

# Brief technical note on linearizing recurrent neural networks (RNNs) before vs after the pointwise nonlinearity

Marino Pagan<sup>1,†</sup>, Adrian Valente<sup>2,†</sup>, Srdjan Ostojic<sup>2,\*</sup>, and Carlos D. Brody<sup>3,\*</sup>

<sup>1</sup>Simons Initiative for the Developing Brain, University of Edinburgh, Edinburgh, UK

<sup>2</sup>Laboratoire de Neurosciences Cognitives et Computationnelles, INSERM U960, Ecole Normale Supérieure - PSL Research University, 75005 Paris, France

<sup>3</sup>Howard Hughes Medical Institute and Princeton Neuroscience Institute, Princeton University, Princeton NJ, USA

<sup>†</sup>equal Contribution

<sup>\*</sup>equal Contribution

February 2023

## Abstract

Linearization of the dynamics of recurrent neural networks (RNNs) is often used to study their properties. The same RNN dynamics can be written in terms of the “activations” (the net inputs to each unit, before its pointwise nonlinearity) or in terms of the “activities” (the output of each unit, after its pointwise nonlinearity); the two corresponding linearizations are different from each other. This brief and informal technical note describes the relationship between the two linearizations, between the left and right eigenvectors of their dynamics matrices, and shows that some context-dependent effects

are readily apparent under linearization of activity dynamics but not linearization of activation dynamics.

## Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>   | <b>2</b> |
| <b>2</b> | <b>Two linearizations for the same discrete-time RNN Dynamics</b> | <b>3</b> |
| <b>3</b> | <b>Left and right eigenvectors of the dynamics matrices</b>       | <b>6</b> |
| <b>4</b> | <b>Linearizations and context-dependence of input vectors</b>     | <b>8</b> |
| <b>5</b> | <b>Conclusion</b>   | <b>9</b> |

## 1 Introduction

Recurrent neural network (RNN) dynamics can be equivalently expressed in two different forms [3]. One form describes the dynamics of the net input, or “activations” of the units, usually interpreted as the membrane potential of biological neurons. A second form describes dynamics in terms of the output, i.e. “activity” or “rate” of the units, often thought of as comparable to spiking rates of biological neurons. A pointwise nonlinearity relates the two, with the activity being the result of the nonlinearity after it is applied to the activation.

Linearization of dynamics is often used to study the properties of dynamical systems. But when considering an RNN, should one linearize the activity dynamics? Or the activation dynamics? The two linearizations produce different linear equations. What is the difference between them and what is the relationship between the two? Do some conclusions depend on which linearization is chosen?

This document explores these questions, and makes the relationship between the two linearizations explicit. The two are related by a simple diagonal linear transform that depends on the gains of each unit.

We additionally briefly consider effects of the two linearizations when considering context-dependent networks [2, 4], in which each “context” is defined by a constant vector of inputs to each unit, and point out that a

modulation by context of the linearized inputs to the RNN is observable only in the activity space linearization, not in the activation space linearization.

## 2 Two linearizations for the same discrete-time RNN Dynamics

Consider the standard recurrent neural network equations

$$\begin{aligned}\hat{\mathbf{x}}^{k+1} &= W\hat{\mathbf{r}}^k + \mathbf{u}^k \\ \hat{\mathbf{r}}^{k+1} &= g(\hat{\mathbf{x}}^{k+1})\end{aligned}\tag{1}$$

where  $\hat{\mathbf{x}}^k$  represents the vector of unit activations at timepoint  $k$ ,  $g(\cdot)$  is a differentiable and invertible pointwise nonlinear function such as  $\tanh(\cdot)$ ,  $\hat{\mathbf{r}}^k$  is the vector of unit activities at timepoint  $k$ ,  $W$  is a square matrix representing recurrent connection weights, and  $\mathbf{u}^k$  is a vector of external inputs at timepoint  $k$ .

The dynamics of (1) can be rewritten entirely in terms of  $\hat{\mathbf{x}}$ . As we do that, let us define the vector-valued dynamics function  $\mathbf{F}_x$ :

$$\mathbf{F}_x(\hat{\mathbf{x}}, \mathbf{u}) = Wg(\hat{\mathbf{x}}) + \mathbf{u},\tag{2}$$

so that

$$\hat{\mathbf{x}}^{k+1} = \mathbf{F}_x(\hat{\mathbf{x}}^k, \mathbf{u}^k).\tag{3}$$

Similarly, we can define the dynamics function  $\mathbf{F}_r$

$$\mathbf{F}_r(\hat{\mathbf{r}}, \mathbf{u}) = g(W\hat{\mathbf{r}} + \mathbf{u})\tag{4}$$

and rewrite the dynamics (1) entirely in terms of  $\hat{\mathbf{r}}$ ,

$$\hat{\mathbf{r}}^{k+1} = \mathbf{F}_r(\hat{\mathbf{r}}^k, \mathbf{u}^k)\tag{5}$$

We will consider the effects of linearizing around a fixed point when the dynamics are written in terms of  $\mathbf{F}_x$  versus when they are written in terms of  $\mathbf{F}_r$ .

To begin, consider a point specified by

$$\begin{aligned}\hat{\mathbf{x}}_0 \\ \mathbf{u}_0 &= \mathbf{0} \\ \hat{\mathbf{r}}_0 &= g(\hat{\mathbf{x}}_0)\end{aligned}\tag{6}$$

which we choose to be a fixed point of the dynamics (1), i.e., it is such that

$$\hat{\mathbf{x}}_0 = \mathbf{F}_x(\hat{\mathbf{x}}_0, \mathbf{u}_0).\tag{7}$$

Linearizing  $\mathbf{F}_x$  around that fixed point, we obtain

$$\hat{\mathbf{x}}^{k+1} = \mathbf{F}_x(\hat{\mathbf{x}}^k, \mathbf{u}^k) \approx \mathbf{F}_x(\hat{\mathbf{x}}_0) + \frac{\partial \mathbf{F}_x}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}^k - \hat{\mathbf{x}}_0) + \frac{\partial \mathbf{F}_x}{\partial \mathbf{u}} \mathbf{u}^k\tag{8}$$

Inserting Eq. 7, we obtain

$$\hat{\mathbf{x}}^{k+1} - \hat{\mathbf{x}}_0 = \frac{\partial \mathbf{F}_x}{\partial \hat{\mathbf{x}}}(\hat{\mathbf{x}}^k - \hat{\mathbf{x}}_0) + \frac{\partial \mathbf{F}_x}{\partial \mathbf{u}} \mathbf{u}^k.\tag{9}$$

Changing variables to

$$\mathbf{x}^k = \hat{\mathbf{x}}^k - \hat{\mathbf{x}}_0,\tag{10}$$

we arrive at

$$\mathbf{x}^{k+1} = \frac{\partial \mathbf{F}_x}{\partial \hat{\mathbf{x}}} \mathbf{x}^k + \frac{\partial \mathbf{F}_x}{\partial \mathbf{u}} \mathbf{u}^k.\tag{11}$$

In index notation, the two matrices involved in (11) are

$$\left[ \frac{\partial \mathbf{F}_x}{\partial \hat{\mathbf{x}}} \right]_{ij} = W_{ij} g'(\hat{x}_{0j})\tag{12}$$

$$\left[ \frac{\partial \mathbf{F}_x}{\partial \mathbf{u}} \right]_{ij} = \delta_{ij}\tag{13}$$

Let us define a diagonal matrix  $D$ , i.e., with zeros on all the non-diagonals, and entries along the diagonal that are each a function of the  $j^{\text{th}}$  element of  $\hat{\mathbf{x}}_0$  :

$$D_{jj} = g'(\hat{x}_{0j}). \quad (14)$$

Since its elements are the gains of  $g$  for each element of  $\hat{x}_0$ , we will call this matrix **the diagonal gain matrix**  $D$ . Then, in matrix notation, we can use  $D$  to rewrite the linearized dynamics (11) as

$$\mathbf{x}^{k+1} = WD\mathbf{x}^k + \mathbf{u}^k \quad (15)$$

The second linearization is obtained by differentiating equation (5) with respect to  $\hat{\mathbf{r}}$  and changing variables to

$$\mathbf{r}^k = \hat{\mathbf{r}}^k - \hat{\mathbf{r}}_0. \quad (16)$$

This requires the derivatives

$$\left[ \frac{\partial \mathbf{F}_r}{\partial \hat{\mathbf{r}}} \right]_{ij} = g'(\hat{x}_{0i})W_{ij} \quad (17)$$

$$\left[ \frac{\partial \mathbf{F}_r}{\partial \mathbf{u}} \right]_{ij} = g'(\hat{x}_{0i})\delta_{ij} \quad (18)$$

which we rewrite in matrix notation as

$$\mathbf{r}^{k+1} = DW\mathbf{r}^k + D\mathbf{u}^k \quad (19)$$

The two linear dynamical systems (15) and (19) might appear at first sight to be quite disparate.  $DW$  represents a scaling of the *rows* of  $W$  by the diagonal elements of  $D$ , while  $WD$  represents a scaling of the *columns* of  $W$  by the diagonal elements of  $D$ . The results of the two scalings could be quite different, suggesting that different conclusions might be drawn from analyzing  $WD$  versus analyzing  $DW$ , even though they are both linearizations of the same dynamics around the same fixed point.

But this is not the case. The two equations describe dynamics in terms of different variables,  $\mathbf{x}$  and  $\mathbf{r}$ , but are in fact intimately related. If we express the dynamics in terms of the same variable, the two different linearizations lead to identical trajectories.

To relate the variables  $\mathbf{x}$  and  $\mathbf{r}$ , consider a linearization of  $g(\cdot)$  around  $\hat{\mathbf{x}}_0$  so that

$$\hat{\mathbf{r}} \approx g(\hat{\mathbf{x}}_0) + g'(\hat{\mathbf{x}}_0)(\hat{\mathbf{x}} - \hat{\mathbf{x}}_0) \quad (20)$$

Then, given that  $\hat{\mathbf{r}}_0 = g(\hat{\mathbf{x}}_0)$ , and using the variable changes (10) and (16), we can find the map relating  $\mathbf{r}$  and  $\mathbf{x}$ :

$$\mathbf{r} \approx g'(\hat{\mathbf{x}}_0) \mathbf{x} = D\mathbf{x} \quad (21)$$

This makes it plain that the two equations (15) and (19) are equivalent, related through the map in (21). That is, we can take equation (15), multiply it on the left by the gain matrix  $D$ , and obtain equation (19):

$$\begin{aligned} \mathbf{x}^{k+1} &= WD\mathbf{x}^k + \mathbf{u}^k \\ D\mathbf{x}^{k+1} &= DW D\mathbf{x}^k + D\mathbf{u}^k \\ \mathbf{r}^{k+1} &= DW\mathbf{r}^k + D\mathbf{u}^k \end{aligned} \quad (22)$$

This means that if we take a trajectory of points  $\mathbf{x}^k$  produced by the linearization of  $\mathbf{F}_x$  in (15), and map each  $\mathbf{x}^k$  onto its corresponding  $\mathbf{r}^k$  using (21), we will get exactly the set of  $\mathbf{r}^k$  that the linearization of  $\mathbf{F}_r$  in (19) would have produced. The two linearizations describe the same trajectories and thus the same dynamics, albeit mapped onto each other through  $D$ , as in (21).

### 3 Left and right eigenvectors of the dynamics matrices

As we have described, (15) and (19) are two views of the same dynamical trajectories. But they have different linearized dynamics matrices, respectively  $WD$  and  $DW$ , which in general have different eigendecompositions. The

right and left eigenvectors of linearized dynamics matrices determine many features of the dynamics, but as shown above, the dynamics are independent of the chosen linearization. This suggests that the eigendecompositions of the two matrices should be closely related, and here we show that indeed they are.

Let  $W$  be a square matrix and  $D$  be a diagonal matrix of the same size as  $W$ .

Let  $\mathbf{s}_r^T$  be a **left** eigenvector of matrix  $DW$ , with corresponding eigenvalue  $\lambda$ . In other words,

$$\mathbf{s}_r^T DW = \lambda \mathbf{s}_r^T \quad (23)$$

Multiplying on the right by  $D$  we obtain

$$\mathbf{s}_r^T DWD = \lambda \mathbf{s}_r^T D \quad (24)$$

which means that the vector  $\mathbf{s}_r^T D$  is a left eigenvector of the matrix  $WD$ , with eigenvalue  $\lambda$ .

In other words,

|  |
|--|
| <p>If <math>\mathbf{s}_r^T</math> is a <b>left</b> eigenvector of <math>DW</math> with eigenvalue <math>\lambda</math>, then</p> $\mathbf{s}_x^T = \mathbf{s}_r^T D \quad (25)$ <p>is a corresponding <b>left</b> eigenvector of <math>WD</math>, also with eigenvalue <math>\lambda</math>.</p> |
|--|

Similarly, let  $\boldsymbol{\rho}_x$  be a **right** eigenvector of  $WD$ , with eigenvalue  $\lambda$ . That is,

$$WD\boldsymbol{\rho}_x = \lambda\boldsymbol{\rho}_x \quad (26)$$

Multiplying on the left by  $D$  we obtain

$$DWD\boldsymbol{\rho}_x = \lambda D\boldsymbol{\rho}_x \quad (27)$$

which means that the vector  $D\boldsymbol{\rho}_x$  is a right eigenvector of the matrix  $DW$ , with eigenvalue  $\lambda$ .

In other words,

If  $\boldsymbol{\rho}_r$  is a **right** eigenvector of  $DW$  with eigenvalue  $\lambda$ , then

$$\boldsymbol{\rho}_x = D^{-1}\boldsymbol{\rho}_r \quad (28)$$

is a corresponding **right** eigenvector of  $WD$ , also with eigenvalue  $\lambda$ .

These relationships imply that the dot product between left and right eigenvectors is preserved:

$$\begin{aligned} \mathbf{s}_x^T \cdot \boldsymbol{\rho}_x &= \\ &= \mathbf{s}_r^T D \cdot D^{-1} \boldsymbol{\rho}_r \\ &= \mathbf{s}_r^T \cdot \boldsymbol{\rho}_r \end{aligned}$$

Note that, except for the case when  $W$  is rank 1, the relationship between the eigenvectors of  $W$  and the eigenvectors of  $WD$  or  $DW$  is in general non-trivial.

## 4 Linearizations and context-dependence of input vectors

Any given RNN will be defined by its weight matrix  $W$ , and trajectories on it will be induced by inputs  $\mathbf{u}^k$ , where  $k$  indexes timepoints. We define  $\mathbf{u}^k = 0$  for  $k < 0$ , and consider the case where the network is simulated over multiple different “runs” or “trials”, each of which begins at a timepoint  $k \ll 0$ , and evolves to some timepoint  $k > 0$ . Let us now consider a situation in which there are additional inputs to the units of the network, constant in time during each run, but potentially different across different runs. That is, during each run  $R$ , the dynamical equations are

$$\begin{aligned} \hat{\mathbf{x}}^{k+1} &= W\hat{\mathbf{r}}^k + \mathbf{u}^k + \mathbf{c}_R \\ \hat{\mathbf{r}}^{k+1} &= g(\hat{\mathbf{x}}^{k+1}) \end{aligned} \quad (29)$$

The inputs  $\mathbf{c}_R$  define what we will call *context*  $R$ .



Let us further suppose that before timepoint  $k = 0$  of each run in context  $R$ , and before any inputs  $\mathbf{u}$  are non-zero in that run, the network has settled into a fixed-point determined by  $\mathbf{c}_R$ . This fixed-point will be such that

$$\hat{\mathbf{x}}_0^R = Wg(\hat{\mathbf{x}}_0(\mathbf{c}_R)) + \mathbf{c}_R \quad (30)$$

and will have a corresponding gain matrix  $D_R$  whose diagonal entries are the elements of  $g'(\hat{\mathbf{x}}_0(\mathbf{c}_R))$ .

Following (15) and (19), let us define the linearization of the network for context  $R$  to be the linear network with dynamics

$$\mathbf{x}^{k+1} = WD_R\mathbf{x}^k + \mathbf{u}^k \quad (31)$$

and

$$\mathbf{r}^{k+1} = D_RW\mathbf{r}^k + D_R\mathbf{u}^k \quad (32)$$

Differences between two contexts  $A$  and  $B$  in how a network behaves will then correspond to different instantiations of the network, one determined by the gain matrix  $D_A$ , the other by the gain matrix  $D_B$ .

Notice that context-dependent modulation of the linearized input  $\mathbf{u}$  is observable only in the activity space linearization (32) (where the linearized input is  $D_R\mathbf{u}$ , and thus depends on the gain matrix  $D_R$ ). In the activation space linearization (31), the linearized input is always  $\mathbf{u}$ , independent of  $D_R$ .

Context-dependent input modulation of recurrent networks with a fixed input vector  $\mathbf{u}$  is studied, for example, in [1], who utilize activity space linearization (32) for this purpose: the linearized inputs  $D_R\mathbf{u}$  depend on context through  $D_R$ . In contrast, [2] used activation space linearization (31) when studying context dependence of RNN dynamics with fixed input vectors, and therefore did not study context-dependent input modulation.

## 5 Conclusion

In a recurrent neural network, the linear dynamics that result from linearization in activation space, and those that result from linearization in activity

space, are different. Nevertheless, the two linear dynamics describe the same underlying trajectories, albeit mapped onto each other through a scaling given by the gain of each of the network’s units.

Despite this close relationship between the two linearizations, the two are not interchangeable. In particular, context-dependent modulations of external inputs that follow from context-dependent changes in unit gains are directly observable as input modulations in the activity space linearization, but not in the activation space linearization.

## References

- [1] Niru Maheswaranathan and David Sussillo. “How recurrent networks implement contextual processing in sentiment analysis”. In: (Apr. 2020). arXiv: 2004.08013 [cs.CL].
- [2] Valerio Mante et al. “Context-dependent computation by recurrent dynamics in prefrontal cortex”. en. In: *Nature* 503.7474 (Nov. 2013), pp. 78–84.
- [3] Kenneth D Miller and Francesco Fumarola. “Mathematical equivalence of two common forms of firing rate models of neural networks”. en. In: *Neural Comput.* 24.1 (Jan. 2012), pp. 25–31.
- [4] Marino Pagan et al. “A new theoretical framework jointly explains behavioral and neural variability across subjects performing flexible decision-making”. en. In: *bioRxiv* (Nov. 2022), p. 2022.11.28.518207.