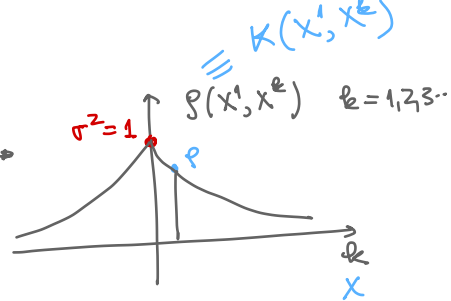
$\qquad$ joint distribution

**Ex1:** Markov chain



$$\rho = \rho(x^1, x^2) = \rho(x^2, x^3) = \cdots$$

$$\sigma^2 = 1 = Var(x^{1:n})$$

$$\Sigma = Cov\left(\begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix}\right) =$$

$$\equiv K(x^1, x^2)$$

$$\rho(x^1, x^k) \quad k = 1,2,3\cdots$$

$\sigma^2 = 1$

$n \to \infty$

**Ex 2.**

$\sigma^2$

$$\|x - x'\| = u$$

$$K(x, x') = K(\|x - x'\|)$$

(time)
(spatially)
invariant kernel

$K(u)$

**Rem:** $x, x'$

$$\|x - x'\| > R \implies$$

$$f(x) \perp\!\!\!\perp f(x')$$

$R$

$$\lim_{} k(u) = 0, \quad u \to \infty$$

$\to$ smoothing param

$$R: \quad supp\ K = [-R, R]$$

# 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

What fctions
do wide NN's
represent?

$\left(\begin{array}{c} \text{wide} \Leftrightarrow m_\ell \to \infty \\ \text{for } \ell=1:L-1 \end{array}\right)$

$$\mathcal{D} = \{(x^i, y^i)\} \text{ fixed}$$

The (Neural )Tangent Kernel (NTK)

any predictor training

What fctions do wide NN's represent?

(wide ⟺ $m_\ell \to \infty$ for $\ell = 1:L-1$)

Wide networks and Gaussian Processes

initialization (before training)

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

← because $\theta$ doesn't change much

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 $z^L \equiv f(x)$

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

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

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

- Parameter vector $\theta = \text{vector}\{w^{1:L}, \beta^{1:L}\} \in \mathbb{R}^p$ initialized i.i.d. $\sim N(0, 1)$
- $\sigma_{w,b}$ 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 $m_{1:L}$ 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

*any $f_\theta$, any $\alpha$*     $\theta \in \mathbb{R}^p$    $(x', y')$ fixed

- Assume training by gradient descent on $\hat{\mathcal{L}} = \sum_i \mathcal{L}(y^i, f(x^i))$   *empirical loss*
- The gradient of $\hat{\mathcal{L}}$

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

*no $y^i$*

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

$\theta^{t+1} - \theta^t \longrightarrow$

$$\dot{\theta} = -\eta \nabla_\theta f_\mathcal{D} \nabla_f \mathcal{L}_\mathcal{D} \quad \in \mathbb{R}^p \tag{4}$$

$f_{\theta^{t+1}} - f_{\theta^t} \longrightarrow$

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

$$\dot{f}_\mathcal{D} = -\eta \underbrace{(\nabla_\theta f_\mathcal{D})^T \nabla_\theta f_\mathcal{D}}_{G} \nabla_f \mathcal{L}_\mathcal{D} \quad \in \mathbb{R}^p \quad \leftarrow \text{at X data} \tag{6}$$

*vector*

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

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

4

$$X = \begin{bmatrix} x^1 \\ \vdots \end{bmatrix}_{\mathbb{R}^n} \xrightarrow{\theta} f_\theta = \begin{bmatrix} f_\theta(x^i) \end{bmatrix} \longrightarrow \nabla_\theta f = \begin{bmatrix} \dfrac{\partial f_\theta(x^1)}{\partial \theta_j} \end{bmatrix}^T \in \mathbb{R}^{p \times n} \qquad G = K(X,X)$$

$$y = \begin{bmatrix} y^1 \\ \vdots \end{bmatrix}_{\mathbb{R}^n} \longrightarrow \mathcal{L} = \begin{bmatrix} \mathcal{L}(f(x^i), y^i) \\ \vdots \end{bmatrix}_{\mathbb{R}^n}$$

# Gradient flow and NTK – summary

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

- $f_X$, $\nabla_\theta f_X$, $G$ 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?

  $\theta \sim N(0, \cdots)$ iid  NN initialization  ①

- First, we will analyze $\kappa$ for very wide neural networks with random parameters (e.g. at initialization)
- Then, we will analyze what happens during training under assumption (A3)  ②

# Wide NN's Gaussian Process (GP) ⓙ

- 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) = $ sum of $m_{L-1}$ i.i.d. random variables, hence $f(x) \sim Normal$ by CLT, for $m_{L-1}$ large
- Randomness is over weights $W^L, b^L$!!!
- We have $E[f(x)] = 0$ and

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

(8)

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

▶ From (8), for layer $l = 1 : L$ we have

$$\kappa^l(x, x') \;=\; E[z_j^l(x) z_j^l(x')] \;=\; \frac{\sigma_w^2}{m_{l-1}} (x^{l-1})^T (x')^{l-1} + \sigma_b^2 \qquad (9)$$

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

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

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

$$\kappa_{x,x'}^{l-1} \;=\; \left[ \begin{array}{cc} \kappa^{l-1}(x, x) & \kappa^{l-1}(x, x') \\ \kappa^{l-1}(x', x) & \kappa^{l-1}(x', x') \end{array} \right]. \qquad (11)$$

▶ 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  $\vartheta$ fixed

[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

$$\kappa^1(x, x') = \Sigma^1(x, x'), \quad \Sigma^1(x, x') = \frac{1}{m} x^T x' + 1 \qquad \leftarrow \text{ see next page} \quad (12)$$
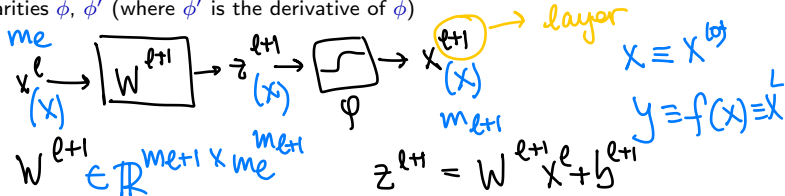
$$\kappa^{l+1}(x, x') = \kappa^l(x, x') \dot{\Sigma}^{l+1}(x, x') + \Sigma^{l+1}(x, x'), \qquad (13)$$

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

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

$$\text{and} \qquad L^{\phi}_{\Sigma} = E[\phi(X)\phi(X')] \text{ with}(X, X') \sim N(0, \begin{bmatrix} \Sigma(X, X) & \Sigma(X, X') \\ \Sigma(X, X') & \Sigma(X', X') \end{bmatrix}) \quad (16)$$

- 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 $\phi$)

# NN random initialization

$$W^{\ell+1}_{j,j'} \sim N\left(0, \frac{1}{m_\ell} \sigma_w^2\right) \quad iid$$

$$b^{\ell+1}_{j} \sim N\left(0, \sigma_b^2\right)$$

## intuition

$\ell = 0$

first layer

$$z^1 = W^0 x + b \sim N(0, \Sigma_1) \quad \in \mathbb{R}^{m_1}$$

$\uparrow$ data point

$x, x'$ data points

1) $j, j' \in 1:m_1$

$$\text{Cov}\left(z^1_j(x), z^1_{j'}(x')\right) = E\left[\left(W^1_{j:} x + b_j\right)\left(W^1_{j':} x' + b_{j'}\right)\right] = 0 \quad \text{for any } x, x'$$

2) $j = j' \in 1:m_1$

$$\text{Cov}\left(z^1_j(x), z^1_j(x')\right) =$$

$b_j, b_{j'}, W_{jk}, W_{j'k'}$  mutually independent  for all $k, k', j, j'$

$$= E\left[\left(W^1_{j:} x + b_j\right)\left(W^1_{j:} x' + b_j\right)\right] = \frac{1}{m_1}\sigma_w^2 I_{m_0}$$

$$= E\left[\underbrace{b_j^2}_{\sigma_b^2} + b_j \underbrace{W^1_{j:}(x' + x)}_{0} + x^T \underbrace{\left(W^{1T}_{j:} W^1_{j:}\right)}_{} x'\right] = \sigma_b^2 + \frac{1}{m_1} x^T x'$$

for any $j = 1:m_1$
depends on $x, x'$ data points

## 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 $\sigma_{b,w}^2$)

What next?
▶ Analysis of training by linearization
▶ Then, the NTK limit for $L \to \infty$ and its relevance for classification and regression

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

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

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

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

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

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

$$\kappa(x,x') = \nabla_\theta f_0(x)^T \nabla_\theta f_0(x') \quad \text{NTK}(\theta_o) \quad \textbf{constant during training} \tag{20}$$

$$G_0 \equiv \kappa_{X,X} \quad \text{Gram matrix at } \theta_o \tag{21}$$

from (4) ⟶
$$\dot{\theta}_t = -\eta \nabla_\theta f_0(X)^T \nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(x)) \tag{22}$$

from (5) ⟶
$$\dot{f}_t^{\text{lin}}(x) = -\eta \underbrace{\kappa(x,X)}_{\text{depends on } \theta_0}{}^T \nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(x)) \tag{23}$$

from (6) ⟶
$$\dot{f}_t^{\text{lin}}(X) = -\eta \, G_o \, \nabla_f \mathcal{L}(Y, f_t^{\text{lin}}(X))$$

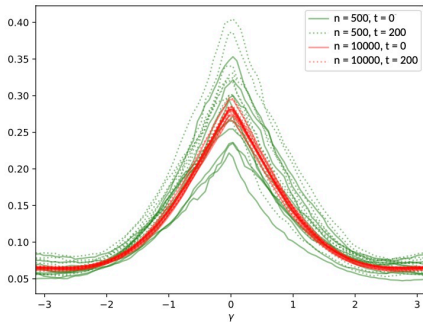# NTK during training – empirical evidence



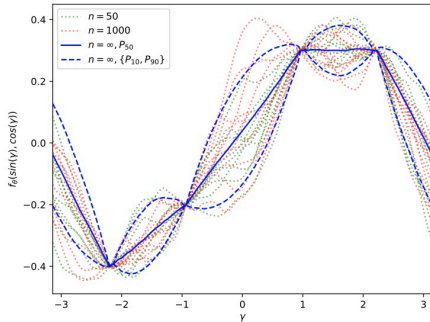Figure 1: Convergence of the NTK to a fixed limit for two widths $n$ and two times $t$.

Figure 2: Networks function $f_\theta$ near convergence for two widths $n$ and 10th, 50th and 90th percentiles of the asymptotic Gaussian distribution.

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

▶ For example, for $\mathcal{L}_{\mathrm{LS}}(y, f) = \frac{1}{2}(f - y)^2$, $\nabla_f \mathcal{L}_{\mathrm{LS}}(f, y) = f - y$. In this case, equations (22),(23) are a linear system and have an analytic solution.

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

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

$$f_t^{\mathrm{lin}}(x) = \underbrace{\kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) Y}_{\mu(x)} + \underbrace{f_0(x) - \kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) f_0(X)}_{\gamma(x)} \tag{26}$$

Notes:
  ▶ if $G_0 \succ 0$ then $e^{-\eta G_0 t} \to 0$ for $t \to \infty$
  ▶ in discrete time $t = 0, 1, 3, \ldots$ 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 $e^{-at} \approx (1 - a)^t$.
  ▶ $f_t^{\mathrm{lin}}(x) = f_0(x) + \kappa(x, X)^T G_0^{-1} \left( I - e^{-\eta G_0 t} \right) (Y - f_0(X))$

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