

Infinite-width Neural Networks: Relationship with Gaussian Process and Neural Tangent Kernel

Richard Xu

January 8, 2021

1 Preamble

In this tutorial, my contribution mainly has been the attempt to summarize the following papers and blogs in a unified and (hopefully) more intuitive for Computer Science researchers. In particular, the blogs below are extremely useful, and I encourage you to read the original blog as well.

Jachoon Lee, Lechao Xiao, Samuel S Schoenholz, Yasaman Bahri, Jascha Sohl-Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models under gradient descent. arXiv preprint arXiv:1902.06720, 2019

Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In Advances in neural information processing systems, pages 8571–8580, 2018

J. H. Lee, Y. Bahri, R. Novak, S. S. Schoenholz, J. Pennington, and J. Sohl-Dickstein. Deep neural networks as gaussian processes. ICLR, 2018

Radford M. Neal. Priors for infinite networks (tech. rep. no. crg-tr-94-1). University of Toronto, 1994
Alexander G de G Matthews, Jiri Hron, Richard E Turner, and Zoubin Ghahramani. Samplethen-optimize posterior sampling for Bayesian linear models. In NeurIPS Workshop on Advances in Approximate Bayesian Inference, 2017.

<https://www.uv.es/gonmagar/blog/2019/01/21/DeepNetworksAsGPs>

<https://bryn.ai/jekyll/update/2019/04/02/neural-tangent-kernel.html>

http://chenyilan.net/files/ntk_derivation.pdf

<http://chenyilan.net/files/linearization.pdf>

1.1 notations

I attempted to unify notations, where I used the following definition for Neural Network functions:

$$z_k^l(x) = b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \times \phi\left(z_j^{l-1}(x)\right) \quad W_{k,j}^l \sim \mathcal{N}\left(0, \frac{1}{\sqrt{N_l}}\right) \quad b_k^l \sim \mathcal{N}(0, \sigma_b) \quad \text{or :} \quad (1)$$

$$z_k^l(x) = \sigma_b b_k^l + \sum_{j=1}^{N_l} \frac{1}{\sqrt{N_l}} W_{k,j}^l \times \phi\left(z_j^{l-1}(x)\right) \quad W_{k,j}^l \sim \mathcal{N}(0, 1) \quad b_k^l \sim \mathcal{N}(0, 1)$$

1. $k \in \{1, \dots, N_{l+1}\}$ indexes elements of z^l
2. $i \in \{1, \dots, N_{l+1}\}$ also indexes elements of z^l , and it is used when k is reserved to a specific index
3. $j \in \{1, \dots, N_k\}$ indexes elements of z^{l-1}
4. $W^l \in \mathcal{R}^{N_{l+1} \times N_l}$
5. x and x' are used to indicate two data points
6. k and k' indexes two functional output of z^l
7. size of data input is $|d_{\text{in}}|$

1.2 Others minor contributions

I made the derivations a bit more verbose for people to follow

To make this tutorial self-contained, I have included a very quick introduction on the relevant topics include Gaussian Process and Central Limit Theorem

2 Gaussian Process

This tutorial makes frequent references to GP, so we talk about it briefly:

- \mathcal{GP} is a (potentially infinite) collection of RVs, s.t., joint distribution of every finite subset of RVs is multivariate Gaussian:

$$f \sim \mathcal{GP}(\mu(x), \mathcal{K}(x, x')) \quad \text{for any arbitrary } x, x'$$

- **prior** defined over $p(f|\mathcal{X})$, instead of $p(x)$ over $\mathcal{X} \equiv \{x_1, \dots, x_k\}$

$$p(f|\mathcal{X}) \equiv p\left(\begin{bmatrix} f(x_1) \\ \vdots \\ f(x_k) \end{bmatrix}\right) = \mathcal{N}\left(0, K\right) = \mathcal{N}\left(0, \begin{bmatrix} k(x_1, x_1) & \dots & k(x_1, x_k) \\ \vdots & \ddots & \vdots \\ k(x_k, x_1) & \dots & k(x_k, x_k) \end{bmatrix}\right)$$

2.1 marginal and conditional marginal under noisy output

- in a regression setting:

$$y_i = f(x_i) + \epsilon_i \quad \epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma_\epsilon^2)$$

- joint $[\mathcal{Y}, y^*]^\top$, after integrate out f :

$$\begin{aligned} p\left(\begin{bmatrix} \mathcal{Y} \\ y^* \end{bmatrix} \middle| \begin{bmatrix} \mathcal{X} \\ x^{*\top} \end{bmatrix}, \sigma_\epsilon^2\right) &= \int p\left(\begin{bmatrix} \mathcal{Y} \\ y^* \end{bmatrix} \middle| \begin{bmatrix} \mathcal{X} \\ x^{*\top} \end{bmatrix}, f\right) p(f|\mathcal{X}, x^*) df \\ &= \int \mathcal{N}\left(\begin{bmatrix} \mathcal{Y} \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f(\mathcal{X}) \\ f(x^{*\top}) \end{bmatrix}, \sigma_\epsilon^2 I\right) p(f|\mathcal{X}, x^*) df \\ &= \mathcal{N}\left(0, \begin{bmatrix} \underbrace{K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 I}_{\Sigma_{1,1}} & \underbrace{K(\mathcal{X}, x^*)}_{\Sigma_{1,2}} \\ \underbrace{K(x^*, \mathcal{X})}_{\Sigma_{2,1}} & \underbrace{K(x^*, x^*) + \sigma_\epsilon^2}_{\Sigma_{2,2}} \end{bmatrix}\right) \end{aligned}$$

- **predictive distribution** of $y^*|\mathcal{Y}$ using conditional formula of multivariate Gaussian:

$$\begin{aligned} p(y^*|\mathcal{Y}, \mathcal{X}, x^*) &= \mathcal{N}\left(\underbrace{\mathbf{0}}_{\mu_2} + \underbrace{K(x^*, \mathcal{X})}_{\Sigma_{2,1}} \underbrace{(K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 I)^{-1}}_{\Sigma_{1,1}^{-1}} (\mathcal{Y} - \underbrace{\mathbf{0}}_{\mu_1}), \right. \\ &\quad \left. \underbrace{k(x^*, x^*) + \sigma_\epsilon^2}_{\Sigma_{2,2}} - \underbrace{K(x^*, \mathcal{X})}_{\Sigma_{2,1}} \underbrace{(K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 I)^{-1}}_{\Sigma_{1,1}^{-1}} \underbrace{K(\mathcal{X}, x^*)}_{\Sigma_{1,2}}\right) \end{aligned}$$

2.2 marginal and conditional marginal under noiseless output

- **posterior** of f given \mathcal{Y} in regression:

$$p\left(\begin{bmatrix} \mathcal{Y} \\ f \end{bmatrix} \middle| \begin{bmatrix} \mathcal{X} \\ \mathbf{x}^\top \end{bmatrix}\right) = p\left(\begin{bmatrix} f(\mathcal{X}) \\ f(\mathbf{x}) \end{bmatrix}\right) = \mathcal{N}\left(0, \begin{bmatrix} K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 \mathbf{I} & K(\mathcal{X}, \mathbf{x}) \\ K(\mathbf{x}, \mathcal{X}) & K(\mathbf{x}, \mathbf{x}) \end{bmatrix}\right)$$

for arbitrary variable \mathbf{x}

conditional marginal of $y^*|\mathcal{Y}$ using conditional formula of multivariate Gaussian:

$$p(f|\mathcal{X}, \mathcal{Y}) = \mathcal{GP}\left(K(\mathbf{x}, \mathcal{X})(K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathcal{Y},\right. \\ \left. k(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathcal{X})(K(\mathcal{X}, \mathcal{X}) + \sigma_\epsilon^2 \mathbf{I})^{-1} K(\mathcal{X}, \mathbf{x}')\right)$$

- **deterministic function** $y_i = f(x_i)$ is used, e.g., neural network's read-out layer $f(x_i)$, data y_i

$p([\mathcal{Y}, y^*]^\top)$ no longer need to integrate f :

$$p\left(\begin{bmatrix} \mathcal{Y} \\ y^* \end{bmatrix} \middle| \begin{bmatrix} \mathcal{X} \\ x^{*\top} \end{bmatrix}\right) = p\left(\begin{bmatrix} f(\mathcal{X}) \\ f(x^*) \end{bmatrix}\right) = \mathcal{N}\left(0, \begin{bmatrix} K(\mathcal{X}, \mathcal{X}) & K(\mathcal{X}, x^*) \\ K(x^*, \mathcal{X}) & K(x^*, x^*) \end{bmatrix}\right)$$

predictive distribution $y^*|\mathcal{Y}$ using conditional formula of multivariate Gaussian:

$$p(y^*|\mathcal{Y}, \mathcal{X}, x^*) = \mathcal{N}\left(K(x^*, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1} \mathcal{Y},\right. \\ \left. k(x^*, x^*) - K(x^*, \mathcal{X})K(\mathcal{X}, \mathcal{X})^{-1} K(\mathcal{X}, x^*)\right)$$

3 sample and optimize

in words can be summarized as follows:

given a linear equation with non-empty null-space, and some prior distribution $w \sim \mathcal{N}()$ defined over $w \in \mathcal{W}$. Since null-space is non-empty, there exist a “subset” of \mathcal{W} , namely w^* satisfy linear equation, and forming a new posterior distribution:

$$w^* \sim p(w \mid \|Aw = y\|_2^2) \quad (2)$$

We intend to draw a sample w^* from this posterior.

sample-then-optimize allows us to sample $w^0 \sim \mathcal{N}()$, and optimize using gradient descend to obtain w^* . The following proves they indeed have the same probability as sampling from posterior directly

3.0.1 proof

$$\begin{aligned} w &\sim \mathcal{N}(0, \lambda I) \\ y &= Aw \quad A \in \mathcal{R}^{N \times M} \end{aligned} \quad (3)$$

when $N < M$, there is a family of solutions and this can be understood as to see the following:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} A_{1,1} & A_{1,2} & \dots & A_{1,M-1} & A_{1,M} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N,1} & A_{N,2} & \dots & A_{N,M-1} & A_{N,M} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \\ w_{N+1} \\ \vdots \\ w_M \end{bmatrix} \quad (4)$$

general solution is:

$$\begin{aligned} \hat{w} &= A^+ y + (I - A^+ A) t \quad \text{for arbitrary } t \in \mathcal{R}^M \\ &= \hat{w}_{||} + \hat{w}_{\text{null}} \end{aligned} \quad (5)$$

simple verification would see that,

$$\begin{aligned} A(\hat{w}_{||}) &= y \\ \implies A(\hat{w}_{||} + \hat{w}_{\text{null}}) &= y \\ \text{because } A(\hat{w}_{||}) + \underbrace{A(\hat{w}_{\text{null}})}_0 &= y \end{aligned} \quad (6)$$

informally $\hat{w}_{||} = A^+ y$ gives the solution \hat{w} for:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} A_{1,1} & A_{1,2} & \dots & A_{1,M-1} & A_{1,M} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ A_{N,1} & A_{N,2} & \dots & A_{N,M-1} & A_{N,M} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \\ w_{N+1} \\ \vdots \\ w_M \end{bmatrix} \quad (7)$$

and $\hat{w}_{\text{null}} = (I - A^+ A) t$ to ensure \hat{w} also simultaneously satisfy:

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 & \dots & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \\ w_{N+1} \\ \vdots \\ w_M \end{bmatrix} \quad (8)$$

pseudo-inverse:

$$\begin{aligned} AA^+A &= A \\ \implies A^+ &= A^\top (AA^\top)^{-1} \quad A^+ \in \mathcal{R}^{M \times N} \quad (AA^\top) \in \mathcal{R}^{N \times N} \quad \text{full rank} \\ &= (A^\top A)^+ A^\top \quad A^\top A \in \mathcal{R}^{M \times M} \quad \text{not full rank} \\ \implies A^+ &= \lim_{\delta^2 \rightarrow 0} A^\top (AA^\top + \delta^2 I)^{-1} \quad \text{limit definition of pseudo inverse} \end{aligned} \quad (9)$$

it is easy to show $(I - A^+A)$ is projection into $\text{null}(A)$ because:

$$\begin{aligned} A(I - A^+A) &= A(I - A^\top (AA^\top)^{-1} A) \\ &= A - (AA^\top)(AA^\top)^{-1}A \\ &= \mathbf{0} \end{aligned} \quad (10)$$

since each of the columns of the matrix $(I - A^+A)$ span null space of A , then:

$\hat{w}_{\text{null}} = (I - A^+A) t$ is some linear combination of columns of $(I - A^+A)$ and hence:

$$\hat{w}_{\text{null}} = (I - A^+A) t \in \text{null}(A) \quad (11)$$

3.0.2 what happen to w during gradient descend?

suppose we have the objective function:

$$\begin{aligned} f(w) &= \|Aw - y\|_2^2 \\ \implies \frac{\partial f(w)}{\partial w} &= \frac{\partial (Aw - y)^\top (Aw - y)}{\partial w} \\ &= \frac{\partial (w^\top A^\top - y^\top)(Aw - y)}{\partial w} \\ &= \frac{\partial (w^\top A^\top Aw - y^\top Aw - w^\top A^\top y + y^\top y)}{\partial w} \\ &= 2A^\top Aw - 2A^\top y \\ \implies \Delta w &= -\eta(A^\top Aw - A^\top y) \end{aligned} \quad (12)$$

therefore, say at time $t = 0$ we let the initial point w^0 :

$$w^0 \sim \mathcal{N}(0, \lambda I) \quad (13)$$

and break up w^0 into two components with axes in different sub-spaces:

$$w^0 = w_{\text{null}}^0 + w_{||}^0 \quad (14)$$

and apply the following transform in gradient descend:

$$\begin{aligned}
\Delta w(0) &\propto A^\top A w^0 - A^\top y \\
&= A^\top A (w_{\text{null}}^0 + w_{\parallel}^0) - A^\top y \\
&= A^\top A w_{\text{null}}^0 + A^\top A w_{\parallel}^0 - A^\top y
\end{aligned} \tag{15}$$

because: $g^\top A^\top f = f^\top A g = 0$ $g \in \mathcal{R}^M$, $f \in \mathcal{R}^N$ $g \in \text{null}(A)$

$$\begin{aligned}
\Delta w(0) &\propto \underbrace{A^\top}_{f} A \underbrace{w_{\text{null}}^0}_{g} + A_f^\top A w_{\parallel}^0 - A^\top y \\
&= 0 + A^\top A w_{\parallel}^0 - A^\top y
\end{aligned} \tag{16}$$

the above means that there is no “movement” for w_{null} component, $\forall t$ in $\text{null}(A)$. only movement occur in the w_{\parallel} subspace, assume gradient descend converge eventually. So at the end of gradient descend at time $t = \tau$:

$$\begin{aligned}
\Delta w(\tau) &= 0 \\
\implies A^\top A w_{\parallel}^\tau - A^\top y &= 0 \\
\implies A^\top A w_{\parallel}^\tau &= A^\top y \\
\implies w_{\parallel}^\tau &= (A^\top A)^+ A^\top y \\
&= A^+ y
\end{aligned} \tag{17}$$

now we look at the general solution in Eq.(5) and replacing arbitrary t with starting point w^0 as a repeating applications of gradient descend:

$$\begin{aligned}
w^* &= \textcolor{red}{w}_{\parallel}^\tau + \textcolor{blue}{w}_{\text{null}}^\tau \\
&= \underbrace{\textcolor{red}{A}^+ y}_{\mu} + \underbrace{(\textcolor{blue}{I} - \textcolor{red}{A}^+ A)}_{\sigma} w^0
\end{aligned} \tag{18}$$

so having $w^0 \sim \mathcal{N}(0, \lambda I)$, w^* is just apply an Affine transform with the following Gaussian distribution:

$$w^* \sim \mathcal{N}(A^+ y, \lambda(I - A^+ A)) \tag{19}$$

so moral of the story, instead of sample from $p(w|y)$ given the above, one can just sample w^0 from prior $p(w)$, and then optimize using gradient:

$$w^* = \arg \min_w \|Aw - y\|_2^2 \tag{20}$$

3.0.3 apply to Gaussian Process

so if you want to sample w^* from the posterior:

$$\begin{aligned}
\mu_*(y^*|x^*, \mathcal{X}, \mathcal{Y}) &= \mathcal{K}(x^*, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \mathcal{Y} \\
\Sigma_*(y^*|x^*, \mathcal{X}, \mathcal{Y}) &= \mathcal{K}(x^*, x^*) - \mathcal{K}(x^*, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \mathcal{K}(\mathcal{X}, x^*)
\end{aligned} \tag{21}$$

one may sample $\theta_0 \sim \mathcal{N}(0, \mathcal{K})$, and then because of the assumption:

in MSE context of read-out only layer:

$$\begin{aligned}
C &= \frac{1}{2} \left\| \underbrace{\bar{a}(\mathcal{X})^\top}_A \underbrace{\theta^{L+1}}_w - \mathcal{Y} \right\|_2^2 \\
&= \frac{1}{2} \left\| \hat{y}_t(\mathcal{X}) - \mathcal{Y} \right\|_2^2
\end{aligned} \tag{22}$$

the only difference, is instead of representing a prior for θ_0^{L+1} , we have a prior for $\hat{y}(x, \theta_0)$, which we know how to express (NNGP) and sample

and, instead of express posterior of $w^* \equiv \theta^{L+1*} \equiv \theta_\infty$, we let:

$$w^* \equiv \hat{y}(x, \theta_\infty) \quad \text{since } \hat{y}(x, \theta_\infty) \text{ is linear in } \theta_{t=\infty} \tag{23}$$

We can derive the expression of posterior for $\hat{y}(x, \theta_\infty)$, but we do not sample from it, and instead we use **sample-then-optimize**

4 Kernel methods

consider the equation:

$$\begin{aligned}
 y &= \phi(x)^\top \mathbf{w} \\
 &= \begin{bmatrix} \phi_1(x) \\ \vdots \\ \phi_m(x) \end{bmatrix}^\top \mathbf{w} \\
 &= [\phi_1(x) \quad \dots \quad \phi_m(x)] \mathbf{w}
 \end{aligned} \tag{24}$$

using definition:

$$\begin{aligned}
 \mathcal{Y} &= [y_1, \dots, y_n]^\top \\
 \Phi &= [\phi(x_1), \dots, \phi(x_n)]^\top \\
 &= \underbrace{\begin{bmatrix} \phi_1(x_1) & \dots & \phi_m(x_1) \\ \vdots & & \vdots \\ \phi_1(x_n) & \dots & \phi_m(x_n) \end{bmatrix}}_{n \times m}
 \end{aligned} \tag{25}$$

Ridge regression can be re-written as:

$$\begin{aligned}
 \mathbf{w}^* &= \arg \min_{\mathbf{w}} \sum_{i=1}^n (y_i - \phi(x_i)^\top \mathbf{w})^2 + \lambda \|\mathbf{w}\|_2^2 \\
 &= \arg \min_{\mathbf{w}} \|\mathcal{Y} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2
 \end{aligned} \tag{26}$$

just like the normal ridge regression, the least-square solution is:

$$\mathbf{w}^* = (\underbrace{\Phi^\top \Phi}_{m \times m} + \lambda I)^{-1} \Phi^\top \mathcal{Y} \tag{27}$$

substitute \mathbf{w}^* back to $y = \phi(x)^\top \mathbf{w}$ for a single pair of data,output (x, y) :

$$\begin{aligned}
 y_{\mathbf{w}^*}(x) &= \phi(x)^\top \mathbf{w}^* \\
 &= \phi(x)^\top (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathcal{Y} \\
 &= \underbrace{\phi(x)^\top \Phi^\top}_{1 \times n} (\underbrace{\Phi \Phi^\top}_{n \times n} + \lambda I)^{-1} \mathcal{Y} \\
 &\text{using identity } (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top = \Phi^\top (\Phi \Phi^\top + \lambda I)^{-1}
 \end{aligned} \tag{28}$$

4.1 Kernel trick

the above looks all good, except we want to avoid computing $\phi(x)$ explicitly, especially when m is large! However, knowing

$$\begin{aligned}
 [\Phi \Phi^\top]_{i,j} &= \phi(x_i)^\top \phi(x_j) = \mathcal{K}(x_i, x_j) \\
 [\phi(x)^\top \Phi^\top]_j &= \phi(x)^\top \phi(x_j) = \mathcal{K}(x, x_j)
 \end{aligned} \tag{29}$$

we dodged the bullet of computing $\phi(x)$ explicitly!

4.2 relationship with Neural Tangent Kernel

Taylor Expansion of $f_{\mathbf{w}}(x)$ around w_0 :

$$f_{\mathbf{w}}(x) \equiv f(\mathbf{w}, x) \approx f(w_0, x) + \underbrace{\nabla_w f(w_0, x)}_{\phi(x)^\top} (w - w_0) + \dots \quad (30)$$

so, in theory, one may solve this using Kernel regression. However, question is **why still using neural networks?**
in here, we have not made any linkage to gradient descend yet.

4.3 relationship with gradient flow

This is a simplified version to Section[8]:

Gradient descend algorithm:

$$\begin{aligned} \theta_{t+1} &= \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_{t+1}) \\ \implies \frac{\theta_{t+1} - \theta_t}{\eta} &= -\nabla_{\theta} \mathcal{L}(\theta_{t+1}) \\ \implies \lim_{\eta \rightarrow 0} \frac{\theta_{t+1} - \theta_t}{\eta} &= \frac{d\theta(t)}{dt} = -\nabla_{\theta} \mathcal{L}(\theta) \end{aligned} \quad (31)$$

let's substitute that into **least square** problem:

$$\begin{aligned} \mathcal{L}(\theta) &= \frac{1}{2} \|\tilde{y}(\theta) - y\|_2^2 \\ \implies \nabla_{\theta} \mathcal{L}(\theta) &= \nabla_{\theta} \tilde{y}(\theta) (\tilde{y}(\theta) - y) \\ \implies \frac{d\theta(t)}{dt} &= -\nabla_{\theta} \tilde{y}(\theta) (\tilde{y}(\theta) - y) \end{aligned} \quad (32)$$

so let's look at $\frac{d\tilde{y}(\theta_t)}{dt}$:

$$\begin{aligned} \frac{d\tilde{y}(\theta_t)}{dt} &= \frac{\partial \tilde{y}(\theta(t))}{\partial \theta(t)}^\top \frac{d\theta(t)}{dt} \\ &= \nabla_{\theta} \tilde{y}(\theta) \left(-\nabla_{\theta} \tilde{y}(\theta) (\tilde{y}(\theta) - y) \right) \\ &= -\underbrace{\nabla_{\theta} \tilde{y}(\theta)^\top \nabla_{\theta} \tilde{y}(\theta)}_{K(\theta)} (\tilde{y}(\theta) - y) \\ &\approx -K(\theta_0) (\tilde{y}(\theta) - y) \end{aligned} \quad (33)$$

5 GP for Neural Network: Direct computation

5.1 neural network function

using parameters:

$$\theta \equiv \{W^L, b^L, \dots, W^1, b^1\} \quad (34)$$

Deep neural network function $f_\theta(X)$ is defined as:

$$\begin{aligned} f_\theta(X) &= W^L \phi^L(X) + b^L \\ &= W^L (\phi^{L-1}(X) W^{L-1} + b^{L-1}) + b^L \\ &\dots \\ &= W^L \dots (W^1 \phi^1(X) + b^1) + \dots + b^L \end{aligned} \quad (35)$$

it should be noted that non-linear output $\phi^l(\cdot)$:

$$\begin{aligned} \phi^L(X) &\equiv \phi^L(X | \theta^1, \dots, \theta^{L-1}) \\ &\equiv \phi^L(X | W^1, b^1, \dots, W^{L-1}, b^{L-1}) \end{aligned} \quad (36)$$

5.2 Apply NN function in predictive distribution

However, applying NN function in predictive distribution: prior is defined over θ instead of over f . i.e., i.i.d noises are injected to each element of θ . The predictive distribution:

$$p\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix}\right) = \int \mathcal{N}\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f_\theta(X) \\ f_\theta(x^*) \end{bmatrix}, \sigma_\epsilon^2 I\right) \mathcal{N}(\theta | 0, \sigma_\theta^2 I) d\theta \quad (37)$$

The integral is **not** analytic!!

5.3 what is the predictive distribution

eventually, we will need to ask an even harder question on, i.e., suppose we let $N^l \equiv |W^l|$, i.e., the “width” of the neural network at each layer l , and we would like to study the effect of:

$$p\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix}\right) \xrightarrow[N^1, \dots, N^L \rightarrow \infty]{d} ? \quad (38)$$

however, firstly, we ask the question on, what is:

$$p\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix}\right) = ? \quad (39)$$

attempt to compute it **directly**, by looking the **mean** and **variance**:

$$\begin{aligned} &\mathbb{E} \left[\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right] \\ &\mathbb{E} \left[\begin{bmatrix} y \\ y^* \end{bmatrix} \begin{bmatrix} y^\top & y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right] \end{aligned}$$

5.3.1 look at the mean:

$$\begin{aligned}
& \mathbb{E} \left[\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right] \\
&= \int_y \int_{y^*} \begin{bmatrix} y \\ y^* \end{bmatrix} p \left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right) dy dy^* \\
&= \int_y \int_{y^*} \begin{bmatrix} y \\ y^* \end{bmatrix} \int_{\theta} p \left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \theta, \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right) p(\theta | \sigma_{\theta}^2) d\theta dy dy^* \\
&= \int_{\theta} \int_y \int_{y^*} \underbrace{\begin{bmatrix} y \\ y^* \end{bmatrix} \mathcal{N} \left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f_{\theta}(X) \\ f_{\theta}(x^*) \end{bmatrix}, \sigma_{\epsilon}^2 I \right)}_{\mathbb{E} \left[\begin{bmatrix} y \\ y^* \end{bmatrix} \right] = \begin{bmatrix} f_{\theta}(X) \\ f_{\theta}(x^*) \end{bmatrix}} dy dy^* \mathcal{N}(\theta | 0, \sigma_{\theta}^2 I) d\theta \\
&= \int \begin{bmatrix} f_{\theta}(X) \\ f_{\theta}(x^*) \end{bmatrix} \mathcal{N}(\theta | 0, \sigma_{\theta}^2 I) d\theta \quad \text{to expand one layer :} \\
&= \int \begin{bmatrix} \phi^L(X) W^L + b^L \\ \phi^L(x^{*\top}) W^L + b^L \end{bmatrix} \mathcal{N}(W^L | 0, \sigma_w^2 I) \mathcal{N}(b^L | 0, \sigma_b^2 I) \mathcal{N}(\theta^1, \dots, \theta^{L-1} | 0, \sigma_{\theta}^2 I) d\theta^1, \dots, \theta^{L-1} dW^L db^L \\
&= \int \left[\underbrace{\phi^L(X) \int W^L \mathcal{N}(W^L | 0, \sigma_w^2 I) dW^L}_{=0} + \underbrace{\int b^L \mathcal{N}(b^L | 0, \sigma_b^2 I) db^L}_{=0} \right] \mathcal{N}(\theta^1, \dots, \theta^{L-1} | 0, \sigma_{\theta}^2 I) d\theta^1, \dots, \theta^{L-1} \\
&\quad \underbrace{\phi^L(x^{*\top}) \int W^L \mathcal{N}(W^L | 0, \sigma_w^2 I) dW^L}_{=0} + \underbrace{\int b^L \mathcal{N}(b^L | 0, \sigma_b^2 I) db^L}_{=0} \\
&= \begin{bmatrix} 0 \\ 0 \end{bmatrix}
\end{aligned} \tag{40}$$

note we are not dealing with infinity at the moment

5.3.2 look at co-variance

$$\mathbb{E} \left[\begin{bmatrix} y \\ y^* \end{bmatrix} \begin{bmatrix} y^{\top} & y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right] \tag{41}$$

Apply same trick as calculating mean, i.e., introducing θ and then integrate it out:

$$\begin{aligned}
&= \int_y \int_{y^*} \int_{\theta} p \left(\begin{bmatrix} y \\ y^* \end{bmatrix} \begin{bmatrix} y^{\top} & y^* \end{bmatrix} \middle| \theta, \begin{bmatrix} X \\ x^{*\top} \end{bmatrix} \right) p(\theta | \sigma_{\theta}^2) d\theta dy dy^* \\
&= \int_{\theta} \int_y \int_{y^*} \underbrace{\begin{bmatrix} y \\ y^* \end{bmatrix} \begin{bmatrix} y^{\top} & y^* \end{bmatrix} \mathcal{N} \left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f_{\theta}(X) \\ f_{\theta}(x^*) \end{bmatrix}, \sigma_{\epsilon}^2 I \right)}_{\mathbb{E}[Z^2] \quad Z \text{ is not mean-subtracted}} dy dy^* \mathcal{N}(\theta | 0, \sigma_{\theta}^2 I) d\theta
\end{aligned} \tag{42}$$

Let $Z = \begin{bmatrix} y \\ y^* \end{bmatrix}$:

$$\begin{aligned}
&\text{Var}[Z] = \mathbb{E}[Z^2] - (\mathbb{E}[Z])^2 \implies \mathbb{E}[Z^2] = \text{Var}[Z] + (\mathbb{E}[Z])^2 \\
&= \int_{\theta} \underbrace{\sigma_{\epsilon}^2 I}_{\text{Var}[Z]} + \underbrace{\begin{bmatrix} f_{\theta}(X) \\ f_{\theta}(x^*) \end{bmatrix} \begin{bmatrix} f_{\theta}(X)^{\top} & f_{\theta}(x^*) \end{bmatrix}}_{(\mathbb{E}[Z])^2} \mathcal{N}(\theta | 0, \sigma_{\theta}^2 I) d\theta \\
&= \sigma_{\epsilon}^2 I + \int_{\theta} \left[\begin{pmatrix} \phi^L(X) W^L + b^L \\ \phi^L(x^{*\top}) W^L + b^L \end{pmatrix} \begin{pmatrix} W^{L\top} \phi^L(X)^{\top} + b^{L\top} \\ W^{L\top} \phi^L(x^{*\top})^{\top} + b^{L\top} \end{pmatrix} + \begin{pmatrix} \phi^L(X) W^L + b^L \\ \phi^L(x^{*\top}) W^L + b^L \end{pmatrix} \begin{pmatrix} W^{L\top} \phi^L(x^{*\top})^{\top} + b^{L\top} \\ W^{L\top} \phi^L(X)^{\top} + b^{L\top} \end{pmatrix} \right] \mathcal{N}(\theta | 0, \sigma_{\theta}^2 I) d\theta
\end{aligned} \tag{43}$$

realize $\mathbf{Cov}(x^L(X)W^L, b^L) = 0$:

$$= \sigma_\epsilon^2 I + \int_{\theta} \begin{bmatrix} \phi^L(X)W^L W^{L\top} x^L(X)^{\top} + b^L b^{L\top} & \phi^L(X)W^L W^{L\top} x^L(x^{\star\top})^{\top} + b^L b^{L\top} \\ \phi^L(x^{\star\top})W^L W^{L\top} \phi^L(X)^{\top} + b^L b^{L\top} & \phi^L(x^{\star\top})W^L W^{L\top} \phi^L(x^{\star\top})^{\top} + b^L b^{L\top} \end{bmatrix} \mathcal{N}(\theta \mid 0, \sigma_\theta^2 I) d\theta \quad (44)$$

factorize $\mathcal{N}(\theta)$ as each element of θ is independent:

$$\mathcal{N}(\theta \mid 0, \sigma_\theta^2 I) d\theta = \mathcal{N}(\theta^L \mid 0, \sigma_\theta^2 I) \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} \quad (45)$$

$$= \int \begin{bmatrix} \sigma_w^2 \phi^L(X) x^L(X)^{\top} + \sigma_b^2 & \sigma_w^2 \phi^L(X) \phi^L(x^{\star\top})^{\top} + \sigma_b^2 \\ \sigma_w^2 \phi^L(x^{\star\top}) \phi^L(X)^{\top} + \sigma_b^2 & \sigma_w^2 \phi^L(x^{\star\top}) \phi^L(x^{\star\top})^{\top} + \sigma_b^2 \end{bmatrix} \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} \quad (46)$$

let's taking the **left corner** element, and expand θ by one:

$$\begin{aligned} & \int \sigma_w^2 \phi^L(X) \phi^L(X)^{\top} \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} + \int \sigma_b^2 \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} \\ &= \sigma_w^2 \int \phi^L(X) \phi^L(X)^{\top} \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} + \sigma_b^2 \end{aligned} \quad (47)$$

as we know $\phi^L(X) \phi^L(X)^{\top} \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} + \sigma_b^2$:

$$= \sigma_b^2 + \sigma_w^2 \int \left[\phi(W^{L-1} \phi^{L-1}(X) + b^{L-1}) \phi(W^{L-1} \phi^{L-1}(X) + b^{L-1})^{\top} \right] \mathcal{N}(\theta^{1,\dots,L-1} \mid 0, \sigma_\theta^2 I) d\theta^{1,\dots,L-1} \quad (48)$$

it's difficult to see what is this distribution is.

6 Single layer neural network

$$\begin{aligned}
 f_k(x) &= b_k + \sum_{j=1}^H v_{jk} h_j(x) \\
 h_j(x) &= \tanh \left(a_j + \sum_{i=1}^I u_{ij} x_i \right)
 \end{aligned} \tag{49}$$

this is very strange way to define neural network, and it defines it to part of the second layer:

$$\begin{aligned}
 f_k(x) &= \underbrace{b_k}_{z_k^l} + \sum_{j=1}^{\overbrace{H}^{N_l}} \underbrace{v_{jk}}_{W_{k,j}^l} \times \underbrace{\tanh}_{\phi} \left(\underbrace{a_j}_{b_j^{l-1}} + \underbrace{u_{:,j}^\top}_{W_{:,j}^{l-1}{}^\top} x \right) \\
 &\quad \underbrace{\hspace{10em}}_{z_j^{l-1}(x)} \\
 \implies z_k^l(x) &= b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \times \phi(z_j^{l-1}(x)) \quad \text{modern notation}
 \end{aligned} \tag{50}$$

6.1 $p(z_k^l(x))$ for single input x

We need CLT for computing this probability.

6.1.1 Central Limit Theorem:

$$X^{(1)}, X^{(2)}, \dots, X^{(n)} \text{ are i.i.d samples} \tag{51}$$

note any **arbitrary** distribution with *bounded variance* for $X^{(i)}$ will do
let \bar{X} be sample mean, and let: $\sigma^2 = \text{Var}[X^{(1)}]$
Limiting form of the distribution:

$$\begin{aligned}
 \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, \sigma^2) \\
 (\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, \frac{\sigma^2}{n}) \\
 \frac{1}{\sigma} \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, 1)
 \end{aligned} \tag{52}$$

Similarly, instead of “**sample mean**”, it can be also be applied to “**sample sum**” of i.i.d random variables:

$$\begin{aligned}
 \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, \sigma^2) \\
 \implies \sqrt{n} \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, \sqrt{n}^2 \sigma^2) = \mathcal{N}(0, n\sigma^2) \\
 \implies n(\bar{X} - \mathbb{E}[X^{(1)}]) &\xrightarrow{d} \mathcal{N}(0, n\sigma^2) \\
 \implies \left(\sum_{i=1}^n X_i - n\mathbb{E}[X^{(1)}] \right) &\xrightarrow{d} \mathcal{N}(0, n\sigma^2)
 \end{aligned} \tag{53}$$

choose one of these conditions to suit the situation

6.1.2 Apply CLT to compute $p(z_k^l(x))$

let's pick any arbitrary x , since we only pick a single x , so the index is **not** important, there is no need to use $x^{(1)}$ like in the literature:
 computing $p(z_k^l(x))$ directly is hard!
 however, $z_k^l(x)$ is $b_k^l + \text{sum of i.i.d elements using CLT notations:}$

$$z_k^l(x) = b_k^l + \underbrace{\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x))}_{\sum_{j=1}^{N_l} X_j}, \quad \text{note we are not taking average} \quad (54)$$

therefore, we can just compute mean and variance of its individual element, i.e., an arbitrary $j = 1$ and then apply CLT!

$$X_j \equiv W_{k,j}^l \phi(z_j^{l-1}(x)) \quad (55)$$

6.1.3 mean and variance of $W_{k,j}^l \phi(z_j^{l-1}(x))$

Expectation

$$\begin{aligned} \mathbb{E}[W_{k,j}^l \phi(z_j^{l-1}(x))] &= \mathbb{E}[W_{k,j}^l] \mathbb{E}[\phi(z_j^{l-1}(x))] \quad \text{since } W_{k,j}^l \text{ and } \phi(z_j^{l-1}(x)) \text{ are independent} \\ &\quad \text{as } z_j^{l-1}(x) \text{ depends on } (W^{l-1}, b^{l-1}) \\ &= 0 \times \mathbb{E}[\phi(z_j^{l-1}(x))] \quad \text{because we choose } W_{k,j}^l \sim \mathcal{N}(0, \sigma_w) \\ &= 0 \end{aligned} \quad (56)$$

Variance

$$\begin{aligned} \text{Var}[W_{k,j}^l \phi(z_j^{l-1}(x))] &= \mathbb{E}\left[\left(W_{k,j}^l \phi(z_j^{l-1}(x))\right)^2\right] \\ &= \mathbb{E}[(W_{k,j}^l)^2] \mathbb{E}[\phi(z_j^{l-1}(x))^2] \quad \text{since } W_{k,j}^l \text{ and } \phi(z_j^{l-1}(x)) \text{ are independent} \\ &= \sigma_w^2 \underbrace{\mathbb{E}[\phi(z_j^{l-1}(x))^2]}_{\text{bounded}} \implies \text{Var}[W_{k,j}^l \phi(z_j^{l-1}(x))] \text{ to be bounded} \\ &= \sigma_w^2 \mathbb{E}[\phi(z_j^{l-1}(x))^2] \end{aligned} \quad (57)$$

we leave in this form, as

$$\mathbb{E}[\phi(z_j^{l-1}(x))^2] \equiv \mathbb{E}_{W^{l-1}, \dots, b^{l-1}, \dots} [\phi(z_j^{l-1}(x))^2] \quad (58)$$

6.1.4 apply CLT:

However, we can apply CLT: making $p(z^l(x))$ distributed as Gaussian where its variance is dependent on variance of previous layer, a recursion.

$$\begin{aligned}
& \text{using } \left(\sum_{i=1}^{\textcolor{red}{n}} X_i - \textcolor{red}{n} \mathbb{E}[X_1] \right) \xrightarrow{d} \mathcal{N}(0, \textcolor{red}{n} \sigma^2) \\
& \Rightarrow \left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) - 0 \right) \sim \mathcal{N} \left(0, \textcolor{red}{N}_l \sigma_w^2 \mathbb{E}[\phi(z_1^{l-1}(x))^2] \right) \quad N_l \rightarrow \infty
\end{aligned} \tag{59}$$

However, variance under this expression $N_l \sigma_w^2 \mathbb{E}[\phi(z_1^{l-1}(x))^2]$ is divergent because of $\textcolor{red}{N}_l$! luckily, we can take control the choice of σ_w^2 , if we let:

$$\sigma(W_{k,j}^l) = \sigma_w = \frac{1}{\sqrt{N_l}} \quad \Rightarrow \quad \sigma_w^2 = \frac{1}{\textcolor{blue}{N}_l} \tag{60}$$

the above is the key, implication is:

$$\begin{aligned}
& \Rightarrow \left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) - 0 \right) \sim \mathcal{N} \left(0, \textcolor{red}{N}_l \frac{1}{\textcolor{blue}{N}_l} \mathbb{E}[\phi(z_1^{l-1}(x))^2] \right) \\
& \quad = \mathcal{N} \left(0, \underbrace{\mathbb{E}[\phi(z_1^{l-1}(x))^2]}_{\text{bounded}} \right)
\end{aligned} \tag{61}$$

finally adding the bias b_k^l :

Note that sum of two **independent** Gaussian random variables is also Gaussian: (not to confuse with GMM!)

$$\begin{aligned}
X & \sim \mathcal{N}(\mu_X, \sigma_X^2) \\
Y & \sim \mathcal{N}(\mu_Y, \sigma_Y^2) \\
Z = X + Y & \quad Z = X + Y \\
\Rightarrow Z & \sim \mathcal{N}(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)
\end{aligned} \tag{62}$$

Therefore:

$$\left(z_k^l(x) = b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) \right) \xrightarrow{d} \mathcal{N} \left(0, \underbrace{\sigma_b^2}_{\sigma_X^2} + \underbrace{\mathbb{E}[\phi(z_1^{l-1}(x))^2]}_{\sigma_Y^2} \right) \quad \text{as } N_l \rightarrow \infty \tag{63}$$

appreciate the recursion here

6.2 given two inputs $\textcolor{red}{x}, \textcolor{blue}{x}'$: compute $\text{Cov}[z_k^l(\textcolor{red}{x}) \ z_k^l(\textcolor{blue}{x}')]]$

To do so, we need to used results from **Multidimensional CLT**:

6.2.1 Multidimensional CLT:

let $\mathbf{X}_i \in \mathcal{R}^d$:

$$\begin{aligned}
\sum_{i=1}^n \mathbf{X}_i &= \underbrace{\begin{bmatrix} X_1^{(1)} \\ \vdots \\ X_1^{(p)} \\ \vdots \\ X_1^{(q)} \\ \vdots \\ X_1^{(k)} \end{bmatrix}}_{\mathbf{X}_1} + \underbrace{\begin{bmatrix} X_2^{(1)} \\ \vdots \\ X_2^{(p)} \\ \vdots \\ X_2^{(q)} \\ \vdots \\ X_2^{(k)} \end{bmatrix}}_{\mathbf{X}_2} + \cdots + \underbrace{\begin{bmatrix} X_n^{(1)} \\ \vdots \\ X_n^{(p)} \\ \vdots \\ X_n^{(q)} \\ \vdots \\ X_n^{(k)} \end{bmatrix}}_{\mathbf{X}_n} = \underbrace{\begin{bmatrix} \sum_{i=1}^n X_i^{(1)} \\ \vdots \\ \sum_{i=1}^n X_i^{(p)} \\ \vdots \\ \sum_{i=1}^n X_i^{(q)} \\ \vdots \\ \sum_{i=1}^n X_i^{(k)} \end{bmatrix}}_{\sum_{i=1}^n \mathbf{X}_i} \\
\Rightarrow \bar{\mathbf{X}} &= \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i = \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n X_i^{(1)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n X_i^{(p)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n X_i^{(q)} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n X_i^{(k)} \end{bmatrix} = \begin{bmatrix} \bar{\mathbf{X}}^{(1)} \\ \vdots \\ \bar{\mathbf{X}}^{(p)} \\ \vdots \\ \bar{\mathbf{X}}^{(q)} \\ \vdots \\ \bar{\mathbf{X}}^{(k)} \end{bmatrix}
\end{aligned} \tag{64}$$

1. **Sample mean version**

$$\begin{aligned}
&\frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i]] \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}[\mathbf{X}_1]) \quad \text{since } p(X_i) = p(X_1) \\
&= \frac{\sqrt{n}}{\sqrt{n}} \frac{1}{\sqrt{n}} \left(\sum_{i=1}^n \mathbf{X}_i \right) - \frac{n}{\sqrt{n}} \mathbb{E}[\mathbf{X}_1] \\
&= \sqrt{n} (\bar{\mathbf{X}} - \mathbb{E}[\mathbf{X}_1])
\end{aligned} \tag{65}$$

$$\begin{aligned}
&\frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i]] = \sqrt{n} (\bar{\mathbf{X}} - \mathbb{E}[\mathbf{X}_1]) \xrightarrow{d} \mathcal{N}_d(0, \boldsymbol{\Sigma}(\mathbf{X}_1)) \\
\Rightarrow \sqrt{n} \mathbb{E} \left[\underbrace{(\bar{\mathbf{X}}^{(p)} - \mathbb{E}[\bar{\mathbf{X}}_1^{(p)}])}_{\text{scalar}} \underbrace{(\bar{\mathbf{X}}^{(q)} - \mathbb{E}[\bar{\mathbf{X}}_1^{(q)}])}_{\text{scalar}} \right] &= \boldsymbol{\Sigma}(\mathbf{X}_1)_{(p),(q)}
\end{aligned} \tag{66}$$

for every elements $(p, q) \in \{1, \dots, k\}$:

2. **Sample sum version:**

$$\begin{aligned}
& \left(\left[\sum_i^n \mathbf{x}_i \right] - n\mathbb{E}[\mathbf{x}_1] \right) \xrightarrow{d} \mathcal{N}_k(0, n\mathbf{\Sigma}) \\
\Rightarrow & \mathbb{E} \left[\left(\left[\sum_i^n \mathbf{x}_i \right]^{(p)} - n\mathbb{E}[\mathbf{x}_1]^{(p)} \right) \left(\left[\sum_i^n \mathbf{x}_i \right]^{(q)} - n\mathbb{E}[\mathbf{x}_1]^{(q)} \right) \right] = n\mathbf{\Sigma}_{(p),(q)} \\
& \Rightarrow \mathbb{E} \left[\left(n\overline{\mathbf{x}}^{(p)} - n\mathbb{E}[X_1^{(p)}] \right) \left(n\overline{\mathbf{x}}^{(q)} - n\mathbb{E}[X_1^{(q)}] \right) \right] = n\mathbf{\Sigma}_{(p),(q)} \\
\Rightarrow & \mathbb{E} \left[\left(\left[\sum_i^n \mathbf{x}_i \right]^{(p)} - n\mathbb{E}[X_1^{(p)}] \right) \left(\left[\sum_i^n \mathbf{x}_i \right]^{(q)} - n\mathbb{E}[X_1^{(q)}] \right) \right] = n\mathbf{\Sigma}_{(p),(q)}
\end{aligned} \tag{67}$$

where $\mathbf{\Sigma}_{(p),(q)} = \text{Cov}(X_1^{(p)}, X_1^{(q)})$

6.2.2 Put into multidimensional CLT structure:

now, let's look at k^{th} dimension of z^l , i.e., z_k^l , and to see in this dimension, how correlation between pair of data input x and x' is. note that what happen to k^{th} dimension, applies to the rest

$$\begin{aligned}
& \begin{bmatrix} \vdots \\ W_{k,1}^l \phi(z_1^{l-1}(x)) \\ \vdots \\ W_{k,1}^l \phi(z_j^{l-1}(x')) \\ \vdots \end{bmatrix} + \cdots + \begin{bmatrix} \vdots \\ W_{k,N_l}^l \phi(z_j^{l-1}(x)) \\ \vdots \\ W_{k,N_l}^l \phi(z_j^{l-1}(x')) \\ \vdots \end{bmatrix} = \underbrace{\begin{bmatrix} \vdots \\ \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) \\ \vdots \\ \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x')) \\ \vdots \end{bmatrix}}_{\begin{bmatrix} \sum_{i=1}^n X_i^{(1)} \\ \vdots \\ \sum_{i=1}^n X_i^{(p)} \\ \vdots \\ \sum_{i=1}^n X_i^{(q)} \\ \vdots \\ \sum_{i=1}^n X_i^{(k)} \end{bmatrix}} = \underbrace{\begin{bmatrix} \vdots \\ z_k^l(x) \\ \vdots \\ z_k^l(x') \\ \vdots \end{bmatrix}}_{\sum_{i=1}^n \mathbf{x}_i}
\end{aligned} \tag{68}$$

compare with the standard notation of Multi-dimensional CLT, and use “sample sum version” of CLT, Eq.[64], and remember $z_k^l(x) = b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x))$, to simplify derivation, let's deliberately not looking at b_k^l for now

$$\begin{aligned}
\sum_i^n X_i^{(p)} &= \left[\sum_i^n X_i \right]^{(p)} \triangleq \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x})) \\
\sum_i^n X_i^{(q)} &= \left[\sum_i^n X_i \right]^{(q)} \triangleq \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x}')) \\
X_1^{(p)} &\triangleq W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x})) \\
X_1^{(q)} &\triangleq W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}')) \\
\implies \Sigma(\mathbf{X}_1)_{(p),(q)} &\triangleq \mathbf{Cov}\left(W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x})), W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}'))\right)
\end{aligned} \tag{69}$$

using above identities in Eq.[69]

$$\begin{aligned}
&\mathbb{E} \left[\left(\left[\sum_i^n \mathbf{X}_i \right]^{(p)} - n \mathbb{E}[X_1^{(p)}] \right) \left(\left[\sum_i^n \mathbf{X}_i \right]^{(q)} - n \mathbb{E}[X_1^{(q)}] \right) \right] = n \Sigma(\mathbf{X}_1)_{(p),(q)} \\
\implies &\mathbb{E} \left[\left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x})) - \underbrace{N_l \mathbb{E}[W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}))]}_{=0} \right) \left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x}')) - \underbrace{N_l \mathbb{E}[W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}'))]}_{=0} \right) \right] \\
&= N_l \mathbf{Cov}\left(W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x})), W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}'))\right) \\
&= N_l \mathbb{E}\left[W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x})) \times W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}'))\right]
\end{aligned} \tag{70}$$

look at $z_k^l(\mathbf{x})$ with b_k^l too:

$$\begin{aligned}
\mathbf{Cov}(z_k^l(\mathbf{x}), z_k^l(\mathbf{x}')) &= \sigma_b^2 + \mathbb{E} \left[\left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x})) \right) \left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(\mathbf{x}')) \right) \right] \\
&= \sigma_b^2 + N_l \mathbf{Cov}(W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x})), W_{k,1}^l \phi(z_1^{l-1}(\mathbf{x}')))) \quad \text{use CLT result above from Eq.[71]} \\
&= \sigma_b^2 + N_l \sigma_w^2 \mathbf{Cov}(\phi(z_1^{l-1}(\mathbf{x})), \phi(z_1^{l-1}(\mathbf{x}')))) \\
&= \sigma_b^2 + N_l \frac{1}{N_l} \mathbf{Cov}(\phi(z_1^{l-1}(\mathbf{x})), \phi(z_1^{l-1}(\mathbf{x}')))) \\
&= \sigma_b^2 + \mathbf{Cov}(\phi(z_1^{l-1}(\mathbf{x})), \phi(z_1^{l-1}(\mathbf{x}')))) \\
&= \sigma_b^2 + \mathbb{E}[\phi(z_1^{l-1}(\mathbf{x})) \times \phi(z_1^{l-1}(\mathbf{x}'))]
\end{aligned} \tag{71}$$

there are many notes this:

1. **note 1** by definition, $\dim(z^l) = N_{l+1}$, so the “entire” $\mathbf{Cov}(z^l, z^l)$ is of size:

$$\underbrace{N_{l+1}}_{\forall k} \underbrace{|\mathcal{X}|}_{\forall x} \times \underbrace{N_{l+1}}_{\forall k'} \underbrace{|\mathcal{X}|}_{\forall x'} \tag{72}$$

exactly how one may arrange this “gigantic” matrix, either N_{l+1} sub-blocks of $\mathbf{Cov}(x, x')$, or $|\mathcal{X}|$ blocks of $\mathbf{Cov}(k, k')$ has the same effect

2. **note 2:** this co-variance is same $\forall k$ in $z_k^l(x)$, so right hand side does not need to keep k index because in this particular setting, since $b_k, b_{k'}, W_{k,j}$ and $W_{k',j'}$ are independent variables, co-variance between any of them are zero:

$$\begin{aligned}
z_k^l(x) &= b_k + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) \\
z_{k'}^l(x) &= b_{k'} + \sum_{j=1}^{N_l} W_{k',j}^l \phi(z_j^{l-1}(x)) \\
\implies \mathbb{E}[W_{k,j}^l \phi(z_j^{l-1}(x)) \times W_{k',j'}^l \phi(z_{j'}^{l-1}(x))] &= 0 \quad \forall \{k, k', j, j'\}
\end{aligned} \tag{73}$$

note 3: in literature, it is written:

$$\begin{aligned}
\mathbb{E}[z_k^l(x) z_{k'}^l(x')] &= \sigma_b^2 + \sigma_w^2 \mathbb{E}\left[\sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))\right] \\
\text{instead of } &= \sigma_b^2 + \mathbb{E}\left[\left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x))\right) \left(\sum_{j=1}^{N_l} W_{k',j}^l \phi(z_j^{l-1}(x'))\right)\right]
\end{aligned} \tag{74}$$

This is because of **note1** above regardless of this special property CLT still apply.

6.2.3 Relationship with Gaussian Process (GP):

let $f(x) \equiv z_k^l(x)$ be some function, and since for every arbitrary point pair, x and x' , we have:

$$\begin{aligned}
\mathbb{E}[f(x)] &= 0 \\
\mathbb{E}[f(x, x')] &= \Sigma_{x, x'} \\
\implies f &\sim \mathcal{GP}(0, \Sigma)
\end{aligned} \tag{75}$$

looking at mean and co-variance as $N_l \rightarrow \infty$

$$\begin{aligned}
\text{Cov}[z_k^l(x), z_k^l(x')] &= \sigma_b^2 + \mathbb{E}[\phi(z_1^{l-1}(x)) \times \phi(z_1^{l-1}(x'))] \quad \text{as } N_l \rightarrow \infty \\
z_k^l(x) &\xrightarrow{d} \mathcal{N}\left(0, \sigma_b^2 + \mathbb{E}[\phi(z_1^{l-1}(x))^2]\right) \quad \text{as } N_l \rightarrow \infty
\end{aligned} \tag{76}$$

putting it in layer specific GP define over some domain \mathcal{X} as $N_l \rightarrow \infty$:

$$\begin{aligned}
\implies z_k^l(\mathcal{X}) &\sim \mathcal{GP}(0, \Sigma^l) \\
\text{where specific co-variance } \Sigma_{x, x'}^l &= \sigma_b^2 + \mathbb{E}[\phi(z_1^{l-1}(x)) \times \phi(z_1^{l-1}(x'))]
\end{aligned} \tag{77}$$

6.3 looking at GP systematically

First let's change for the rest of the tutorial:

$$\Sigma^l \rightarrow K^l \quad (78)$$

$K^l(x, x')$ in terms of pre-activation $z_k^l(x)$ in this section, it will be changed later to post-activation. instead of letting $\sigma(W_{k,j}^l) = \frac{1}{\sqrt{N_l}}$ in previous section, we let it be more generically:

$$\sigma(W_{k,j}^l) = \frac{\sigma_w}{\sqrt{N_l}} \quad (79)$$

we look at all GP kernel K^l relate to K^{l-1} :

$$\begin{aligned} K^l(x, x') &= \mathbb{E}[z_k^l(x) z_k^l(x') | z^{l-1}] \\ &= \sigma_b^2 + \sigma_w^2 \mathbb{E}[\phi(z_1^{l-1}(x)) \times \phi(z_1^{l-1}(x'))] \quad \text{apply CLT } N_l \rightarrow \infty \\ &= \sigma_b^2 + \sigma_w^2 \underbrace{\mathbb{E}_{z_1^{l-1}(\mathcal{X}) \sim \mathcal{GP}(0, K^{l-1})} [\phi(z_1^{l-1}(x)) \phi(z_1^{l-1}(x'))]}_{F_\phi(K^{l-1})} \end{aligned} \quad (80)$$

since $\mathbb{E}[\phi(z)] = \mathbb{E}_{z \sim p(z)}[\phi(z)]$ and $p(z_1^{l-1}(\mathcal{X})) = \mathcal{GP}(0, K^{l-1})$ just as Eq.[77], and $\phi(z_1^{l-1}(x))$ is function on a specific point x , keep on going:

$$\begin{aligned} &= \sigma_b^2 + \sigma_w^2 \underbrace{F_\phi(K^{l-1}(x, x'), K^{l-1}(x, x), K^{l-1}(x', x'))}_{F_\phi(K^{l-1})} \\ &= \sigma_b^2 + \sigma_w^2 F_\phi(K^{l-1}(x, x')) \end{aligned} \quad (81)$$

using properties of point Marginals of Gaussian Process:

$$\begin{aligned} F_\phi(K^{l-1}(x, x')) &= \mathbb{E}_{z_j^{l-1} \sim \mathcal{GP}(0, K^{l-1})} [\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))] \\ &= \mathbb{E} \left(\underbrace{(z_j^{l-1}(x), z_j^{l-1}(x'))}_{\substack{\text{2 points on function } z_j^{l-1} \\ \text{2D Gaussian}}} \sim \mathcal{N}(0, K^{l-1}(x, x')) \right) [\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))] \end{aligned} \quad (82)$$

$$\begin{bmatrix} z_j^{l-1}(x) \\ z_j^{l-1}(x') \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K^{l-1}(x, x) & K^{l-1}(x, x') \\ K^{l-1}(x, x') & K^{l-1}(x', x') \end{bmatrix}\right) \quad (83)$$

assume z^{l-1} can be integrated out:

$$= F_\phi(K^{l-1}(x, x'), K^{l-1}(x, x), K^{l-1}(x', x')) \quad (84)$$

6.4 in summary

this is how K^l relates to K^{l-1} :

$$K^l(x, x') = \sigma_b^2 + \sigma_w^2 \mathbb{E}_{(z_j^{l-1}(x), z_j^{l-1}(x')) \sim \mathcal{N}(0, K^{l-1}(x, x'))} [\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))] \quad (85)$$

we will see the same recursion also applies in NTK, except $\phi \rightarrow \phi'$

7 Expand GP across all layers

7.1 Overall objective

Looking the probability of the final layer output z^L depending on input x :

$$\begin{aligned} p(z^L|x) &= \int p(z^L, K^0, K^1, \dots, K^L|x) dK^0, \dots, L \\ &= \int p(z^L|K^L) \left(\prod_{l=1}^L p(K^l|K^{l-1}) \right) p(K^0|x) dK^0, \dots, L \end{aligned} \quad (86)$$

7.2 $p(z^L|K^L)$: conditions on $K^l \equiv \{\phi(z^{l-1})(x))\phi(z^{l-1})(x'))\}_{p,q}$

(J. H. Lee et. all 2018) presents an **alternative** definition of K^l , where no longer define K from pre-activation:

$$K^l(x, x') = \mathbb{E}[z_k^l(x)z_k^l(x') | z^{l-1}] \quad (87)$$

instead it define K^l in terms of post-activation of previous later $\phi(z^{l-1})$ for reason illustrated later look at Neural Network function:

$$z_k^l(x) = b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) \quad (88)$$

let's make it dependent on $\{\phi(z_j^{l-1}(x))\}_j^{N_l}$, i.e.:
Conditional Marginal

$$\begin{aligned} z_k^l(x) | \{\phi(z_j^{l-1}(x))\}_j^{N_l} &= b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \underbrace{\phi(z_j^{l-1}(x))}_{\text{constant}} \\ \implies z_k^l(x) | \{\phi(z_j^{l-1}(x))\}_j^{N_l} &\sim \mathcal{N}\left(0, \sigma_b^2 + \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x))^2 \text{Var}[W_{k,j}^l]\right) \\ &= \mathcal{N}\left(0, \sigma_b^2 + \frac{\sigma_w^2}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x))^2\right) \end{aligned} \quad (89)$$

using property of weighted sum of Gaussian:

$$\begin{aligned} X_i &\sim \mathcal{N}(\mu_i, \sigma_i^2), \quad i = 1, \dots, \\ \implies \sum_{i=1}^n a_i X_i &\sim \mathcal{N}\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \text{Var}[X_i]\right) \end{aligned} \quad (90)$$

Conditional Co-variance

$$\begin{aligned}
& \text{Cov} \left[z_k^l(x), z_k^l(x') \mid \left\{ \phi(z_j^{l-1}(x)), \phi(z_j^{l-1}(x')) \right\}_{j=1}^{N_l} \right] \\
&= \mathbb{E} \left[z_k^l(x) z_k^l(x') \mid \left\{ \phi(z_j^{l-1}(x)), \phi(z_j^{l-1}(x')) \right\}_{j=1}^{N_l} \right] \\
&= \sigma_b^2 + \mathbb{E}_{W_{k,j}^l} \left[\sum_{j=1}^{N_l} W_{k,j}^l \underbrace{\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))}_{\text{constant, used as condition}} \right] \\
&= \sigma_b^2 + \sum_{j=1}^{N_l} \text{Var}[W_{k,j}^l] \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x')) \\
&= \sigma_b^2 + \frac{\sigma_w^2}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))
\end{aligned} \tag{91}$$

not using property of weighted sum of Gaussian:
Combine all together

$$\begin{aligned}
\text{Cov} \left[z_k^l(x), z_k^l(x') \mid \left\{ \phi(z_j^{l-1}(x)), \phi(z_j^{l-1}(x')) \right\}_{j=1}^{N_l} \right] &= \sigma_b^2 + \sigma_w^2 \frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x')) \\
z_k^l(x) \mid \left\{ \phi(z_j^{l-1}(x)) \right\}_j^{N_l} &\sim \mathcal{N} \left(0, \sigma_b^2 + \sigma_w^2 \frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x))^2 \right) \\
\Rightarrow \begin{bmatrix} z^l(x) \\ z^l(x') \end{bmatrix} \mid \begin{bmatrix} \phi(z_j^{l-1}(x)) \\ \phi(z_j^{l-1}(x')) \end{bmatrix}_j^{N_l} &\sim \mathcal{N} \left(\mathbf{0}, G \left(\begin{bmatrix} K^l(x, x) & K^l(x, x') \\ K^l(x, x') & K^l(x', x') \end{bmatrix} \right) \right)
\end{aligned} \tag{92}$$

in GP paradigm:

$$z^l(x) \mid K^l \sim \mathcal{GP}(z^l; \mathbf{0}, G(K^l)) \tag{93}$$

where

$$\begin{aligned}
K^l(x, x') &= \frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x')) \\
G(K^l(x, x')) &= \sigma_b^2 + \sigma_w^2 K^l(x, x')
\end{aligned} \tag{94}$$

Conveniently, we use K^l as a short-notation collection of $\phi(z_j^{l-1}(x)), \phi(z_j^{l-1}(x')) \quad \forall x, x', j$
also taking care of the layer one, which is just input x :

$$K_{p,q}^l \equiv K^l(x, x') = \begin{cases} \frac{1}{d_{\text{in}}} \sum_{j=1}^{d_{\text{in}}} x_j x'_j = \frac{1}{d_{\text{in}}} x^\top x' & l = 0 \\ \frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x')) & l > 0 \end{cases} \tag{95}$$

to reflect:

$$\text{Cov}(z_k^l, z_{k'}^l) = 0 \quad \forall k, k' \in \{1, \dots, N_{l+1}\} \tag{96}$$

note that

$$K^0(x, x') = \frac{1}{d_{\text{in}}} x^\top x' \quad \text{appears again in NTK} \quad (97)$$

one may construct giant co-variance matrix with $N_{l+1} \times N_{l+1}$ diagonal blocks:

$$\begin{aligned} \mathbf{z}^l &= \underbrace{\begin{bmatrix} \color{red}{z_1^l(x^{(1)})} & \color{red}{z_1^l(x^{(2)})} & \dots & z_1^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_j^l(x^{(1)}) & z_j^l(x^{(2)}) & \dots & z_j^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_{N_{l+1}}^l(x^{(1)}) & z_{N_{l+1}}^l(x^{(2)}) & \dots & z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix}}_{|\mathcal{D}|} \left. \vphantom{\begin{bmatrix} \color{red}{z_1^l(x^{(1)})} & \color{red}{z_1^l(x^{(2)})} & \dots & z_1^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_j^l(x^{(1)}) & z_j^l(x^{(2)}) & \dots & z_j^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_{N_{l+1}}^l(x^{(1)}) & z_{N_{l+1}}^l(x^{(2)}) & \dots & z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix}} \right\} \text{width} \Rightarrow \text{vec}(\mathbf{z}^l) = \begin{bmatrix} \color{red}{z_1^l(x^{(1)})} \\ z_2^l(x^{(1)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(1)}) \\ \color{red}{z_1^l(x^{(2)})} \\ z_2^l(x^{(2)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(2)}) \\ \vdots \\ z_1^l(x^{(|\mathcal{D}|)}) \\ z_2^l(x^{(|\mathcal{D}|)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix} \\ \\ \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} G(K_{1,1}^l) & \dots & 0 & \dots & G(K_{1,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{1,1}^l) & \dots & 0 & 0 & G(K_{1,|\mathcal{D}|}^l) \\ \color{red}{G(K_{2,1}^l)} & \dots & 0 & \dots & G(K_{2,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{2,1}^l) & \dots & 0 & 0 & G(K_{2,|\mathcal{D}|}^l) \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ G(K_{|\mathcal{D}|,1}^l) & \dots & 0 & \dots & G(K_{|\mathcal{D}|,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{|\mathcal{D}|,1}^l) & \dots & 0 & 0 & G(K_{|\mathcal{D}|,|\mathcal{D}|}^l) \end{bmatrix} \right) \\ \Rightarrow p(\mathbf{z}^l | K^l) &= \mathcal{N}(\mathbf{0}, G(K^l) \otimes \mathbf{I}_{N_{l+1} \times N_{l+1}}) \\ &= \mathcal{GP}(\mathbf{z}^l; \mathbf{0}, G(K^l)) \end{aligned} \quad (98)$$

7.3 $p(K^l | K^{l-1})$

Use marginal property of GP and look at: $p(K^l | K^{l-1})$:

$$\begin{aligned}
p(K^l | K^{l-1}) &= \int_{z^{l-1}} p(K^l | z^{l-1}) p(z^{l-1} | K^{l-1}) \\
&= \int_{z^{l-1}} p(K^l | z^{l-1}) \mathcal{GP}(z^{l-1}; 0, G(K^{l-1}))
\end{aligned} \tag{99}$$

using GP property, and just look at two points x, x' :

$$\begin{aligned}
p(K_{p,q}^l | K_{p,q}^{l-1}) &= \int_{z^{l-1}(x), z^{l-1}(x')} p\left(\frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^l(x)) \phi(z_j^l(x'))\right) \\
&\quad \mathcal{N}\left(\begin{bmatrix} z^{l-1}(x) \\ z^{l-1}(x') \end{bmatrix}; 0, G\left(\begin{bmatrix} K^{l-1}(x, x) & K^{l-1}(x, x') \\ K^{l-1}(x, x') & K^{l-1}(x', x') \end{bmatrix}\right)\right)
\end{aligned} \tag{100}$$

7.3.1 what happen to sum $\sum_{j=1}^{N_l} \phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))$ as $N_l \rightarrow \infty$ using CLT:

look at $K_{p,q}^l$ and notice it's sum of iid random variable $K_{p,q}^{l,j}$:

$$\begin{aligned}
\underbrace{K_{p,q}^l}_{\bar{X}} &= \frac{1}{N_l} \sum_{j=1}^{N_l} \underbrace{\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x'))}_{X_j \equiv K_{p,q}^{l,j}} \\
\Rightarrow p(K_{p,q}^{l,1} | K_{p,q}^{l-1}) &= \int_{z^{l-1}(x), z^{l-1}(x')} p(\phi(z_j^l(x)) \phi(z_j^l(x'))) \\
&\quad \mathcal{N}\left(\begin{bmatrix} z^{l-1}(x) \\ z^{l-1}(x') \end{bmatrix}; 0, G\left(\begin{bmatrix} K^{l-1}(x, x) & K^{l-1}(x, x') \\ K^{l-1}(x, x') & K^{l-1}(x', x') \end{bmatrix}\right)\right) \\
&= (F \circ G)(K_{p,q}^{l-1})
\end{aligned} \tag{101}$$

using CLT, pick the most appropriate definition:

$$(\bar{X} - \mathbb{E}[X_1]) \xrightarrow{d} \mathcal{N}\left(0, \frac{\text{Var}[X_1]}{n}\right) \tag{102}$$

let's see what is $\lim_{N_l \rightarrow \infty} p(K^l | K^{l-1})$:

$$\begin{aligned}
&(\bar{X} - \mathbb{E}[X_1]) \xrightarrow{d} \mathcal{N}\left(0, \frac{\text{Var}[X_1]}{n}\right) \\
\Rightarrow (K_{p,q}^l - \mathbb{E}[K_{p,q}^{l,1}]) &\xrightarrow{d} \mathcal{N}\left(0, \frac{\text{Var}[K_{p,q}^{l,1}]}{N_l}\right) \\
\Rightarrow (K_{p,q}^l - (F \circ G)(K_{p,q}^{l-1})) &\xrightarrow{d} \mathcal{N}\left(0, \frac{\text{Var}[K_{p,q}^{l,1}]}{N_l}\right) \\
\Rightarrow (K_{p,q}^l | K_{p,q}^{l-1}) &\xrightarrow{d} \mathcal{N}\left((F \circ G)(K^{l-1}), \frac{\text{Var}[K_{p,q}^{l,1}]}{N_l}\right) \\
\Rightarrow \lim_{N_l \rightarrow \infty} p(K^l | K^{l-1}) &= \delta(K^l - (F \circ G)(K^{l-1})) \quad \text{entire matrix}
\end{aligned} \tag{103}$$

note using CLT, sample mean converge to δ_μ , can be exploited for other application
note that this single step conditional is quite easy

7.4 putting in the overall objective function

let width of all layers to $\rightarrow \infty$:

$$\begin{aligned}
p(z^L|x) &= \int p(z^L, K^0, K^1, \dots, K^L|x) \, dK^{0,\dots,L} \\
&= \int p(z^L|K^L) \left(\prod_{l=1}^L p(K^l|K^{l-1}) \right) p(K^0|x) \, dK^{0,\dots,L} \\
\lim_{N_L \rightarrow \infty, \dots, N_1 \rightarrow \infty} p(z^L|x) &= \int p(z^L|K^L) \left(\prod_{l=1}^L \delta(K^l - (F \circ G)(K^{l-1})) \right) p(K^0|x) \, dK^{0,\dots,L} \\
&= \int \mathcal{GP}(z^L; 0, G(K^L)) \underbrace{\left(\prod_{l=1}^L \delta(K^l - (F \circ G)(K^{l-1})) \right) \delta\left(K^0 - \frac{1}{d_{\text{in}}} x^\top x\right)}_{= \begin{cases} 1 & \text{if } K^L = (F \circ G)(K^{L-1}) \\ & = (F \circ G)^2(K^{L-2}) \dots \\ & = (F \circ G)^L\left(\frac{1}{d_{\text{in}}} x^\top x\right) \\ 0 & \text{otherwise} \end{cases}} \, dK^{0,\dots,L} \\
&= \mathcal{GP}\left(z^L; 0, G \circ (F \circ G)^L\left(\frac{1}{d_{\text{in}}} x^\top x\right)\right)
\end{aligned} \tag{104}$$

8 Gradient Flow in terms of Neural Tangent Kernel

8.1 The problem

since Cost (or output layer) can be defined in convex function terms of post-activation last layer neurons $\phi(z^L(x))$, for example:

$$C = \|y - \phi(z^L(x))\|^2 \quad (105)$$

there must be a global minimal if we were to optimize it in term of $\phi(z^L(x))$
however, current training regime:

1. gradient descend are **not** optimized using $\frac{\delta C}{\delta \phi(z^L(x))}$
2. but, it is computed through $\frac{\partial \delta C}{\partial \theta}$

so it's unclear if a tiny step taken when $\theta(t) \rightarrow \theta(t + \epsilon)$ is to lead towards a negative gradient value in $\frac{\delta C}{\delta \phi(z^L(x))}$

8.2 what do you hope for functional gradient $\frac{\delta C}{\delta f(\theta)}$

Under any training regime, there will be parameter dynamics (gradient flow) $\frac{d\theta}{dt}$
what you hope: under **gradient descend** with infinitesimal step size (a.k.a. gradient flow)

$$\frac{d\theta}{dt} = -\frac{\partial C}{\partial \theta} \quad (106)$$

functional gradient $\frac{\delta C}{\delta f(\theta)}$ is **negative all the time!**

because C is convex functional of $f(\theta)$, and if gradient is negative all the time, it will eventually reach the global minima
and **no!** this doesn't work all the time, it only occur under specific conditions listed below:

8.2.1 $\frac{\delta C}{\delta f(\theta)}$ under arbitrary infinitesimal step change $\theta \rightarrow \theta + \epsilon \eta$

So the **question** is, when θ undertakes infinitesimal step change in a **direction vector** η , i.e.,:

$$\theta \rightarrow \theta + \epsilon \eta \quad (107)$$

how does $\frac{\delta C}{\delta f(\theta)}$ change. Formally, we want to compute the following limit:

$$\lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon \eta)] - C[f(\theta)]}{\epsilon} \quad (108)$$

it is a mathematical traditional to write functional $C[f]$ is in square bracket
Since $C[f]$ is a functional, we need to use **Riesz-Markov-Kakutani Representation Theorem**:

$$\int_{\mathbf{x}} \frac{\delta J}{\delta g}(x)^\top \phi(x) dx = \lim_{\epsilon \rightarrow 0} \frac{J[g + \epsilon \phi] - J[g]}{\epsilon} \quad (109)$$

if g was a variable instead of a function, then, the above is analogous to:

$$\phi^\top \nabla_g J \quad (110)$$

i.e., directional derivative of J in the direction of ϕ , and there is no integral $\int_{\mathbf{x}} dx$!
we can **not** substitute into RMK Representation directly, because our changes $\epsilon\eta$ occur in f 's argument:

$$\lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon\eta)] - C[f(\theta)]}{\epsilon} \quad (111)$$

But we must get it in the form of $C[f(\theta) + \epsilon\eta]$. Therefore, we need to use Taylor Expansion:

$$\begin{aligned} & C \left[\underbrace{f(\theta)}_g + \epsilon \underbrace{\eta \cdot \frac{\partial f(\theta)}{\partial \theta}}_{\phi} + O(\epsilon^2) \right] - C[f(\theta)] \\ \Rightarrow & \lim_{\epsilon \rightarrow 0} \frac{\quad}{\epsilon} \quad \text{matching with RMK representation} \\ = & \int_{\mathbf{x}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)^\top \left(\eta \cdot \frac{\partial f(\theta)}{\partial \theta} \right) dx \\ = & \sum_{d=1}^{|\theta|} \int_{\mathbf{x}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)^\top \left(\eta \cdot \frac{\partial f(\theta)}{\partial \theta_d} \right) dx \\ = & \sum_{d=1}^{|\theta|} \int_{\mathbf{x}} \sum_{i=1}^N \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\eta \cdot \frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \\ = & \sum_{d=1}^{|\theta|} \sum_{i=1}^N \eta \int_{\mathbf{x}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \quad \text{change order of integral and sum} \end{aligned} \quad (112)$$

8.2.2 $\frac{\delta C}{\delta f(\theta)}$ under gradient flow in gradient descent training

above tells how much does C change if $\theta \rightarrow \theta + \epsilon\eta$

since we can choose any direction η , we can equally (and meaningfully) choose a direction to be **gradient flow**, i.e.:

$$\eta \equiv \frac{\partial \theta}{\partial t} \quad (113)$$

which correspond to the training regime used
by substitution:

$$\lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon\eta)] - C[f^\theta]}{\epsilon} = \sum_{d=1}^{|\theta|} \sum_{i=1}^N \left(\frac{\partial \theta}{\partial t} \right) \int_{\mathbf{x}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \quad (114)$$

if **gradient descent** training regime is used, then:

$$\begin{aligned} \frac{\partial \theta}{\partial t} &= - \frac{\partial C[f(\theta)]}{\partial \theta} \\ &= - \lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon \mathbf{I})] - C[f^\theta]}{\epsilon} \\ &= - \sum_{d'=1}^{|\theta|} \sum_{k=1}^N \int_{\mathbf{x}'} \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left(\frac{\partial f_k(\theta)}{\partial \theta_{d'}} \right)_k dx' \quad \text{change index to } k \text{ and } x \rightarrow x' \end{aligned} \quad (115)$$

substitution:

$$\begin{aligned}
& \lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon \eta)] - C[f^\theta]}{\epsilon} \\
&= \sum_{d=1}^{|\theta|} \sum_{i=1}^N \left(\frac{\partial \theta}{\partial t} \right) \int_{\mathbf{X}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \\
&= - \sum_{d=1}^{|\theta|} \sum_{i=1}^N \left[\sum_{d'=1}^{|\theta|} \sum_{k=1}^N \int_{\mathbf{X}'} \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left(\frac{\partial f(\theta)}{\partial \theta_{d'}} \right)_k dx' \right] \left[\int_{\mathbf{X}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \right] \\
&= - \sum_{d=1}^{|\theta|} \sum_{i=1}^N \left[\sum_{k=1}^N \int_{\mathbf{X}'} \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_k dx' \right] \int_{\mathbf{X}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \\
&= - \sum_{d=1}^{|\theta|} \sum_{i=1}^N \sum_{k=1}^N \left(\int_{\mathbf{X}'} \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_k dx' \right) \left(\int_{\mathbf{X}} \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx \right) \\
&\quad \text{can take sum } \sum_k \text{ out of bracket because second term has no } k \\
&= - \sum_{d=1}^{|\theta|} \sum_{i=1}^N \sum_{k=1}^N \int_{\mathbf{X}'} \int_{\mathbf{X}} \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_k \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i dx dx' \\
&= - \int_{\mathbf{X}'} \int_{\mathbf{X}} \underbrace{\sum_{i=1}^N \sum_{k=1}^N \left(\frac{\delta C}{\delta f(\theta)}(x') \right)_k \left[\sum_{d=1}^{|\theta|} \left(\frac{\partial f_k(\theta)}{\partial \theta_d} \right)_k \left(\frac{\partial f(\theta)}{\partial \theta_d} \right)_i \right] \left(\frac{\delta C}{\delta f(\theta)}(x) \right)_i}_{\sum_i \sum_j x_i x_j M_{i,j} = \bar{x}^\top M \bar{x}} dx dx' \\
&= - \int_{\mathbf{X}'} \int_{\mathbf{X}} \underbrace{\left(\frac{\delta C}{\delta f(\theta)}(x') \right)^\top}_{\Theta(x, x')} \left(\frac{\delta C}{\delta f(\theta)}(x) \right) dx dx'
\end{aligned} \tag{116}$$

note $\Theta(x, x')$ above has nothing to do Neural Networks, i.e., the above is true under gradient descent regardless of $f(\theta)$ used

8.2.3 What happens $\Theta(x, x')$ is positive definite

the above implies that **if** NTK is positive definite (which is the NTK paper is all about):

$$\lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon \frac{\partial \theta}{\partial t})] - C[f^\theta]}{\epsilon} = \text{negative value} \tag{117}$$

cost will converge to a global optima.

it is important to know the term **inside** the integral is actually **not** guaranteed to be positive.

It is only become positive when the integrals are taken. To make it clear, we rewrite the following using simple notations:

$$\int_x \int_{x'} \underbrace{\bar{f}(x)^\top \Theta(x, x') \bar{f}(x')} \tag{118}$$

for a specific term

$$\bar{f}(x)^\top K(x, x') \bar{f}(x') \tag{119}$$

it may not be positive as left vector $\bar{f}(x)$ and right vector $\bar{f}(x')$ may not equate. However, by summing all **four** elements concerning the co-efficient of $\Theta(i, j) \equiv \Theta_{i,j}(x, x')$:

$$\begin{aligned}
A &\equiv \textcolor{red}{f_i(x)}\Theta(i, j)f_j(x) + \textcolor{red}{f_i(x')}\Theta(i, j)f_j(x') = \textcolor{red}{f_i(x)}\Theta(i, j)(f_j(x) + f_j(x')) \\
B &\equiv \textcolor{blue}{f_i(x')}\Theta(i, j)f_j(x) + \textcolor{blue}{f_i(x')}\Theta(i, j)f_j(x') = \textcolor{blue}{f_i(x')}\Theta(i, j)(f_j(x) + f_j(x')) \\
A + B &= \underbrace{(\textcolor{red}{f_i(x)} + \textcolor{blue}{f_i(x')})}_{g_i(x, x')} \Theta(i, j) \underbrace{(f_j(x) + f_j(x'))}_{g_j(x, x')}
\end{aligned} \tag{120}$$

since g is non-specific to value in x and x' , as both are used.
therefore, $K(x, x')$ is positive definitely **condition** on the fact that x and x' are distributed from the same distribution, e.g., p^{in} .
formally, we can write it as:

$$\begin{aligned}
&K \text{ is positive definite with respect to } \|\cdot\|_{p^{\text{in}}} \quad \text{if} \\
&\mathbb{E}_{x, x' \sim p^{\text{in}}} [f(x)^\top f(x')] > 0 \implies \mathbb{E}_{x, x' \sim p^{\text{in}}} [f(x)^\top K f(x')] > 0
\end{aligned} \tag{121}$$

8.2.4 What does NTK paper aims to prove

NTK paper is all about, Proving under gradient descend training regime/gradient field

and with the following conditions:

1. $f(\theta)$ is a neural network
2. θ has appropriate Gaussian initialization is applied
3. having $N_1, \dots, N_L \rightarrow \infty$:

Then,

1. NTK is indeed positive definite, in a Scalar matrix form: “some positive scalar” $\times \mathbf{I}_{N_{l+1}}$
2. remains approximately constant throughout training

Consequently, leading $\lim_{\epsilon \rightarrow 0} \frac{C[f(\theta + \epsilon \frac{\partial \theta}{\partial t})] - C[f(\theta)]}{\epsilon}$ to stay negative, i.e., cost always going down in a convex function, so it will eventually reach global minimum.

8.3 NTK in Neural Networks

we use the re-parameterization version of NN function:

$$z_k^{(l)} = \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_k^l \tag{122}$$

where $W_{k,j}^l, b_k^l \sim \mathcal{N}(0, 1)$

Neural Tangent Kernel at each Layer l :

$$\begin{aligned}
\Theta^l(x, x') &= \sum_{d=1}^{|\theta|} \frac{\partial z^l(x)}{\partial \theta_d} \otimes_{\text{outer}} \frac{\partial z^l(x')}{\partial \theta_d} \\
&= \sum_{d=1}^{|\theta|} \begin{bmatrix} \frac{\partial z_1^l(x)}{\partial \theta_d} & \dots & \frac{\partial z_{N_{l+1}}^l(x)}{\partial \theta_d} \end{bmatrix} \begin{bmatrix} \frac{\partial z_1^l(x')}{\partial \theta_d} & \dots & \frac{\partial z_{N_{l+1}}^l(x')}{\partial \theta_d} \end{bmatrix}^\top \\
&= \sum_{d=1}^{|\theta|} \begin{bmatrix} \frac{\partial z_1^l(x)}{\partial \theta_d} \frac{\partial z_1^l(x')}{\partial \theta_d} & \dots & \frac{\partial z_{N_{l+1}}^l(x)}{\partial \theta_d} \frac{\partial z_{N_{l+1}}^l(x')}{\partial \theta_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_{N_{l+1}}^l(x)}{\partial \theta_d} \frac{\partial z_1^l(x')}{\partial \theta_d} & \dots & \frac{\partial z_{N_{l+1}}^l(x)}{\partial \theta_d} \frac{\partial z_{N_{l+1}}^l(x')}{\partial \theta_d} \end{bmatrix}
\end{aligned} \tag{123}$$

note that size of $\Theta^l(x, x')$ is $N_{l+1} \times N_{l+1}$, it has nothing to do with $|\theta|$ (it is used in the sum)

however, the **entire NTK** Θ^l will be the size $(N_{l+1} \times |\mathcal{D}|) \times (N_{l+1} \times |\mathcal{D}|)$
loosely speaking:

1. NTK studies “pseudo-correlations” between a pair of output (k, k') of a vector function z^l by summing over their derivatives over all parameters from two data x and x'
(derivative correlations between **function’s output**)

which is different to fisher information matrix:

$$\mathbf{I}_{i,j} = \mathbb{E} \left[\left(\frac{\partial}{\partial \theta_i} \log f(X; \theta) \right) \left(\frac{\partial}{\partial \theta_j} \log f(X; \theta) \right) \right] \quad (124)$$

2. where FIM studies correlation between log derivative of pair of parameters (θ_i, θ_j) from a scalar function f .
(derivative correlations between **function’s parameters**)

note that symbol here \otimes above is outer product as oppose to kronecker product everywhere else in this tutorial. But the two are related:

$$\mathbf{u} \otimes_{\text{Kron}} \mathbf{v}^T = \mathbf{u} \mathbf{v}^T = \mathbf{u} \otimes_{\text{outer}} \mathbf{v} \quad (125)$$

9 NTK at initialization

Given a single input x , we show the following is the relationship between two adjacent layers $z^{l-1}(x) \rightarrow z^l(x)$:

$$\begin{aligned}
 & \begin{bmatrix} \frac{1}{\sqrt{N_l}} W_{1,1}^l \phi(z_1^{l-1}(x)) + \sigma_b b_1 \\ \vdots \\ \frac{1}{\sqrt{N_l}} W_{k,1}^l \phi(z_1^{l-1}(x)) + \sigma_b b_k \\ \vdots \\ \frac{1}{\sqrt{N_l}} W_{N_{l+1},1}^l \phi(z_j^{l-1}(x)) + \sigma_b b_{N_{l+1}} \end{bmatrix} + \dots + \begin{bmatrix} \frac{1}{\sqrt{N_l}} W_{1,N_l}^l \phi(z_1^{l-1}(x)) + \sigma_b b_1^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} W_{k,N_l}^l \phi(z_1^{l-1}(x)) + \sigma_b b_k^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} W_{N_{l+1},N_l}^l \phi(z_j^{l-1}(x)) + \sigma_b b_{N_{l+1}}^l \end{bmatrix} \\
 &= \begin{bmatrix} \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{1,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_1^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_k^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{N_{l+1},j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_{N_{l+1}}^l \end{bmatrix} = \begin{bmatrix} z_1^l(x) \\ \vdots \\ z_k^l(x) \\ \vdots \\ z_{N_{l+1}}^l(x) \end{bmatrix}
 \end{aligned} \tag{126}$$

9.1 prove by induction

9.1.1 prove by induction

induction works by proving value at $l = 1$, then show relationship between l and $l - 1$ in general. Finally, it shows what value is at an arbitrary index L

9.1.2 for NTK

we need to show by induction,

1. assume for a small network, at $l = 1$ we prove:

$$\Theta_{k,k'}^1(x, x') = \underbrace{\left(\frac{1}{d_{\text{in}}} x^\top x' + \sigma_b^2 \right)}_{K^1} \delta_{k,k'}$$

even better, no need to show: $\Theta_{k,k'}^1(x, x') \rightarrow K^1 \delta_{k,k'}$, it is actually equal! (127)

2. then by assuming:

$$\Theta_{k,k'}^{l-1}(x, x') = \frac{\partial z_k^{l-1}(x, \theta)}{\partial \theta^l}^\top \frac{\partial z_k^{l-1}(x', \theta)}{\partial \theta^l} \xrightarrow{N_l \rightarrow \infty} \Theta_{\infty}^{l-1}(x, x') \delta_{k,k'} \tag{128}$$

3. we prove:

$$\Theta_{k,k'}^l(x, x') = \frac{\partial z_k^l(x, \theta)}{\partial \theta^l}^\top \frac{\partial z_k^l(x', \theta)}{\partial \theta^l} \xrightarrow{N_{l+1} \rightarrow \infty} \Theta_{\infty}^l(x, x') \delta_{k,k'} \tag{129}$$

9.2 when $l = 1$: $\Theta_{k,k'}^1(x, x') = \left(\frac{1}{d_{\text{in}}} x^\top x' + \sigma_b^2 \right) \delta_{k,k'}$

From the Eq.(126), we have:

$$\begin{bmatrix} \frac{1}{\sqrt{d_{\text{in}}}} \sum_{j=1}^{d_{\text{in}}} W_{1,j}^1 \phi(x_1) + \sigma_b b_1^1 \\ \vdots \\ \frac{1}{\sqrt{d_{\text{in}}}} \sum_{j=1}^{d_{\text{in}}} W_{k,j}^1 \phi(x_k) + \sigma_b b_k^1 \\ \vdots \\ \frac{1}{\sqrt{d_{\text{in}}}} \sum_{j=1}^{d_{\text{in}}} W_{N_1,j}^1 \phi(x_{N_1}) + \sigma_b b_{N_1}^1 \end{bmatrix} = \begin{bmatrix} z_1^1(x) \\ \vdots \\ z_k^1(x) \\ \vdots \\ z_{N_1}^1(x) \end{bmatrix} \quad (130)$$

when computing: $\frac{\partial z_{\bar{k}}^1(x)}{\partial W_{i,j}}$, here, we use i to index entries of W , because \bar{k} is used by $z_{\bar{k}}^1(x)$. note when computing $\frac{\partial z_{\bar{k}}^1(x)}{\partial W_{i,j}}$ only \bar{k} th row going to return a gradient, i.e., $\frac{\partial z_{\bar{k}}^1(x)}{\partial W_{i,j}} = 0$ if $i \neq k$

$$\begin{aligned} \frac{\partial z_k^1(x)}{\partial W_{i,j}} &= \begin{cases} \frac{1}{\sqrt{d_{\text{in}}}} x_i & \text{if } i = k \text{ i.e., row } k \\ 0 & \text{otherwise} \end{cases} \\ &= \frac{1}{\sqrt{d_{\text{in}}}} \delta_{i,k} x_i \\ \Rightarrow \frac{\partial z_{k'}^1(x)}{\partial W_{i,j}} &= \frac{1}{\sqrt{d_{\text{in}}}} \delta_{i,k'} x_i \end{aligned} \quad (131)$$

now, taking pair of data x and x' , each element of the outer product matrix $\Theta^l(x, x') = \sum_{d=1}^{|\theta|} \frac{\partial F_k^l(x)}{\partial \theta_d} \otimes \frac{\partial F_{k'}^l(x')}{\partial \theta_d}$ at k, k' is:

$$\begin{aligned} \Theta_{k,k'}^1(x, x') &= \sum_{d=1}^{|\theta^1|} \frac{\partial F_k^1(x)}{\partial \theta_d^1} \frac{\partial F_{k'}^1(x')}{\partial \theta_d^1} \quad \theta^1 = \{W^1, b^1\} \\ &= \sum_{d=1}^{|W^1|} \frac{\partial F_k^1(x)}{\partial W_d^1} \frac{\partial F_{k'}^1(x')}{\partial W_d^1} + \sum_{d=1}^{|b^1|} \frac{\partial F_k^1(x)}{\partial b_d^1} \frac{\partial F_{k'}^1(x')}{\partial b_d^1} \\ &= \sum_{i=1}^{N_1} \sum_{j=1}^{d_{\text{in}}} \frac{\partial z_k^1(x)}{\partial W_{i,j}} \frac{\partial z_{k'}^1(x')}{\partial W_{i,j}} + \sum_{i=1}^{N_1} \frac{\partial z_k^1(x)}{\partial b_i} \frac{\partial z_{k'}^1(x')}{\partial b_i} \\ &= \sum_{i=1}^{N_1} \sum_{j=1}^{d_{\text{in}}} \frac{1}{\sqrt{d_{\text{in}}}} x_i \delta_{i,k'} \frac{1}{\sqrt{d_{\text{in}}}} x'_i \delta_{i,k} + \sum_{i=1}^{N_1} \sigma_b \delta_{i,k} \sigma_b \delta_{i,k'} \quad \text{only one } i \in \{1, \dots, N_1\} \text{ in outer sum remain} \\ &= \sum_{j=1}^{d_{\text{in}}} \frac{1}{d_{\text{in}}} x_i x'_i \delta_{k,k'}^2 + \sigma_b^2 \delta_{k,k'} \quad \delta_{i,k'} \delta_{i,k} = \delta_{k,k'} \\ &= \frac{1}{d_{\text{in}}} x^\top x' \delta_{k,k'} + \sigma_b^2 \delta_{k,k'} \\ &= \underbrace{\left(\frac{1}{d_{\text{in}}} x^\top x' + \sigma_b^2 \right)}_{K^1} \delta_{k,k'} \\ &\equiv K^1(x, x') \delta_{k,k'} \quad \text{conform to notation in NNGP section} \end{aligned} \quad (132)$$

now we have each element $\Theta_{k,k'}^l$, the final Θ^l is:

$$\begin{aligned} \Rightarrow \Theta^1(x, x') &= \left\{ \underbrace{\begin{bmatrix} K^1(x, x') & \dots & 0 & \dots & 0 \\ 0 & K^1(x, x') & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & \dots & K^1(x, x') & 0 \\ 0 & 0 & 0 & 0 & K^1(x, x') \end{bmatrix}}_{k \in \{1, \dots, N_l\}} \right\}_{k' \in \{1, \dots, N_l\}} \\ &= \text{repeating diagonal with } K^1(x, x') \delta_{k,k'} \\ &= \underbrace{K^1(x, x')}_{\text{scalar}} \otimes_{\text{outer}} \mathbf{I}_{N_l+1 \times N_l+1} \end{aligned} \quad (133)$$

Θ^1 matrix of square the size of input $|z^1|$, importantly, there is no limit to take for Θ^1

9.3 when $l > 1$

$$\begin{bmatrix} \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{1,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_1^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_k^l \\ \vdots \\ \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{N_l+1,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_{N_l+1}^l \end{bmatrix} = \begin{bmatrix} z_1^l(x) \\ \vdots \\ z_k^l(x) \\ \vdots \\ z_{N_l+1}^l(x) \end{bmatrix} \quad (134)$$

split sum into two parts: $\{W^l, b^l\}$ and θ^{l-1}

$$\begin{aligned} \Theta_{k,k'}^l(x, x') &= \sum_{d=1}^{|\theta^l|} \frac{\partial z_k^1(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{k'}^l(x')}{\partial \theta_d^{l-1}} \\ &= \sum_{d=1}^{|W^l, b^l|} \frac{\partial z_k^1(x)}{\partial \{W^l, b^l\}} \frac{\partial z_{k'}^l(x')}{\partial \{W^l, b^l\}} + \sum_{d=1}^{|\theta^{l-1}|} \frac{\partial z_k^1(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{k'}^l(x')}{\partial \theta_d^{l-1}} \end{aligned} \quad (135)$$

looking at this specific term: $\frac{\partial z_k^1(x)}{\partial \theta_d^{l-1}}$, write $x \equiv x$, and definition again:

$$\begin{aligned} z_k^l &= \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x)) + \sigma_b b_k^l \\ &= \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi\left(\frac{1}{\sqrt{N_{l-1}}} \sum_{i=1}^{N_{l-1}} W_{j,i}^{l-1} \phi(z_i^{l-1}(x)) + \sigma_b b_j^{l-1}\right) + \sigma_b b_k^l \end{aligned} \quad (136)$$

$$\begin{aligned}
\frac{\partial z_k^l(x)}{\partial \theta_d^{l-1}} &= \frac{\partial z_k^1(x)}{\partial \phi(z^{l-1}(x))} \frac{\partial \phi(z^{l-1}(x))}{\partial z^{l-1}(x)} \frac{\partial z^{l-1}(x)}{\partial \theta_d^{l-1}} \quad \text{drop index for the last two terms} \\
&= \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \frac{\partial \phi(z_j^{l-1}(x))}{\partial z_j^{l-1}(x)} \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \\
&= \frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi'(z_j^{l-1}(x)) \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \quad \text{leave last derivative as is, in "recursion"}
\end{aligned} \tag{137}$$

substitution:

$$\begin{aligned}
&\sum_{d=1}^{|\theta^{l-1}|} \frac{\partial z_k^1(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{k'}^l(x')}{\partial \theta_d^{l-1}} \\
&= \sum_{d=1}^{|\theta^{l-1}|} \left(\frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k,j}^l \phi'(z_j^{l-1}(x)) \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \right) \times \left(\frac{1}{\sqrt{N_l}} \sum_{j=1}^{N_l} W_{k',j}^l \phi'(z_j^{l-1}(x')) \frac{\partial z_j^{l-1}(x')}{\partial \theta_d^{l-1}} \right) \\
&= \sum_{d=1}^{|\theta^{l-1}|} \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{j'=1}^{N_l} \left(W_{k,j}^l \phi'(z_j^{l-1}(x)) \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \right) \times \underbrace{\left(W_{k',j'}^l \phi'(z_{j'}^{l-1}(x')) \frac{\partial z_{j'}^{l-1}(x')}{\partial \theta_d^{l-1}} \right)}_{j \rightarrow j'} \\
&= \sum_{d=1}^{|\theta^{l-1}|} \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{j'=1}^{N_l} W_{k,j}^l W_{k',j'}^l \phi'(z_j^{l-1}(x)) \phi'(z_{j'}^{l-1}(x')) \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{j'}^{l-1}(x')}{\partial \theta_d^{l-1}} \\
&= \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{j'=1}^{N_l} W_{k,j}^l W_{k',j'}^l \phi'(z_j^{l-1}(x)) \phi'(z_{j'}^{l-1}(x')) \underbrace{\sum_{d=1}^{|\theta^{l-1}|} \frac{\partial z_j^{l-1}(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{j'}^{l-1}(x')}{\partial \theta_d^{l-1}}}_{\text{definition } \Theta_{j,j'}^{l-1}(x,x')} \\
&= \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{j'=1}^{N_l} W_{k,j}^l W_{k',j'}^l \phi'(z_j^{l-1}(x)) \phi'(z_{j'}^{l-1}(x')) \Theta_{j,j'}^{l-1}(x,x') \\
&= \frac{1}{N_l} \sum_{j=1}^{N_l} \sum_{j'=1}^{N_l} W_{k,j}^l W_{k',j'}^l \phi'(z_j^{l-1}(x)) \phi'(z_{j'}^{l-1}(x')) \Theta_{\infty}^{l-1}(x,x') \delta_{j,j'} \\
&\quad \text{substitute induction criteria: } \Theta_{j,j'}^{l-1}(x,x') \rightarrow \underbrace{\Theta_{\infty}^{l-1}(x,x') \delta_{j,j'}}_{\text{deterministic and diagonal limit}} \\
&= \frac{1}{N_l} \sum_{j=1}^{N_l} W_{k,j}^l W_{k',j}^l \phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \Theta_{\infty}^{l-1}(x,x') \quad \text{change } j' \rightarrow j \text{ and remove } \sum_{j'=1}^{N_l} W_{k,j}^l \\
&\tag{138}
\end{aligned}$$

apply CLT, we know that:

$$\begin{aligned}
& \mathbb{E} \left[\sum_{d=1}^{|\theta^{l-1}|} \frac{\partial z_k^l(x)}{\partial \theta_d^{l-1}} \frac{\partial z_{k'}^l(x')}{\partial \theta_d^{l-1}} \right] \\
&= \mathbb{E} \left[W_{k,j}^l W_{k',j}^l \phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \Theta_\infty^{l-1}(x, x') \right] \quad N_l \rightarrow \infty \\
&= \mathbb{E} \left[\underbrace{W_{k,j}^l W_{k',j}^l}_{\text{constant}} \right] \mathbb{E} \left[\phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \right] \Theta_\infty^{l-1}(x, x') \\
&= \delta_{k,k'} \mathbb{E} \left[\phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \right] \Theta_\infty^{l-1}(x, x')
\end{aligned} \tag{139}$$

we have seen previously Eq. (85):

$$K^l(x, x') = \sigma_b^2 + \sigma_w^2 \mathbb{E}_{(z_j^{l-1}(x), z_j^{l-1}(x')) \sim \mathcal{N}(0, K^{l-1}(x, x'))} \left[\phi(z_j^{l-1}(x)) \phi(z_j^{l-1}(x')) \right] \tag{140}$$

however, this time we need to define a similar auxiliary variable \dot{K}^l , notice it **has no** σ_b^2 **term**, describing expectation of $\phi'()$

$$\begin{aligned}
\dot{K}^l(x, x') &= \sigma_w^2 \mathbb{E}_{(z_j^{l-1}(x), z_j^{l-1}(x')) \sim \mathcal{N}(0, K^{l-1}(x, x'))} \left[\phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \right] \\
&= \mathbb{E}_{(z_j^{l-1}(x), z_j^{l-1}(x')) \sim \mathcal{N}(0, K^{l-1}(x, x'))} \left[\phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \right] \quad \text{assume } \sigma_w = 1
\end{aligned} \tag{141}$$

also notice the above equation is **not a recursion**, i.e., $\dot{K}^l(x, x')$ and $K^{l-1}(x, x')$ are not the same thing.

$$\begin{aligned}
&= \delta_{k,k'} \mathbb{E}_{(z_j^{l-1}(x), z_j^{l-1}(x')) \sim \mathcal{N}(0, K^{l-1}(x, x'))} \left[\phi'(z_j^{l-1}(x)) \phi'(z_j^{l-1}(x')) \right] \Theta_\infty^{l-1}(x, x') \\
&= \delta_{k,k'} \dot{K}^l(x, x') \Theta_\infty^{l-1}(x, x')
\end{aligned} \tag{142}$$

look at $\{W^l, b^l\}$ part:

$$\sum_{d=1}^{|W^l, b^l|} \frac{\partial z_k^l(x)}{\partial \{W^l, b^l\}} \frac{\partial z_{k'}^l(x')}{\partial \{W^l, b^l\}} \tag{143}$$

and compare that with for $l = 1$:

$$\sum_{d=1}^{|\theta^1|} \frac{\partial F_k^1(x)}{\partial \theta_d^1} \frac{\partial F_{k'}^1(x')}{\partial \theta_d^1} \quad \theta^1 = \{W^1, b^1\} \tag{144}$$

it's the same if we replace

$$\left(K^1(x, x') \equiv \frac{1}{d_{\text{in}}} x^\top x' + \sigma_b^2 \right) \delta_{k,k'} \rightarrow \left(K^l(x, x') \equiv \frac{1}{N_l} \phi(z^l(x))^\top \phi(z^l(x')) + \sigma_b^2 \right) \delta_{k,k'} \tag{145}$$

$$\begin{aligned}
\Theta_{k,k'}^l(x,x') &= \sum_{d=1}^{|W^l,b^l|} \frac{\partial z_k^l(x)}{\partial \{W^l,b^l\}} \frac{\partial z_{k'}^l(x')}{\partial \{W^l,b^l\}} + \sum_{d=1}^{|\theta^{l-1}|} \frac{\frac{\partial z_k^1(x)}{\partial \theta_d^{l-1}}}{\frac{\partial \theta_d^{l-1}}{\partial \theta_d^{l-1}}} \frac{\partial z_{k'}^l(x')}{\partial \theta_d^{l-1}} \\
&= K^l(x,x') \delta_{k,k'} + \delta_{k,k'} \dot{K}^l(x,x') \Theta_\infty^{l-1}(x,x') \\
&= \left(K^l(x,x') + \dot{K}^l(x,x') \Theta_\infty^{l-1}(x,x') \right) \delta_{k,k'} \\
&\text{repeating diagonal with } K^l(x,x') + \dot{K}^l(x,x') \Theta_\infty^{l-1}(x,x') \\
&= \underbrace{\left(K^l(x,x') + \dot{K}^l(x,x') \Theta_\infty^{l-1}(x,x') \right)}_{\text{scalar}} \otimes_{\text{outer}} \mathbf{I}_{N_{l+1} \times N_{l+1}}
\end{aligned} \tag{146}$$

10 NTK during training

Looking at training the Last-layer:

10.1 single data x under mean-square error

for a single data x in $\mathcal{R}^{d_{\text{in}}}$, and its associated label y , imagine last layer parameter is:

$$\theta^{L+1} = (W^{L+1}, b^{L+1}) \quad (147)$$

then, objective is:

$$\begin{aligned} C &= \frac{1}{2} \|f(x) - y\|_2^2 \\ &= \frac{1}{2} \left\| \left(\frac{\sigma_w}{\sqrt{N_t}} W^{L+1} \phi(z^L(x)) + \sigma_b b^{L+1} \right) - y \right\|_2^2 \end{aligned} \quad (148)$$

the above defines last layer as if it is the linear layers in NN, non-standard part is to re-parameterization $\frac{\sigma_w}{\sqrt{N_t}}$ and σ_b are added to allow:

$$W^{L+1} \sim \mathcal{N}(0, 1) \quad \text{and} \quad b^{L+1} \sim \mathcal{N}(0, 1) \quad (149)$$

then, the above is written as:

$$\begin{aligned} C &= \frac{1}{2} \left\| \underbrace{\begin{bmatrix} W^{L+1} & b^{L+1} \end{bmatrix}}_{\theta^{L+1}} \underbrace{\begin{bmatrix} \frac{\sigma_w}{\sqrt{N_t}} \phi(z^L(x)) & \sigma_b \end{bmatrix}^\top}_{\bar{a}(x)} - y \right\|_2^2 \\ &= \frac{1}{2} \left\| (\bar{a}(x)^\top \theta^{L+1}) - y \right\|_2^2 \\ &= \frac{1}{2} \left\| \hat{y}_t(x) - y \right\|_2^2 \\ \implies \frac{\partial C}{\partial \theta^{L+1}} &= \bar{a}(x)^\top (\bar{a}(x) \theta^{L+1} - y) \end{aligned} \quad (150)$$

reason to write this way is to express derivative in θ^{L+1} jointly, instead of writing W^{L+1} and b^{L+1} separately

also note that if only read-out layer is trained then $\bar{a}(x)$ is constant in t :

$$\begin{aligned} \frac{\partial \hat{y}_t(x)}{\partial \theta_t} &= \bar{a}(x) \\ \implies \frac{\partial \hat{y}_t(x)}{\partial \theta_t}^\top \frac{\partial \hat{y}_t(x)}{\partial \theta_t} &= \bar{a}(x)^\top \bar{a}(x) = \Theta_t(x, x) = \Theta_0(x, x) \end{aligned} \quad (151)$$

10.2 entire dataset \mathcal{X} :

10.2.1 softmax:

$$\hat{y}_t(\mathcal{X}) = \begin{bmatrix} \hat{y}_t^1(x_1) \\ \vdots \\ \hat{y}_t^{N^{L+1}}(x_1) \\ \vdots \\ \hat{y}_t^1(x_k) \\ \vdots \\ \hat{y}_t^{N^{L+1}}(x_k) \\ \vdots \\ \hat{y}_t^1(x_{|D|}) \\ \vdots \\ \hat{y}_t^{N^{L+1}}(x_{|D|}) \end{bmatrix} = \text{vec}([\hat{y}_t^i(x)]_{x \in \mathcal{X}}) \in \mathcal{R}^{N^{L+1} \times |D| \times 1} \quad (152)$$

10.2.2 mean-square error:

we focus on MSE:

$$\begin{aligned} C &= \frac{1}{2} \|(\bar{a}(\mathcal{X})^\top \theta^{L+1}) - \mathcal{Y}\|_2^2 \\ &= \frac{1}{2} \|\hat{y}_t(\mathcal{X}) - \mathcal{Y}\|_2^2 \\ \implies \frac{\partial C}{\partial \theta^{L+1}} &= \bar{a}(\mathcal{X})^\top (\bar{a}(\mathcal{X}) \theta^{L+1} - \mathcal{Y}) \end{aligned} \quad (153)$$

$$\begin{aligned}
\hat{y}_t(\mathcal{X}) &= \begin{bmatrix} \hat{y}_t(x_1) \\ \vdots \\ \hat{y}_t(x_k) \\ \vdots \\ \hat{y}_t(x_{|D|}) \end{bmatrix} = \text{vec}([\hat{y}_t(x)]_{x \in \mathcal{X}}) \in \mathcal{R}^{|D|} \times 1 \\
\Rightarrow \frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} &= \begin{bmatrix} \frac{d\hat{y}_t(x_1)}{d\theta_1} & \dots & \frac{d\hat{y}_t(x_1)}{d\theta_{|\theta|}} \\ \vdots & \ddots & \vdots \\ \frac{d\hat{y}_t(x_k)}{d\theta_1} & \dots & \frac{d\hat{y}_t(x_k)}{d\theta_{|\theta|}} \\ \vdots & \ddots & \vdots \\ \frac{d\hat{y}_t(x_{|D|})}{d\theta_1} & \dots & \frac{d\hat{y}_t(x_{|D|})}{d\theta_{|\theta|}} \end{bmatrix} \\
\Rightarrow \hat{\Theta}(\mathcal{X}, \mathcal{X}) &= \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta_i} \frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta_i}^\top \quad \text{empirical Tangent Kernel} \\
&= \begin{bmatrix} \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_1)}{\partial \theta_i} \frac{\partial \hat{y}_t(x_1)}{\partial \theta_i} & \dots & \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_1)}{\partial \theta_i} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_k)}{\partial \theta_i} \frac{\partial \hat{y}_t(x_k)}{\partial \theta_i} & \dots & \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_k)}{\partial \theta_i} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} & \dots & \sum_{i=1}^{|\theta|} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} \frac{\partial \hat{y}_t(x_{|D|})}{\partial \theta_i} \end{bmatrix}
\end{aligned} \tag{154}$$

10.3 what is linearized model: $f(x, \theta_t) \rightarrow f_t^{\text{lin}}(x, \theta) \equiv f^{\text{lin}}(x, \theta_t)$

$$\begin{aligned}
f_t^{\text{lin}}(x) &\equiv f_0(x) + \underbrace{\frac{\partial f(x, \theta)}{\partial \theta} \Big|_{\theta \rightarrow \theta_0}}_{\text{constant in } t} \underbrace{(\theta_t - \theta_0)}_{\Delta \theta(t)} \\
&= f_0(x) + \frac{\partial f(x, \theta)}{\partial \theta} \Big|_{\theta \rightarrow \theta_0} \Delta \theta(t) \quad \text{note } \Delta \theta(t) \text{ refer to change, irrespective of starting position } \theta_0
\end{aligned} \tag{155}$$

10.4 Sketch of Proof to compute $f_t^{\text{lin}}(x, \theta)$

We are interested to study the behavior of $f_t^{\text{lin}}(x)$ for a singular data x as θ_t evolves

the expression $f_t^{\text{lin}}(x)$ can be misleading: it should be written instead as $f_t^{\text{lin}}(x|\mathcal{X})$ as it depends on training dataset $(\mathcal{X}, \mathcal{Y})$, and θ as well.

Note that we can interchangeably write:

$$f_t^{\text{lin}}(\mathcal{X}, \theta) \equiv f_t^{\text{lin}}(\mathcal{X}) \equiv f^{\text{lin}}(\mathcal{X}, \theta_t) \tag{156}$$

Also, as we do not have expression for θ_t , we must start from $\frac{d\theta}{dt}$ obtained from gradient descent:

1. find expression for $f_t^{\text{lin}}(\mathcal{X}, \theta)$: it has two versions:
using gradient descent:

$$\frac{d\theta}{dt} = -\eta \left(\frac{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}{\partial \theta} \right)^\top \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)} \tag{157}$$

then,

version 1: assume $f_t^{\text{lin}}(\mathcal{X})$ is a linear model

- (a) from $\frac{d\theta}{dt}$, use ODE to obtain $\triangle\theta(t)$
- (b) then, obtain $f_t^{\text{lin}}(\mathcal{X}, \theta) = \left(\frac{\partial f_t^{\text{lin}}(\mathcal{X})}{\partial \theta}\right)^\top \triangle\theta(t)$

version 2:

- (a) obtain expression for:

$$\frac{d\hat{y}}{dt} = \left(\frac{\partial f_t^{\text{lin}}(\mathcal{X})}{\partial \theta}\right)^\top \frac{d\theta}{dt} \quad (158)$$

- (b) then, use ODE to obtain $f_t^{\text{lin}}(\mathcal{X}, \theta)$

2. now we have expression of $f_t^{\text{lin}}(\mathcal{X}, \theta)$, substitute it back to gradient descend:

$$\frac{d\triangle\theta(t)}{dt} = \frac{d\theta}{dt} = -\eta \left(\frac{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}{\partial \theta}\right)^\top \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)} \quad (159)$$

3. from $\frac{d\triangle\theta(t)}{dt}$, use ODE or straight integration to obtain:

$$\triangle\theta(t) = \theta_t - \theta_0 \quad (160)$$

4. Finally obtain how change of parameter contribute to last layer of single data in linear (in terms of θ_t) model:

$$f_t^{\text{lin}}(x, \theta_t) = f(x, \theta_0) + \frac{df(x, \theta_t)}{d\theta} \Big|_{\theta \rightarrow \theta_0} \triangle\theta(t) \quad (161)$$

think above as: instead of taking Euclidean step, i.e., $\theta_t = \theta_0 + h$, it now follows gradient descent path:

$$\theta_t = \theta_0 + \triangle\theta(t) \quad (162)$$

dependency chain is:

$$(\mathcal{X}, \mathcal{Y}) \rightarrow \frac{d\triangle\theta(t)}{dt} \rightarrow \triangle\theta(t) \rightarrow \hat{y}_t^{\text{lin}}(x, \theta) \quad (163)$$

note the following:

- We need all training data pairs $(\mathcal{X}, \mathcal{Y})$ to determine the change in θ
- so $\frac{d\theta}{dt}$ is a function of $(\mathcal{X}, \mathcal{Y})$, more specifically, $f(\mathcal{X}, \theta)$ and \mathcal{Y} , this is in Section (10.6)
- then we can work out how $\frac{d\theta}{dt}$ may impact $f_t^{\text{lin}}(x)$
- simply obtained expression of $f_t^{\text{lin}}(\mathcal{X})$ won't give you expression for $f_t^{\text{lin}}(x)$

10.5 General linear ODE solution

We need basic tools on ODE solution:

looking at the equation, treating everything in 1-d:

$$\begin{aligned}
 \dot{x} &= Ax + b \\
 \Rightarrow \frac{\dot{x}}{Ax + b} &= 1 \\
 \Rightarrow \frac{\dot{x}}{x + \frac{b}{A}} &= A \\
 \Rightarrow \frac{d}{dt} \log \left(x + \frac{b}{A} \right) &= A \quad \text{easy to see: } \frac{d}{dt} \log \left(x + \frac{b}{A} \right) = \frac{\frac{dx}{dt}}{x + \frac{b}{A}} \\
 \Rightarrow \int_t \frac{d}{dt} \log \left(x + \frac{b}{A} \right) dt &= \int_t A dt \\
 \log \left(x + \frac{b}{A} \right) &= At + h \\
 x + \frac{b}{A} &= \exp \left(At + h \right) \\
 \Rightarrow x(t) &= -\frac{b}{A} + \textcolor{red}{C} \exp \left(At \right)
 \end{aligned} \tag{164}$$

when things are in multi-dimensions:

$$x(t) = -A^{-1}b + C \exp \left(At \right) \quad b \text{ is column vector} \tag{165}$$

let $t = 0$:

$$\begin{aligned}
 x(0) &= -A^{-1}b + C \\
 \Rightarrow C &= (x(0) + A^{-1}b)
 \end{aligned} \tag{166}$$

substitute in C :

$$x(t) = -A^{-1}b + (x(0) + A^{-1}b) \exp \left(At \right) \tag{167}$$

10.6 find expression for $f_t^{\text{lin}}(\mathcal{X}, \theta)$:

10.6.1 generic MSE loss for f_t^{lin}

we use **version 2**:

$$\begin{aligned}
 \frac{d\theta(t)}{dt} &= -\eta \left(\frac{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}{\partial \theta} \right)^\top \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)} \\
 &= -\eta \left(\frac{\partial f_0(\mathcal{X}, \theta)}{\partial \theta} \right)^\top \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)} \quad \text{using property of linearised model Eq.(155)} \\
 \Rightarrow \frac{df_t^{\text{lin}}(\mathcal{X}, \theta)}{dt} &= \frac{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}{\partial \theta_t} \frac{d\theta(t)}{dt} \\
 &= -\eta \frac{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}{\partial \theta} \left(\frac{\partial f_0(\mathcal{X}, \theta)}{\partial \theta} \right)^\top \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)} \\
 &= -\eta \Theta_0(\mathcal{X}, \mathcal{X}) \frac{\partial C}{\partial f_t^{\text{lin}}(\mathcal{X}, \theta)}
 \end{aligned} \tag{168}$$

assume C is MSE:

$$\begin{aligned}\frac{df_t^{\text{lin}}(\mathcal{X}, \theta)}{dt} &= -\eta \Theta_0(\mathcal{X}, \mathcal{X}) (f_t^{\text{lin}}(\mathcal{X}) - \mathcal{Y}) \\ &= \underbrace{-\eta \Theta_0(\mathcal{X}, \mathcal{X})}_{A} f_t^{\text{lin}}(\mathcal{X}) + \underbrace{\eta \Theta_0(\mathcal{X}, \mathcal{X}) \mathcal{Y}}_b\end{aligned}\quad (169)$$

so by substitution, using Eq. (167), we have:

$$\begin{aligned}A &= -\eta \Theta_0(\mathcal{X}, \mathcal{X}) \\ \implies A^{-1} &= \frac{-1}{\eta} (\Theta_0(\mathcal{X}, \mathcal{X}))^{-1} \\ b &= \eta \Theta_0(\mathcal{X}, \mathcal{X}) \mathcal{Y}\end{aligned}\quad (170)$$

we have:

$$\begin{aligned}x(t) &= -A^{-1}b + (x(0) + A^{-1}b) \exp(At) \\ \implies \theta^{L+1}(t) &= -\frac{-1}{\eta} (\Theta_0(\mathcal{X}, \mathcal{X}))^{-1} (\eta \Theta_0(\mathcal{X}, \mathcal{X}) \mathcal{Y}) \\ &\quad + \left[f_0(\mathcal{X}) + \frac{-1}{\eta} (\Theta_0(\mathcal{X}, \mathcal{X}))^{-1} (\eta \Theta_0(\mathcal{X}, \mathcal{X}) \mathcal{Y}) \right] \exp(-\eta \Theta_0(\mathcal{X}, \mathcal{X}) t) \\ &= \mathcal{Y} + (f_0(\mathcal{X}) - \mathcal{Y}) \exp(-\eta \Theta_0(\mathcal{X}, \mathcal{X}) t)\end{aligned}\quad (171)$$

solution is:

$$f_t^{\text{lin}}(\mathcal{X}, \theta) = \mathcal{Y} + \exp(-\eta \Theta_0(\mathcal{X}, \mathcal{X}) t) (f_0(\mathcal{X}) - \mathcal{Y}) \quad (172)$$

10.6.2 read-out layer only under MSE

we are using **version 1** (link 1) to get expression for $\theta^{L+1}(t)$ from $\frac{d\theta^{L+1}(t)}{dt}$:
in MSE context:

$$\begin{aligned}C &= \frac{1}{2} \|(\bar{a}(\mathcal{X})^\top \theta^{L+1}) - \mathcal{Y}\|_2^2 \\ &= \frac{1}{2} \|\hat{y}_t(\mathcal{X}) - \mathcal{Y}\|_2^2\end{aligned}\quad (173)$$

$$\begin{aligned}\frac{d\theta^{L+1}(t)}{dt} &= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X}, \theta^{L+1})}{\partial \theta^{L+1}} \right)^\top \frac{\partial C}{\partial \hat{y}_t(\mathcal{X}, \theta^{L+1})} \\ &= -\eta \bar{a}(\mathcal{X})^\top (\bar{a}(\mathcal{X}) \theta^{L+1} - \mathcal{Y}) \\ &= \underbrace{-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X})}_A \theta^{L+1}(t) + \underbrace{\eta \bar{a}(\mathcal{X})^\top \mathcal{Y}}_b\end{aligned}\quad (174)$$

so by substitution, using Eq. (167), we have:

$$\begin{aligned}A &= -\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \\ \implies A^{-1} &= \frac{-1}{\eta} (\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}))^{-1} \\ b &= \eta \bar{a}(\mathcal{X})^\top \mathcal{Y}\end{aligned}\quad (175)$$

$$\begin{aligned}
x(t) &= -A^{-1}b + (x(0) + A^{-1}b) \exp(At) \\
\Rightarrow \theta^{L+1}(t) &= -\frac{1}{\eta} \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\eta \bar{a}(\mathcal{X})^\top \mathcal{Y} \right) \\
&\quad + \left[\theta^{L+1}(0) + \frac{1}{\eta} \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\eta \bar{a}(\mathcal{X})^\top \mathcal{Y} \right) \right] \exp \left(-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) t \right) \\
&= \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\bar{a}(\mathcal{X})^\top \mathcal{Y} \right) \\
&\quad + \left[\theta^{L+1}(0) - \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\bar{a}(\mathcal{X})^\top \mathcal{Y} \right) \right] \exp \left(-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) t \right)
\end{aligned} \tag{176}$$

look at MSE again: $\hat{y}_t(x) = \bar{a}(x)^\top \theta^{L+1}$:

$$\begin{aligned}
\hat{y}_t(\mathcal{X}) &= \theta^{L+1}(t) \bar{a}(\mathcal{X}) = \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\bar{a}(\mathcal{X})^\top \mathcal{Y} \bar{a}(\mathcal{X}) \right) \\
&\quad + \left[\theta^{L+1}(0) \bar{a}(\mathcal{X}) - \left(\bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) \right)^{-1} \left(\bar{a}(\mathcal{X})^\top \mathcal{Y} \bar{a}(\mathcal{X}) \right) \right] \exp \left(-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) t \right) \\
&= \mathcal{Y} + (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp \left(-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) t \right) \\
&= \mathcal{Y} + (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp \left(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X}) \right)
\end{aligned} \tag{177}$$

$\hat{y}_t(\mathcal{X})$ can be written in various forms:

$$\begin{aligned}
\hat{y}_t(\mathcal{X}) &= \underbrace{\mathcal{Y}} + \hat{y}_0(\mathcal{X}) \exp \left(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X}) t \right) - \underbrace{\mathcal{Y} \exp \left(-\eta \bar{a}(\mathcal{X})^\top \bar{a}(\mathcal{X}) t \right)} \\
&= (\mathbf{I} - \exp \left(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X}) t \right)) \mathcal{Y} + \hat{y}_0(\mathcal{X}) \exp \left(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X}) t \right)
\end{aligned} \tag{178}$$

relationship between $\hat{y}_t(\mathcal{X})$ and $\frac{d\hat{y}_t(\mathcal{X})}{dt}$:

$$\begin{aligned}
\hat{y}_t(\mathcal{X}) &= \mathcal{Y} + (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp \left(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X}) t \right) \\
\frac{d\hat{y}_t(\mathcal{X})}{dt} &= -\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X}) (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp \left(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X}) t \right) \\
&= -\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X}) (\hat{y}_t(\mathcal{X}) - \mathcal{Y})
\end{aligned} \tag{179}$$

no close-form solution when $\hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})$ is also a function of time, however since we are dealing with last layer, $\hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})$ is constant, therefore:

10.6.3 Two are equal!

looking at Eq. (177) and compare it with Eq. (172) side by side:

$$\begin{aligned}
\hat{y}_t(\mathcal{X}) &= \mathcal{Y} + (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp \left(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X}) \right) \\
f_t^{\text{lin}}(\mathcal{X}, \theta) &= \mathcal{Y} + \exp \left(-\eta \Theta_0(\mathcal{X}, \mathcal{X}) t \right) (f_0(\mathcal{X}) - \mathcal{Y})
\end{aligned} \tag{180}$$

it looks the same as long as $\mathcal{K}(\mathcal{X}, \mathcal{X}) \equiv \Theta_0(\mathcal{X}, \mathcal{X})$, but wait, if one look at Eq. (132)

$$\Theta_{k,k'}^1(x, x') \equiv K^1(x, x') \delta_{k,k'}$$

since we fixed all the preceeding layers, and only train readout layers, so it's like just K^1 using last-layer post-activation as oppose to input x and x'

this means all thereafter sections we can using either one in Eq.(180)

10.7 computing $\Delta\theta(t)$

using gradient descend, but substituting $\hat{y}(\mathcal{X})$:

$$\begin{aligned}
\frac{d\theta}{dt} &= -\eta \frac{\partial C}{\partial \theta} \\
&= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top \frac{\partial C}{\partial \hat{y}_t(\mathcal{X})} \\
&= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_t(\mathcal{X}) - \mathcal{Y}) \\
&= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \quad \text{using equation (179)}
\end{aligned} \tag{181}$$

to work out $\Delta\theta(t)$:

$$\begin{aligned}
\Delta\theta(t) &= \int_{\tau=0}^t \frac{d\theta}{d\tau} d\tau \\
&= \int_{\tau=0}^t -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})\tau) d\tau \\
&= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \underbrace{\int_{\tau=0}^t \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})\tau) d\tau}_{\text{}}
\end{aligned} \tag{182}$$

looking at:

$$\begin{aligned}
&\int_{\tau=0}^t \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})\tau) d\tau \\
&= \left[-\frac{1}{\eta} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})\tau) \right]_{\tau=0}^t \\
&= -\frac{1}{\eta} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) + \frac{1}{\eta} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \\
&= \underbrace{\frac{1}{\eta} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} (\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t))}_{\text{}}
\end{aligned} \tag{183}$$

Finally $\Delta\theta(t)$:

$$\begin{aligned}
\Delta\theta(t) &= -\eta \left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \frac{1}{\eta} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} (\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t)) \\
&= -\left(\frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} (\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t))
\end{aligned} \tag{184}$$

10.8 putting together for $\hat{y}(x, \theta_t)$

$$\begin{aligned}
\hat{y}(x, \theta_t) &= \hat{y}(x, \theta_0) + \frac{\partial \hat{y}(x, \theta)}{\partial \theta} \Big|_{\theta \rightarrow \theta_0} \Delta \theta(t) \\
&= \hat{y}(x, \theta_0) + \bar{a}(x) \left(- \frac{\partial \hat{y}_t(\mathcal{X})}{\partial \theta} \right)^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right) \\
&= \hat{y}(x, \theta_0) - \bar{a}(x) \bar{a}(\mathcal{X})^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right) \\
&= \hat{y}(x, \theta_0) - \bar{a}(x) \bar{a}(\mathcal{X})^\top (\hat{y}_0(\mathcal{X}) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right) \\
&= \hat{y}(x, \theta_0) - \hat{\mathcal{K}}(x, \mathcal{X}) (\hat{y}(\mathcal{X}, \theta_0) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right)
\end{aligned} \tag{185}$$

10.8.1 similarly for f_t^{lin}

looking at pattern, you will see that:

$$f_t^{\text{lin}}(x, \theta) = f(x, \theta_0) - \Theta_0(x, \mathcal{X}) (f(\mathcal{X}, \theta_0) - \mathcal{Y}) \Theta_0(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \Theta_0(\mathcal{X}, \mathcal{X})t) \right) \tag{186}$$

10.9 expectation and variance

mean

$$\begin{aligned}
\mathbb{E}[\hat{y}(x, \theta_t)] &= \mathbb{E}[\hat{y}(x, \theta_0) - \hat{\mathcal{K}}(x, \mathcal{X}) (\hat{y}(\mathcal{X}, \theta_0) - \mathcal{Y}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right)] \\
&= \underbrace{\mathbb{E}[\hat{y}(x, \theta_0)]}_{=0} - \underbrace{\mathbb{E}[\hat{y}(\mathcal{X}, \theta_0)]}_{=0} \mathbb{E}[\hat{\mathcal{K}}(x, \mathcal{X}) \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right)] \\
&\quad + \mathbb{E}[\hat{\mathcal{K}}(x, \mathcal{X}) \mathcal{Y} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right)] \\
&= \mathbb{E}[\hat{\mathcal{K}}(x, \mathcal{X}) \mathcal{Y} \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \hat{\mathcal{K}}(\mathcal{X}, \mathcal{X})t) \right)] \\
&= \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{Y} \quad \text{deterministic in infinite width}
\end{aligned} \tag{187}$$

variance

$$\begin{aligned}
&\hat{y}(x, \theta_t) - \mathbb{E}[\hat{y}(x, \theta_t)] \quad \text{let infinite width} \\
&= \hat{y}(x, \theta_0) - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) (\hat{y}(\mathcal{X}, \theta_0) - \mathcal{Y}) \\
&\quad - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{Y} \\
&= \hat{y}(x, \theta_0) - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0) \\
&\quad + \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{Y} \\
&\quad - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{Y} \\
&= \hat{y}(x, \theta_0) - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0)
\end{aligned} \tag{188}$$

then:

$$\begin{aligned}
& \text{Var}[\hat{y}(x, \theta_t)] \\
&= \mathbb{E} \left[\left(\hat{y}(x, \theta_t) - \mathbb{E}[\hat{y}(x, \theta_t)] \right)^\top \left(\hat{y}(x, \theta_t) - \mathbb{E}[\hat{y}(x, \theta_t)] \right) \right] \\
&= \mathbb{E} \left[\left(\hat{y}(x, \theta_0) - \underbrace{\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0)}_{\text{symmetric}} \right)^\top \right. \\
&\quad \left. \left(\hat{y}(x, \theta_0) - \underbrace{\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0)}_{\text{symmetric}} \right) \right] \quad (189)
\end{aligned}$$

knowing that when $t = 0$:

$$\begin{aligned}
& \text{Cov}[\hat{y}(x, \theta_0), \hat{y}(\mathcal{X}, \theta_0)] = \mathbb{E}[\hat{y}(x, \theta_0) \hat{y}(\mathcal{X}, \theta_0)^\top] = \mathcal{K}(x, \mathcal{X}) \\
& \text{Cov}[\hat{y}(\mathcal{X}, \theta_0), \hat{y}(x, \theta_0)] = \mathbb{E}[\hat{y}(\mathcal{X}, \theta_0) \hat{y}(x, \theta_0)^\top] = \mathcal{K}(\mathcal{X}, x) \\
& \text{Var}[\hat{y}(x, \theta_0), \hat{y}(x, \theta_0)] = \mathbb{E}[\hat{y}(x, \theta_0) \hat{y}(x, \theta_0)^\top] = \mathcal{K}(x, x) \\
& \quad \left(\underbrace{\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0)}_{\text{symmetric}} \right)^\top \quad (190) \\
& \quad = \left(\hat{y}(\mathcal{X}, \theta_0)^\top \underbrace{\mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x)}_{\text{symmetric}} \right)
\end{aligned}$$

$$\begin{aligned}
& \text{Var}[\hat{y}(x, \theta_t)] = \mathbb{E}[\hat{y}(x, \theta_0) \hat{y}(x, \theta_0)^\top] \\
& \quad - \mathbb{E}[\hat{y}(x, \theta_0) \hat{y}(\mathcal{X}, \theta_0)^\top \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x)] \\
& \quad - \mathbb{E}[\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0) \hat{y}(x, \theta_0)^\top] \\
& \quad + \mathbb{E}[\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \hat{y}(\mathcal{X}, \theta_0) \hat{y}(\mathcal{X}, \theta_0)^\top \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x)] \\
& = \mathcal{K}(x, x) \\
& \quad - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \\
& \quad - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \\
& \quad + \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \\
& = \mathcal{K}(x, x) - 2\mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \\
& \quad + \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \left(\mathbf{I} - \exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \\
& = \mathcal{K}(x, x) - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\underbrace{2\mathbf{I} - 2\exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t)}_{\text{red}} \right) \mathcal{K}(\mathcal{X}, x) \\
& \quad + \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\underbrace{\mathbf{I} - 2\exp(-\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t)}_{\text{red}} + \exp(-2\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \quad (191)
\end{aligned}$$

terms outside of the red bits are the same:

$$\text{Var}[\hat{y}(x, \theta_t)] = \mathcal{K}(x, x) - \mathcal{K}(x, \mathcal{X}) \mathcal{K}(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-2\eta \mathcal{K}(\mathcal{X}, \mathcal{X})t) \right) \mathcal{K}(\mathcal{X}, x) \quad (192)$$

10.9.1 similar in f_t^{lin}

- mean:

$$\mathbb{E}[f_t^{\text{lin}}(x, \theta_t)] = \Theta(x, \mathcal{X}) \Theta_0(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-\eta \Theta_0(\mathcal{X}, \mathcal{X}) t) \right) \mathcal{Y} \quad (193)$$

- variance:

$$\text{Var}[f_t^{\text{lin}}(x, \theta_t)] = \Theta_0(x, x) - \Theta_0(x, \mathcal{X}) \Theta_0(\mathcal{X}, \mathcal{X})^{-1} \left(\mathbf{I} - \exp(-2\eta \Theta_0(\mathcal{X}, \mathcal{X}) t) \right) \Theta_0(\mathcal{X}, x) \quad (194)$$

11 Infinite width networks are linearized networks

for every $x \in \mathcal{R}^{N_0}$ with $\|x\|_2 \leq 1$ with probability close to 1 over random initialization:

$$\left. \begin{aligned} & \sup_{t \geq 0} \|f_t(x) - f_t^{\text{lin}}\|_2 \\ & \sup_{t \geq 0} \frac{\|\theta_t - \theta_0\|_2}{\sqrt{n}} \\ & \sup_{t \geq 0} \|\hat{\Theta}_t - \hat{\Theta}_0\|_F \end{aligned} \right\} = \mathcal{O}(n^{-\frac{1}{2}}) \quad \text{as } n \rightarrow \infty \quad (195)$$

check relevant publications for the proof