

Correlated Weights in Infinite Limits of Deep Convolutional Neural Networks (Supplementary material)

Adrià Garriga-Alonso¹

Mark van der Wilk²

¹Department of Engineering, University of Cambridge, UK

²Department of Computer Science, Imperial College London, UK

A PATCH FUNCTIONS AND DISCRETE CONVOLUTIONS

Usually, convolutions are defined explicitly by subtracting the indices of one input tensor from the other one, and not using patch functions. To make this paper clearer, it is convenient to abstract the details of a convolution, so we introduced the patch function.

Definition A.1 (Discrete convolution). Let $D \in \mathbb{N}$ be a number of spatial dimensions, the tensor-valued weights $\mathbf{W} \in \mathbb{R}^P$, input $\mathbf{X} \in \mathbb{R}^F$, and output $\mathbf{Y} \in \mathbb{R}^{F'}$. The tensor sizes P (patch size) and F, F' (feature sizes) are each a D -tuple, $P, F, F' \in \mathbb{N}^D$. We say that \mathbf{Y} is the result of the convolution operation $\mathbf{Y} = \mathbf{W} * \mathbf{X}$, if

$$Y_{q_1, \dots, q_D} = \sum_{p_1=1}^{P_1} \cdots \sum_{p_D=1}^{P_D} W_{p_1, \dots, p_D} X_{\tilde{q}_1(p_1), \dots, \tilde{q}_D(p_D)}. \quad (1)$$

Here, $\tilde{q}_d(\cdot) : [P_d] \rightarrow [F'_d]$ are the *patch functions* for a given output location \mathbf{q} . Using D -tuples \mathbf{p}, \mathbf{q} as indices, we may also write

$$Y_{\mathbf{q}} = \sum_{\mathbf{p}=1}^P W_{\mathbf{p}} X_{\tilde{\mathbf{q}}(\mathbf{p})}. \quad (2)$$

Counting from $\mathbf{1}$ to \mathbf{P} is done in such a way that \mathbf{p} takes all the values in $[\mathbf{P}]$.

Definition A.2 (Patch function). For each dimension $d \in [D]$, layer $\ell \in [L]$, fix a stride $s \in \mathbb{N}$, and dilation $h \in \mathbb{N}$. For output position \mathbf{q} , the patch function of the d th dimension $\tilde{q}_d(\cdot) : [P_d] \rightarrow [F'_d]$ is

$$\tilde{q}_d(p_d) = sq_d - h \left(p_d - \left\lceil \frac{P_d}{2} \right\rceil \right). \quad (3)$$

For a D -tuple index \mathbf{p} , we may compactly write $\tilde{\mathbf{q}}(\mathbf{p}) \triangleq (\tilde{q}_1(p_1), \dots, \tilde{q}_D(p_D))$.

It is possible to verify that definition A.2 overall yields the usual definition of a convolution in deep learning [Goodfellow et al., 2016, Section 9.1].

Remark A.3. The concatenation of two patch functions $\tilde{q}(\cdot), \tilde{q}'(\cdot)$ is also a patch function, with argument in $[P]^2$. That is, for $[\mathbf{p}, \mathbf{p}'] = \mathbf{s} \in \mathbb{N}^{P \times P}$ and $[\mathbf{q}, \mathbf{q}'] = \mathbf{r}$,

$$\begin{aligned} (\tilde{\mathbf{q}}(\mathbf{p}), \tilde{\mathbf{q}}'(\mathbf{p}')) &= (\tilde{q}_1(p_1), \dots, \tilde{q}_D(p_D), \tilde{q}'_1(p'_1), \dots, \tilde{q}'_D(p'_D)) \\ &= \tilde{\mathbf{r}}(\mathbf{s}). \end{aligned} \quad (4)$$

B PROOF THAT A CNN WITH CORRELATIONS IN THE WEIGHTS CONVERGES TO A GP

In this section, we formally prove that a CNN with correlated weights converges in distribution to a Gaussian process in the limit of infinite width. Using the NETSOR programming language due to Yang [2019], most of the work in the proof is done by one step: describe a CNN with correlated weights in NETSOR.

For the reader's convenience, we informally recall the NETSOR programming language [Yang, 2019] and key properties of its programs (theorem B.7 and corollary B.8). The outline of our presentation here also closely follows Yang [2019]. Readers familiar with NETSOR should skip to appendix B.3, where we show the program that proves theorem B.10.

B.1 DEFINITION OF A NETSOR PROGRAM

A NETSOR program expresses numerical computations, such as those used to define the output of a neural network. Each line of a NETSOR program is simply the definition of a new variable, in terms of previously defined variables.

There are three types of variables: $G(n)$ -vars, $A(n_1, n_2)$ -vars, and $H(n)$ -vars (henceforth called ‘‘NETSOR variables’’). Each of these have one or two parameters, which are the widths we will take to infinity. For a given index in $[n]$ (or $[n_1] \times [n_2]$), each NETSOR variable is a random

scalar. To represent vectors that do not grow to infinity, we need to use collections of NETSOR variables.

G-vars, A-vars and H-vars are all random when the program is run. To accomodate non-random variables that may change (like the input \mathbf{X} to a neural network (NN)) we must define a different NETSOR program, defining \mathbf{X} either as a constant or a G-var with variance zero.

What follows is an explanation of the three kinds of NETSOR variables, and example uses of them. The program indicates the type of a variable using “var : Type”.

G-vars (Gaussian-vars) are n -wise *approximately* independent and identically distributed (i.i.d.) and Gaussian. By “ n -wise (approximately) independent” we mean that there can be correlations between G-vars, but only within a single index $i \in 1, \dots, n$. G-vars will converge in distribution to an n -wise independent, identically distributed Gaussian in the limit of $n \rightarrow \infty$, if all widths are n . They are used, for example, to define the biases of a fully connected neural network (FCNN).

A-vars represent matrices, like the weight matrices of a dense neural network. Their entries are always i.i.d. Gaussian with zero mean, even for finite instantiations of the program (finite n). There are no correlations between different A-vars, or elements of the same A-var. They may be used to define the weight matrices of a FCNN.

H-vars represent variables that become n -wise i.i.d. (not necessarily Gaussian) in the infinite limit. G is a subtype of H, so all G-vars are also H-vars. Post-nonlinearity activations are H-vars.

O-vars (Output-vars) are used to define the output of the NETSOR program. A $O(n)$ -var behaves like you would expect a hypothetical $A(n, 1)$ -var to behave: its elements are i.i.d. Gaussian with mean zero, and it is independent of all other variables in the program.

Yang [2019] does not define O-vars, instead choosing to consider them part of the G-vars, since they both converge to i.i.d. Gaussians.

Definition B.1 (Netsor program). A NETSOR program consists of:

Input: A set that may contain G-vars, A-vars, and O-vars.

Body: Each line of the program defines a new variable in terms of existing ones. New variables may be defined using the following rules:

- **MatMul:** $A(n_1, n_2) \times H(n_2) \rightarrow G(n_1)$. Given an $A(n_1, n_2)$ -var (i.i.d. Gaussian matrix) and an $H(n_2)$ -var (i.i.d. vector), its multiplication is a $G(n_1)$ -var (that is, it converges to a Gaussian vector in the limit $n_2 \rightarrow \infty$).
- **LinComb:** Given constants $\alpha_1, \dots, \alpha_K$, and G-vars

x_1, \dots, x_K of type $G(n_1)$, their linear combination $\sum_{k=1}^K \alpha_k x_k$ is a G-var.

- **Nonlin:** applying an elementwise nonlinear function $\phi : \mathbb{R}^K \rightarrow \mathbb{R}$, we map several G-vars x_1, \dots, x_K to one H-var.

Output: A tuple of scalars $(o_1^\top x_1, \dots, o_K^\top x_K / \sqrt{n_K})$. The variables $o_k : O(n_k)$ are O-vars. It may be the case that $o_j = o_k$ for different j, k (that is, the list $[v_1, \dots, v_K]$ has repeated entries). Each $x_k : H(n_k)$ is a H-var.

Outside of these rules, NETSOR does not have conditionals or loops¹. In practice, we may use loops and conditionals to write a NETSOR program, so long as they do not access the values of NETSOR variables. Conceptually, these behave like a LISP-style “macro” that generates a NETSOR program.

B.2 THE OUTPUT OF A NETSOR PROGRAM CONVERGES TO A GAUSSIAN PROCESS

For simplicity, we assume that the width of all the NETSOR variables is n . Yang [2019] also considers the case where each n_k is different. First, the necessary assumptions.

Definition B.2 (Controlled function [Yang, 2019]). A function $\phi : \mathbb{R}^k \rightarrow \mathbb{R}$ is *controlled* if it is measurable and

$$|\phi(\mathbf{x})| \leq \exp\left(C\|\mathbf{x}\|_2^{(2-\epsilon)}\right) + c$$

for some $C, c, \epsilon > 0$, where $\|\cdot\|_2$ is the L2 norm.

If a function is controlled, it is L2 integrable with a Gaussian. That is, if the argument x of the function is Gaussian, the variance of $\phi(x)$ is finite. This in turn ensures that the NN function has finite variance. All common nonlinearities (ReLU, tanh, SiLU, ...) are controlled. This is a very weak assumption, it is vanishingly unlikely that future nonlinearities will grow as fast as $O(e^{x^2})$.

Assumption B.3. All nonlinear functions $\phi(\cdot)$ in the NETSOR program are controlled.

Assumption B.4 (Distribution of A-var inputs). Consider each $A(n, n)$ -var in the program, \mathbf{W} . Each of its elements $W_{i,j}$, where $i, j \in [n]$, is sampled from the zero-mean, i.i.d. Gaussian, $W_{i,j} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_w^2/n)$.

Assumption B.5 (Distribution of G-var inputs). Consider the input vector of all $G(n)$ -vars for each channel $i \in [n]$, that is the vector $\mathbf{z}_i \triangleq [x_i : x \text{ is input G-var}]$. It is drawn from a Gaussian, $\mathbf{z}_i \stackrel{iid}{\sim} \mathcal{N}(\boldsymbol{\mu}^{\text{in}}, \boldsymbol{\Sigma}^{\text{in}})$. The covariance $\boldsymbol{\Sigma}^{\text{in}}$ may be singular.

¹Of course, a nonlinearity ϕ may be internally defined using loops and conditionals, so long as it satisfies assumption B.3.

Assumption B.6 (Distribution of O-vars). Each $O(n_k)$ -var v_k in the program is an independent Gaussian for each channel. Different O-vars may have different variances. That is, for each $k \in [K], i \in [n], v_{k,i} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_k^2/n_k)$.

Theorem B.7 (NETSOR master theorem, Yang, 2019). Fix any NETSOR program satisfying assumptions B.3 to B.6. If $g^{(1)}, \dots, g^{(M)}$ are all the G-vars in the entire program, then for any controlled $\psi : \mathbb{R}^M \rightarrow \mathbb{R}$, as $n \rightarrow \infty$,

$$\frac{1}{n} \sum_{i=1}^n \psi(g_i^{(1)}, \dots, g_i^{(M)}) \xrightarrow{a.s.} \mathbb{E}_{z \sim \mathcal{N}(\mathbf{m}, \mathbf{K})} [\psi(z^{(1)}, \dots, z^{(M)})]. \quad (5)$$

Here $\xrightarrow{a.s.}$ is almost sure convergence [Rosenthal, 2006, sec. 5.2]. The mean \mathbf{m} and covariance \mathbf{K} are calculated under the assumption that all the G-vars are jointly Gaussian, like in section 3.

Proof sketch. The proof is by induction on the number of G-vars included in the output, added in order of definition. The induction invariant is that, for some $m < M$, eq. 5 holds; and that a subset of G-vars in $[m]$ which form a basis have a non-singular distribution. The detailed proof is in Yang [2019, Appendix H]. \square

The following corollary is a consequence of the Master theorem (B.7) and the Central Limit Theorem.

Corollary B.8 (Corollary 5.5, abridged, Yang, 2019). Fix any NETSOR program which satisfies assumptions B.3 to B.6. For simplicity, fix the widths of all the variables to n . The program outputs are $(o_1^\top x_1, \dots, v_K^\top x_K)$, where each x_k is an H-var, and each o_k is a O-var. Then, as $n \rightarrow \infty$, the output tuple converges in distribution to a Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{K})$. The covariance \mathbf{K} is given by doing calculations like section 3, assuming that G-vars are jointly Gaussian.

B.3 NETSOR PROGRAM AND GP BEHAVIOUR: CNN WITH CORRELATED WEIGHTS

NETSOR only has native support for matrix-vector multiplications and linear combinations with constants. How can we represent a convolution operation for convolutional neural network (CNN)? Consider the convolutional layer definition (eq. 1). Changing the sum order, we obtain

Expanding the convolution into a sum, and changing the sum order, we obtain

$$Z_{i,q}^{(\ell)}(\mathbf{X}) = \sum_{p=1}^{P^{(\ell)}} \sum_{j=1}^{C^{(\ell-1)}} W_{i,j,p}^{(\ell)} A_{j,\tilde{q}(p)}^{(\ell-1)}(\mathbf{X}), \quad (6)$$

which is just a spatial sum of matrix multiplications

$$\mathbf{Z}_{:,q}^{(\ell)}(\mathbf{X}) = \sum_{p=1}^{P^{(\ell)}} \mathbf{W}_{:::,p}^{(\ell)} \mathbf{A}_{:, \tilde{q}(p)}^{(\ell-1)}(\mathbf{X}). \quad (7)$$

Thus, we may express a convolution with multiple filters as a sum of matrix-vector multiplications. This is the canonical way to represent convolutional filters in NETSOR [Yang, 2019, NETSOR program 4].

Here we run into a problem. Equation 9 states that CNN filters are spatially correlated, but assumption B.4 states that A-vars have to be independent. To solve this, we will use the the following well-known lemma, which is the $\mathbf{R}\epsilon$ expression of a Gaussian random variable with mean zero. The tensor \mathbf{R} is a square root of the covariance.

Lemma B.9. Let $\Sigma \in \mathbb{R}^{P^2}$ be an arbitrary real-valued covariance tensor. Then there exists another real-valued tensor $\mathbf{R} \in \mathbb{R}^{P^2}$ such that $\Sigma_{q,q'} = \sum_{p=1}^P R_{q,p} R_{q',p}$. Next, let $\mathbf{u}, \mathbf{w} \in \mathbb{R}^P$ be real-valued tensors, such that $\mathbf{w} = \mathbf{R}\mathbf{u}$. Suppose the elements of \mathbf{u} are i.i.d. standard Gaussian variables, $\{u_p\}_{p \in [P]} \stackrel{iid}{\sim} \mathcal{N}(0, 1)$. Then, \mathbf{w} has a multivariate Gaussian distribution with mean zero and covariance tensor Σ .

Proof. Let $K = |P|$, and $\tilde{\Sigma}$ be $K \times K$ matrices, obtained by flattening the dimensions of Σ respectively. Then $\tilde{\Sigma}$ is a real-valued covariance matrix, so it is positive semi-definite and a square matrix $\tilde{\mathbf{R}}$ s.t. $\tilde{\mathbf{R}}\tilde{\mathbf{R}}^\top = \tilde{\Sigma}$ always exists. Un-flattening $\tilde{\mathbf{R}}$ we obtain \mathbf{R} . The variable \mathbf{w} is Gaussian because it is a linear transformation of the Gaussian \mathbf{u} . Calculating the second moment of \mathbf{w} finishes the proof. \square

Thus, to express convolution in NETSOR with correlated weights \mathbf{w} , we can use the following strategy. First, express several convolutions with uncorrelated weights \mathbf{u} , using eq. 7. Then, combine the output of the convolutions using `LinComb` and coefficients of the tensor \mathbf{R} .

Given a collection of A-vars $\{\mathbf{U}_{:::,p}^{(\ell)}\}_{p \in [P^{(\ell)}]}$, we can express the convolutional weights $\mathbf{W}^{(\ell)}$ which have covariance $\Sigma^{(\ell)} = \mathbf{R}^{(\ell)}(\mathbf{R}^{(\ell)})^\top$ as

$$\mathbf{W}_{:::,p}^{(\ell)} = \sum_{s=1}^{P^{(\ell)}} R_{p,s}^{(\ell)} \mathbf{U}_{:::,s}^{(\ell)}. \quad (8)$$

Substituting this into eq. 7 we obtain

$$\mathbf{Z}_{:,q}^{(\ell)}(\mathbf{X}) = \sum_{p=1}^{P^{(\ell)}} \sum_{s=1}^{P^{(\ell)}} R_{p,s}^{(\ell)} \mathbf{U}_{:::,s}^{(\ell)} \mathbf{A}_{:, \tilde{q}(p)}^{(\ell-1)}(\mathbf{X}). \quad (9)$$

To express this computation with NETSOR rules we may write

$$\text{MatMul: } \mathbf{H}_{:,s,p}^{(\ell)}(\mathbf{X}) \triangleq \mathbf{U}_{:,s}^{(\ell)} \mathbf{A}_{:,p}^{(\ell-1)}(\mathbf{X})$$

$$\text{for } \mathbf{s} \in [\mathbf{P}^{(\ell)}], \mathbf{p} \in [\mathbf{F}^{(\ell-1)}], \quad (10)$$

$$\text{LinComb: } \mathbf{Z}_{:,q}^{(\ell)}(\mathbf{X}) \triangleq + \sum_{p=1}^{\mathbf{P}^{(\ell)}} \sum_{s=1}^{\mathbf{P}^{(\ell)}} R_{p,s}^{(\ell)} \mathbf{H}_{:,s,\tilde{q}(p)}^{(\ell)}(\mathbf{X})$$

$$\text{for } \mathbf{q} \in [\mathbf{F}^{(\ell)}]. \quad (11)$$

Algorithm 1 uses this construction for every layer to express an L -layer CNN with correlated weights, applied to an input data set $\mathcal{X} \triangleq [\mathbf{X}_1, \dots, \mathbf{X}_M]$.

Theorem B.10 (Correlated CNN behaves like a GP). *Consider a countable set of input points $\tilde{\mathcal{X}}$, and a fixed number of layers L . Apply the L -layer convolutional neural network (eqs. 1 and 2) with correlated weights (eq. 9) to \mathcal{X} . Assume its nonlinearities are controlled (assumption B.3). For simplicity, fix all layers to have the same number of channels: $C = C^{(1)} = \dots = C^{(L)}$. Then, as the number of channels $C \rightarrow \infty$, the activations $\mathbf{Z}^{(L)}(\tilde{\mathcal{X}})$ converge in distribution to a Gaussian process with mean function $\mathbb{E} [Z_i^{(L)}(\tilde{\mathcal{X}})] = m^{(L)}(\tilde{\mathcal{X}})$ and covariance function $\mathbb{C} [Z_{i,q}^{(L)}(\tilde{\mathcal{X}}), Z_{i',q'}^{(L)}(\tilde{\mathcal{X}})] = \delta_{i,i'} K_{q,q'}^{(L)}(\tilde{\mathcal{X}}, \tilde{\mathcal{X}})$ (section 3).*

Proof. We need to show

1. that algorithm 1, including the postprocessing part, implements a correlated-weight CNN correctly,
2. that the full program converges weakly to a Gaussian process (GP) on $\tilde{\mathcal{X}}$,
3. that the moments of this GP match the ones in section 3.

For the first claim, the key is the equivalence between a convolutional layer with correlated weights (eq. 1), and a spatial outer product followed by linear combination (eqs. 10 and 11). Keeping this in mind, we can verify by inspection that the steps of algorithm 1, including the output postprocessing, implement the recursive CNN equations (eqs. 1 and 2).

The second claim is somewhat more involved. Invoking the Kolmogorov extension theorem [Tao, 2011, Thm. 2.4.3] we restrict our attention to finite subsets $\mathcal{X} \subseteq \tilde{\mathcal{X}}$, which are going to be compatible distributions if claim 3 is true. Since $\mathbf{X} \in \mathcal{X}$ are tensors, we may use the Euclidean metric. We then show by theorem B.7 that the output tuple of the CNN NETSOR program converges weakly to a GP as $C \rightarrow \infty$. Each activation in the postprocessing is defined as a linear combination of a Gaussian random variable (RV) (the bias) and a RV that converges weakly to a Gaussian, and thus the

Algorithm 1 NETSOR description of an L -layer CNN with correlated weights, with input \mathcal{X} .

```

/* G-vars for layer 1 activations, for
   all spatial locations  $\mathbf{p}$  and input
   points  $\mathbf{X}_m$ . */
Input :  $Z_p^{(1)}(\mathbf{X}_m) : \mathbf{G}(C^{(1)})$ 
         for  $\mathbf{p} \in [\mathbf{F}^{(1)}]$  and  $m = 1, \dots, M$ .
/* A-vars for the independent
   convolutional weights */
Input :  $U_p^{(\ell)} : \mathbf{A}(C^{(\ell)}, C^{(\ell-1)})$ 
         for  $\mathbf{p} \in [\mathbf{P}^{(\ell)}]$  and  $\ell = 2, \dots, L-1$ .
/* O-vars for the output, for every
   patch location  $\mathbf{s}$  and channel  $i$  */
Input :  $o_{i,s} : \mathbf{O}(C^{(L-1)})$ 
         for  $\mathbf{s} \in [\mathbf{P}^{(L)}]$  and  $i = 1, \dots, C^{(L)}$ .

for  $m = 1, \dots, M$  (data points of  $m$ ) do
  for  $\ell = 2, \dots, L-1$  (layer  $\ell$ ) do
    for  $\mathbf{p} = 1, \dots, \mathbf{F}^{(\ell-1)}$  do
      Nonlin:  $\mathbf{H}(C^{(\ell-1)})$ 
       $\mathbf{A}_{:,p}^{(\ell-1)}(\mathbf{X}_m) \triangleq \phi(\mathbf{Z}_{:,p}^{(\ell-1)}(\mathbf{X}_m))$ 
      for  $\mathbf{s} = 1, \dots, \mathbf{P}^{(\ell)}$  (patch location  $\mathbf{s}$ ) do
        MatMul:  $\mathbf{G}(C^{(\ell)})$ 
         $\mathbf{H}_{:,s,p}^{(\ell)}(\mathbf{X}_m) \triangleq U_s^{(\ell)} \mathbf{A}_p^{(\ell-1)}(\mathbf{X}_m)$ 
      end
    end
    for  $\mathbf{q} = 1, \dots, \mathbf{F}^{(\ell)}$  (spatial location  $\mathbf{q}$ ) do
      LinComb:  $\mathbf{G}(C^{(\ell)}) : \mathbf{Z}_{:,q}^{(\ell)}(\mathbf{X}_m)$ 
       $\triangleq \sum_{p=1}^{\mathbf{P}^{(\ell)}} \sum_{s=1}^{\mathbf{P}^{(\ell)}} R_{p,s}^{(\ell)} \mathbf{H}_{:,s,\tilde{q}(p)}^{(\ell)}(\mathbf{X}_m)$ 
    end
  end
  for  $\mathbf{p} \in [\mathbf{F}^{(L-1)}]$  (spatial location  $\mathbf{p}$ ) do
    Nonlin:  $\mathbf{H}(C^{(L-1)})$ 
     $\mathbf{A}_{:,p}^{(L-1)}(\mathbf{X}_m) \triangleq \phi(\mathbf{Z}_{:,p}^{(L-1)}(\mathbf{X}_m))$ 
  end
end
/* One output for every spatial
   location  $\mathbf{p}$ , patch location  $\mathbf{s}$ ,
   channel  $i$  and data point  $m$ . */
Output :  $(o_{i,s}^\top \mathbf{A}_{:,p}^{(L-1)}(\mathbf{X}_m)) : \text{for } \mathbf{p} \in [\mathbf{F}^{(L)}], \mathbf{s} \in [\mathbf{P}^{(L)}],$ 
          $i \in [C^{(L)}]$  and  $m \in [M]$ 

Output postprocessing: correlate the outputs and add bi-
ases (not part of NETSOR)
for  $m \in [M], i \in [C^{(L)}]$  and  $\mathbf{q} \in \mathbf{F}^{(L)}$  do
   $Z_{i,q}^{(L)}(\mathbf{X}_m) \triangleq \sum_{p=1}^{\mathbf{P}^{(L)}} \sum_{s=1}^{\mathbf{P}^{(L)}} R_{p,s}^{(L)} (o_{i,s}^\top \mathbf{A}_{:,q(p)}^{(L-1)}(\mathbf{X}_m))$ 
end

```

resulting distribution on \mathcal{X} converges weakly to a Gaussian too.

Finally, we have to show that the postprocessed output has the correct kernel. The output tuple and the activations have mean zero, which is correct. We thus compute the covariance of the output tuple in algorithm 1, for $\mathbf{X}, \mathbf{X}' \in \tilde{\mathcal{X}}$:

$$\begin{aligned} \mathbb{C} \left[\mathbf{o}_{i,s}^\top \mathbf{A}_{:,p}^{(L-1)}(\mathbf{X}), \mathbf{o}_{i',s'}^\top \mathbf{A}_{:,p'}^{(L-1)}(\mathbf{X}') \right] \\ = \delta_{i,i'} \delta_{s,s'} V_{p,p'}^{(L-1)}(\mathbf{X}, \mathbf{X}'). \end{aligned} \quad (12)$$

The delta functions appear because the O-vars and their elements are all independent. Using this, we can calculate the covariance function of the activations

$$\begin{aligned} \mathbb{C} \left[Z_{i,q}^{(L)}(\mathbf{X}), Z_{i',q'}^{(L)}(\mathbf{X}') \right] &= \sum_{\mathbf{p}, \mathbf{p}'}^{\mathbf{P}^{(L)2}} \sum_{\mathbf{s}, \mathbf{s}'}^{\mathbf{P}^{(L)2}} \\ R_{\mathbf{p},\mathbf{s}}^{(L)} R_{\mathbf{p}',\mathbf{s}'}^{(L)} \mathbb{C} \left[\mathbf{o}_{i,s}^\top \mathbf{A}_{:,\tilde{q}(\mathbf{p})}^{(L-1)}(\mathbf{X}), \mathbf{o}_{i',s'}^\top \mathbf{A}_{:,\tilde{q}'(\mathbf{p}')}^{(L-1)}(\mathbf{X}') \right]. \end{aligned} \quad (13)$$

Substitute the value of the expectations and eliminate one of the sums due to $\delta_{s,s'}$,

$$\begin{aligned} \mathbb{C} \left[Z_{i,q}^{(L)}(\mathbf{X}), Z_{i',q'}^{(L)}(\mathbf{X}') \right] \\ = \delta_{i,i'} \left(\sum_{\mathbf{p}, \mathbf{p}'}^{\mathbf{P}^{(L)2}} \sum_{\mathbf{s}=1}^{\mathbf{P}^{(L)}} R_{\mathbf{p},\mathbf{s}}^{(L)} R_{\mathbf{p}',\mathbf{s}}^{(L)} V_{\tilde{q}(\mathbf{p}), \tilde{q}'(\mathbf{p}')}^{(L-1)}(\mathbf{X}, \mathbf{X}') \right). \end{aligned} \quad (14)$$

Finally, noting that $\Sigma_{\mathbf{p},\mathbf{p}'}^{(L)} = \sum_{\mathbf{s}=1}^{\mathbf{P}^{(L)}} R_{\mathbf{p},\mathbf{s}}^{(L)} R_{\mathbf{p}',\mathbf{s}}^{(L)}$ (lemma B.9), and using the definition of $K_{q,q'}^{(L)}(\mathbf{X}, \mathbf{X}')$ (eq. 11), we obtain the claim.

$$\begin{aligned} &= \delta_{i,i'} \left(\sum_{\mathbf{p}, \mathbf{p}'}^{\mathbf{P}^{(L)2}} \Sigma_{\mathbf{p},\mathbf{p}'}^{(L)} V_{\tilde{q}(\mathbf{p}), \tilde{q}'(\mathbf{p}')}^{(L-1)}(\mathbf{X}, \mathbf{X}') \right) \\ &= \delta_{i,i'} K_{q,q'}^{(L)}(\mathbf{X}, \mathbf{X}'). \end{aligned} \quad (15)$$

□

References

- Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. URL <https://www.deeplearningbook.org/>.
- Jeffrey S Rosenthal. *A First Look At Rigorous Probability Theory*. World Scientific Publishing Company, 2006.
- Terence Tao. *An introduction to measure theory*, volume 126. American Mathematical Society Providence, RI, 2011.
- Greg Yang. Wide feedforward or recurrent neural networks of any architecture are Gaussian processes. In *Advances in Neural Information Processing Systems 32 (NeurIPS)*. 2019.