

Neural Networks as a Gaussian Process

Richard Xu

August 15, 2020

1 Preamble

This talk is heavily referenced to the following:

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

I tries to unify notations of the above references

2 Gaussian Process

- if one is to perform a predictive distribution $p(y^*|y, X, x^*)$ through GP:

$$\begin{aligned} p\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix}\right) &= \int p\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} X \\ x^{*\top} \end{bmatrix}, f\right) p(f|X) \mathrm{d}f \\ &= \int \mathcal{N}\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f(X) \\ f(x^{*\top}) \end{bmatrix}, \sigma_\epsilon^2 I\right) p(f|X, x^*) \mathrm{d}f \end{aligned}$$

- This is the **key**: prior $p(f|X, x^*)$ is defined over function $f(X)$ instead of X
- Imagine, if instead, prior is defined over X , i.e., $p(X)$ is the prior:

$$\int \mathcal{N}\left(\begin{bmatrix} y \\ y^* \end{bmatrix} \middle| \begin{bmatrix} f(X) \\ f(x^{*\top}) \end{bmatrix}, \sigma_\epsilon^2 I\right) p(X) \mathrm{d}X$$

Then, non-linear f is **not** making integral tractable!

3 GP for Neural Network: Directly computation

3.1 neural network function

using parameters:

$$\omega \equiv \{W^L, b^L, \dots, W^1, b^1\}$$

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

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

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

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

3.2 Apply NN function in predictive distribution

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

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

- The integral is **not** analytic!!

3.3 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:

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

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

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

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

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

3.3.1 look at the mean:

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

note we are not dealing with infinity at the moment

3.3.2 look at co-variance

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

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

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

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

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

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

$$= \sigma_\epsilon^2 I + \int_\omega \begin{bmatrix} \phi^L(X)W^LW^{L\top}x^L(X)^\top + b^Lb^{L\top} & \phi^L(X)W^LW^{L\top}x^L(x^{\star\top})^\top + b^Lb^{L\top} \\ \phi^L(x^{\star\top})W^LW^{L\top}\phi^L(X)^\top + b^Lb^{L\top} & \phi^L(x^{\star\top})W^LW^{L\top}\phi^L(x^{\star\top})^\top + b^Lb^{L\top} \end{bmatrix} \mathcal{N}(\omega | 0, \sigma_\omega^2 I) d\omega$$

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

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

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

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

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

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

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

4 Single layer neural network

$$f_k(x) = b_k + \sum_{j=1}^H v_{jk} h_j(x)$$

$$h_j(x) = \tanh \left(a_j + \sum_{i=1}^I u_{ij} x_i \right)$$

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

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

4.1 $p(z_k^l(x))$ for single input x

We need CLT for computing this probability.

4.1.1 Central Limit Theorem:

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

- 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:

$$\sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

$$(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, \frac{\sigma^2}{n})$$

$$\frac{1}{\sigma} \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, 1)$$

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

$$\sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

$$\implies \sqrt{n} \sqrt{n}(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, \sqrt{n}^2 \sigma^2) = \mathcal{N}(0, n\sigma^2)$$

$$\implies n(\bar{X} - \mathbb{E}[X^{(1)}]) \xrightarrow{d} \mathcal{N}(0, n\sigma^2)$$

$$\implies \left(\sum_{i=1}^n X_i - n\mathbb{E}[X^{(1)}] \right) \xrightarrow{d} \mathcal{N}(0, n\sigma^2)$$

choose one of these conditions to suit the situation

4.1.2 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 $b_k^l +$ sum of i.i.d elements using CLT notations:

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

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

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

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

- **Expectation**

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

- **Variance**

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

we leave in this form, as

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

4.1.4 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.

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

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

$$\sigma_w = \frac{C_w}{\sqrt{N_l}} \quad \Rightarrow \quad \sigma_w^2 = \frac{C_w^2}{N_l}$$

- the above is the key, implication is:

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

- finally adding the bias b_k^l :

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

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

Therefore:

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

- appreciate the recursion here

4.2 given two inputs $x^{(p)}$, $x^{(q)}$: compute $\text{Cov}[z_k^l(x^{(p)}) z_k^l(x^{(q)})]$

To do so, we need to use **Multidimensional CLT**

4.2.1 Multidimensional CLT:

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

Therefore:

$$\begin{aligned}
&\frac{1}{\sqrt{n}} \sum_{i=1}^n [\mathbf{X}_i - \mathbb{E}[\mathbf{X}_i]] \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n (\mathbf{X}_i - \mathbb{E}[\mathbf{X}_1]) = \frac{\sqrt{n}}{\sqrt{n}} \frac{1}{\sqrt{n}} \left(\sum_{i=1}^n \mathbf{X}_i \right) - n \mathbb{E}[\mathbf{X}_1] \\
&= \sqrt{n} (\bar{\mathbf{X}} - \mathbb{E}[\mathbf{X}_1])
\end{aligned}$$

- Sample mean version:

$$\Rightarrow \sqrt{n} \mathbb{E} \left[\underbrace{\left(\bar{\mathbf{X}}^{(p)} - \mathbb{E}[\bar{\mathbf{X}}_1^{(p)}] \right)}_{\text{scalar}} \underbrace{\left(\bar{\mathbf{X}}^{(q)} - \mathbb{E}[\bar{\mathbf{X}}_1^{(q)}] \right)}_{\text{scalar}} \right] = \Sigma_{(p),(q)}$$

for each co-variance/non-diagonal elements $(p, q) \in \{1, \dots, k\}$:

- Sample sum version:

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

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

4.2.2 put in Multidimensional CLT structure:

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

Since we already know that:

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

for any arbitrary $j = 1$, and then:

$$\begin{aligned}
& \mathbb{E} \left[\left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x^{(p)})) \right) \left(\sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x^{(q)})) \right) \right] \\
&= N_l \Sigma_{(p),(q)} \\
&= N_l \text{Cov} \left(W_{k,1}^l \phi(z_1^{l-1}(x^{(p)})), W_{k,1}^l \phi(z_1^{l-1}(x^{(q)})) \right) \\
&= N_l \mathbb{E} \left[W_{k,1}^l \phi(z_1^{l-1}(x^{(p)})) \times W_{k,1}^l \phi(z_1^{l-1}(x^{(q)})) \right]
\end{aligned}$$

add b_k^l into, and look at $z_k^l(x)$:

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

- **note 1:** 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:

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

- **note 2:** in literature, it is written:

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

This is because of **note1** above

- regardless of this special property CLT still apply.

4.2.3 Relationship with Gaussian Process (GP):

let $f_k(x) \equiv z_k^l(x)$ be some function, and since for every arbitrary point pair, $x^{(p)}$ and $x^{(q)}$, we have:

$$\begin{aligned}\mathbb{E}[f(x)] &= 0 \\ \mathbb{E}[f(x^{(p)}), f(x^{(q)})] &= \mathbf{\Sigma}_{(p),(q)} \\ \implies f &\sim \mathcal{GP}(0, \mathbf{\Sigma})\end{aligned}$$

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

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

- putting it in layer specific GP:

$$\begin{aligned}\implies z_k^l(x) &\sim \mathcal{GP}(0, \mathbf{\Sigma}) \\ \text{where } \mathbf{\Sigma}_{p,q} &= \sigma_b^2 + C_w^2 \mathbb{E}[\phi(z_1^{l-1}(x^{(p)})) \times \phi(z_1^{l-1}(x^{(q)}))] \quad \text{as } N_l \rightarrow \infty\end{aligned}$$

4.3 more on GP

first define $K^l(x^{(p)}, x^{(q)})$ in terms of pre-activation $z_k^l(x)$ in this section, it will be changed later to post-activation

$$\begin{aligned}K^l(x^{(p)}, x^{(q)}) &= \mathbb{E}[z_k^l(x^{(p)})z_k^l(x^{(q)}) | z^{l-1}] \\ &= \mathbb{E}\left[\left(b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x^{(p)}))\right) \times \left(b_k^l + \sum_{j=1}^{N_l} W_{k,j}^l \phi(z_j^{l-1}(x^{(q)}))\right)\right] \\ &= \sigma_b^2 + \sigma_w^2 \mathbb{E}\left[\sum_{j=1}^{N_l} \phi(z_j^{l-1}(x^{(p)})) \times \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x^{(q)}))\right] \\ &= \sigma_b^2 + \sigma_w^2 \mathbb{E}[\phi(z_1^{l-1}(x^{(p)})) \times \phi(z_1^{l-1}(x^{(q)}))] \quad \text{apply CLT } N_l \rightarrow \infty \\ &= \sigma_b^2 + \sigma_w^2 \underbrace{\mathbb{E}_{z_1^{l-1} \sim \mathcal{GP}(0, K^{l-1})} [\phi(z_1^{l-1}(x^{(p)})) \phi(z_1^{l-1}(x^{(q)}))]}_{\text{since } \mathbb{E}[\phi(z)] = \mathbb{E}_{z \sim p(z)} [\phi(z)]} \\ &= \sigma_b^2 + \sigma_w^2 \underbrace{F_\phi(K^{l-1}(x^{(p)}), x^{(q)}), K^{l-1}(x^{(p)}), x^{(p)}), K^{l-1}(x^{(q)}), x^{(q)})}_{G(F_\phi(K^{l-1}))} \\ &= G \circ F_\phi(K^{l-1}(x^{(p)}), x^{(q)})\end{aligned}$$

using properties of point Marginals of Gaussian Process:

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

$$\begin{bmatrix} z_j^{l-1}(x^{(\textcolor{red}{p})}) \\ z_j^{l-1}(x^{(\textcolor{blue}{q})}) \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K^{l-1}(x^{(p)}, x^{(p)}) & K^{l-1}(x^{(p)}, x^{(q)}) \\ K^{l-1}(x^{(p)}, x^{(q)}) & K^{l-1}(x^{(q)}, x^{(q)}) \end{bmatrix}\right)$$

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

$$= F_\phi(K^{l-1}(x^{(p)}, x^{(q)}), K^{l-1}(x^{(p)}, x^{(p)}), K^{l-1}(x^{(q)}, x^{(q)}))$$

5 Expand GP across all layers

5.1 Overall objective

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

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

5.2 $p(z^L|K^L)$: conditions on $K^l \equiv \{\phi(z^{l-1})(x^{(p)}))\phi(z^{l-1})(x^{(q)})\}_{p,q}$

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

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

instead it define K^l in terms of post-activation of previous layer $\phi(z^{l-1})$ for reason illustrated later

- look at Neural Network function:

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

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

- Conditional Marginal

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

using property of weighted sum of Gaussian:

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

- Conditional Co-variance

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

not using property of weighted sum of Gaussian:

- Combine all together

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

- in GP paradigm:

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

where

$$K^l(x^{(p)}, x^{(q)}) = \frac{1}{N_l} \sum_{j=1}^{N_l} \phi(z_j^{l-1}(x^{(p)})) \phi(z_j^{l-1}(x^{(q)}))$$

$$G(K^l(x^{(p)}, x^{(q)})) = \sigma_b^2 + \sigma_w^2 K^l(x^{(p)}, x^{(q)})$$

Conveniently, we use K^l as a short-notation collection of $\phi(z_j^{l-1}(x^{(p)}))$, $\phi(z_j^{l-1}(x^{(q)})) \quad \forall p, q, j$

- also taking care of the layer one, which is just input x :

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

- to reflect:

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

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

$$\mathbf{z}^l = \underbrace{\begin{bmatrix} \color{red}{z_1^l(x^{(1)})} & \color{red}{z_1^l(x^{(2)})} & \dots & z_1^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_j^l(x^{(1)}) & z_j^l(x^{(2)}) & \dots & z_j^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_{N_{l+1}}^l(x^{(1)}) & z_{N_{l+1}}^l(x^{(2)}) & \dots & z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix}}_{|\mathcal{D}|} \left. \vphantom{\begin{bmatrix} \color{red}{z_1^l(x^{(1)})} & \color{red}{z_1^l(x^{(2)})} & \dots & z_1^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_j^l(x^{(1)}) & z_j^l(x^{(2)}) & \dots & z_j^l(x^{(|\mathcal{D}|)}) \\ \vdots & \vdots & \ddots & \vdots \\ z_{N_{l+1}}^l(x^{(1)}) & z_{N_{l+1}}^l(x^{(2)}) & \dots & z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix}} \right\} \text{width} \Rightarrow \text{vec}(\mathbf{z}^l) = \begin{bmatrix} \color{red}{z_1^l(x^{(1)})} \\ z_2^l(x^{(1)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(1)}) \\ \color{red}{z_1^l(x^{(2)})} \\ z_2^l(x^{(2)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(2)}) \\ \vdots \\ \vdots \\ z_1^l(x^{(|\mathcal{D}|)}) \\ z_2^l(x^{(|\mathcal{D}|)}) \\ \vdots \\ z_{N_{l+1}}^l(x^{(|\mathcal{D}|)}) \end{bmatrix}$$

$$\sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} G(K_{1,1}^l) & \dots & 0 & \dots & G(K_{1,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{1,1}^l) & \dots & 0 & 0 & G(K_{1,|\mathcal{D}|}^l) \\ \color{red}{G(K_{2,1}^l)} & \dots & 0 & \dots & G(K_{2,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{2,1}^l) & \dots & 0 & 0 & G(K_{2,|\mathcal{D}|}^l) \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ G(K_{|\mathcal{D}|,1}^l) & \dots & 0 & \dots & G(K_{|\mathcal{D}|,|\mathcal{D}|}^l) & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & G(K_{|\mathcal{D}|,1}^l) & \dots & 0 & 0 & G(K_{|\mathcal{D}|,|\mathcal{D}|}^l) \end{bmatrix} \right)$$

$$\Rightarrow p(\mathbf{z}^l | K^l) = \mathcal{N}(\mathbf{0}, G(K^l) \otimes \mathbf{I}_{N_{l+1} \times N_{l+1}})$$

$$= \mathcal{GP}(\mathbf{z}^l; \mathbf{0}, G(K^l))$$

5.3 $p(K^l | K^{l-1})$

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

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

- using GP property, and just look at two points $x^{(p)}, x^{(q)}$:

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

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

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

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

- using CLT, pick the most appropriate definition:

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

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

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

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

5.4 putting in the overall objective function

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

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