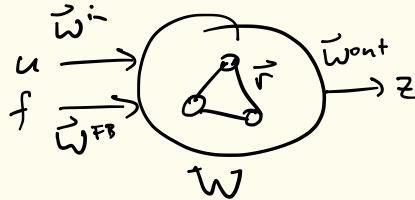


## Stability in trained RNNs

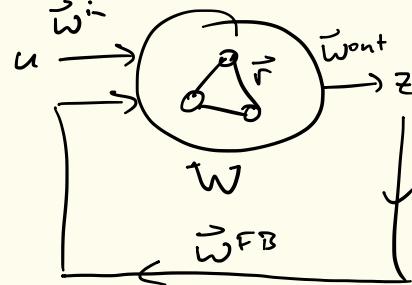
Ref: Rieck and Barak (PRL, 2017)

Train RNN with the echo state approach similar to Jaeger + Haas (2004).

Train in open loop:



Test in closed loop:



$$\dot{\vec{x}} = -\vec{x} + \vec{W} \cdot \vec{r} + \vec{w}^{FB} f + \vec{w}^{in} u \quad (1)$$

$$\dot{\vec{x}} = -\vec{x} + \vec{W} \cdot \vec{r} + \vec{w}^{FB} z + \vec{w}^{in} u \quad (2)$$

where  $\vec{r} = \phi(\vec{x})$ ,  $z = \vec{w}^{out} \cdot \vec{r}$ , and  $W_{ij} \sim \mathcal{N}(0, \frac{q^2}{N})$ .

Training is done with least squares so  $\vec{z} \approx \vec{f}$ .

Note: This differs from FORCE learning, where RLS is used rather than LS, and the training is done in closed loop.

Main question: If we train the RNN to produce a fixed point,  $f(t) = A$ , will it be stable when we close the loop?

(2)

Assume  $u=0$ . Then the open-loop fixed point from (1) is, for  $f(t)=A$ ,

$$\vec{\bar{x}} = W \cdot \phi(\vec{\bar{x}}) + \vec{w}^{FB} A \quad (1)$$

Is this solution stable? Let  $\vec{x} = \vec{\bar{x}} + \delta \vec{x}(t)$  in (1) and linearize:

$$\begin{aligned} \delta \dot{x}_i &= -(\bar{x}_i + \delta x_i) + \sum_j \omega_{ij} \phi(\bar{x}_j + \delta x_j) + \omega_{ii}^{FB} A \\ &\approx -(\bar{x}_i + \delta x_i) + \sum_j \omega_{ij} \phi(\bar{x}_j) + \sum_j \omega_{ij} \phi'(\bar{x}_j) \delta x_j \\ &\quad + \omega_{ii}^{FB} A \\ (3) \quad &= -\delta x_i + \sum_j \omega_{ij} \phi'(\bar{x}_j) \delta x_j \\ &= \sum_j P_{ij} \delta x_j, \quad P_{ij} = -\delta_{ij} + \omega_{ij} \phi'(\bar{x}_j) \quad (4) \end{aligned}$$

The solution is stable if all evals of  $P$  have real parts  $\leq 0$ . From random matrix theory, this happens if  $g^2 \langle |\phi'(\bar{x}_j)|^2 \rangle \leq 1$ , where  $\langle \dots \rangle$  is over neurons.

Is the solution (3) stable if we close the loop? Then rather than (4) we have

$$\delta \dot{x}_i = \underbrace{\sum_j [-\delta_{ij} + (\omega_{ij} + \omega_{ii}^{FB} \omega_j^{\text{out}}) \phi'(\bar{x}_j)]}_{Q_{ij}} \delta x_j \quad (5)$$

(3)

Now we need to compute evals of Q.  
 To do this, first return to the open-loop case and compute the response  $\delta \vec{x}(t)$  to a perturbation  $\delta f(t) = f(t) - A$  (cf. Eq. (4)):

$$\delta \dot{x}_i = -\delta x_i + \sum_j \omega_{ij} \phi'(\bar{x}_j) \delta x_j + \omega_{iF}^B \delta f \quad (6)$$

Let

$$\delta \vec{x}(t) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} e^{i\omega t} \underline{\underline{x}}_i(\omega) \quad (7)$$

$$\delta f(t) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} e^{i\omega t} F(\omega)$$

$$\delta z(t) = \int_{-\infty}^{\infty} \frac{d\omega}{2\pi} e^{i\omega t} \underline{z}(\omega)$$

$$\begin{aligned} \stackrel{(6)}{\Rightarrow} : \omega \underline{\underline{x}}_i(\omega) &= -\underline{\underline{x}}_i(\omega) + \sum_j \omega_{ij} \phi'(\bar{x}_j) \underline{\underline{x}}_j(\omega) \\ &\quad + \omega_{iF}^B F(\omega) \end{aligned} \quad (8)$$

The perturbation of the output is

$$\delta z(t) = \sum_i \omega_i^{\text{out}} \phi'(\bar{x}_i) \delta x_i(t) \quad (9)$$

$$\Rightarrow \underline{z}(\omega) = \sum_i \omega_i^{\text{out}} \phi'(\bar{x}_i) \underline{\underline{x}}_i(\omega) \quad (10)$$

(4)

The open-loop gain:

$$G^{OL}(\omega) = \frac{\mathcal{Z}(\omega)}{F(\omega)}$$

$$= \frac{1}{F(\omega)} \sum_i \omega_i^{\text{out}} \phi'(\bar{x}_i) \bar{x}_i(\omega) \quad (11)$$

The closed-loop gain is then

$$G^{CL}(\omega) = \frac{\mathcal{Z}}{F - \mathcal{Z}} = \frac{G^{OL}(\omega)}{1 - G^{OL}(\omega)} \quad (12)$$

We want to find the poles of  $G^{CL}(\omega)$ .  
First we need  $G^{OL}(\omega)$ .

In (3), let

$$\bar{x}_i = \bar{x}_i^0 + \bar{x}_i^1 \quad (13)$$

where

$$\bar{x}_i^0 = \omega_i^{\text{FB}} A \quad (14)$$

$$\bar{x}_i^1 = \sum_j \omega_{ij} \phi(\bar{x}_j) \quad (15)$$

Then, as in DMFT, assume  $\bar{x}_i^1 \sim \mathcal{N}(0, \sigma^2)$ ,  
where, from (15),  $\sigma$  is given self-consistently:

$$\langle \bar{x}_i^1 \rangle = \left\langle \sum_{jk} \omega_{ij} \omega_{ik} \phi(\bar{x}_j^0 + \bar{x}_j^1) \phi(\bar{x}_k^0 + \bar{x}_k^1) \right\rangle$$

$$= \sum_{jk} \langle \omega_{ij} \omega_{ik} \rangle \langle \phi(\bar{x}_j^0 + \bar{x}_j^1) \phi(\bar{x}_k^0 + \bar{x}_k^1) \rangle$$

$$= g^2 \langle \phi^2 (\omega_j^{\text{FB}} A + \bar{x}_j^1) \rangle$$

$$\textcircled{5} \Rightarrow \sigma^2 = g^2 \langle \phi^2 (\omega_j^{FB} A + \sigma_y) \rangle \quad (16)$$

$$= g^2 \int d\omega^{FB} \rho(\omega^{FB}) \int_{-\infty}^{\infty} dy e^{-y^2/2} \phi^2 (\omega^{FB} A + \sigma_y)$$

Similarly, let

$$\vec{\bar{X}}(\omega) = \vec{\bar{X}}^0 + \vec{\bar{X}}^1 \quad (17)$$

where, from (8),

$$\vec{\bar{X}}_j^0(\omega) = \frac{1}{1+i\omega} \omega_j^{FB} F(\omega) \quad (18)$$

$$\vec{\bar{X}}_j^1(\omega) = \frac{1}{1+i\omega} \sum_j W_{ij} \phi'(\bar{x}_j) \vec{\bar{X}}_j(\omega) \quad (19)$$

Now compute the overlap of  $\vec{\bar{X}}^1$  with  $\vec{\bar{x}}^1$ :

$$(1+i\omega) \sum_i \vec{\bar{x}}_i^1 \cdot \vec{\bar{X}}_i^1 \stackrel{(15, 19)}{=} \sum_{ijk} W_{ij} \phi(\bar{x}_j) W_{ik} \phi'(\bar{x}_k) \vec{\bar{X}}_k$$

$$\stackrel{(18, 19)}{=} \sum_{ijk} W_{ij} W_{ik} \phi(\bar{x}_j) \phi'(\bar{x}_k) \left[ \omega_k^{FB} \frac{F(\omega)}{1+i\omega} + \vec{\bar{X}}_k^1 \right] \quad (20)$$

Further decompose  $\vec{\bar{X}}^1 = \vec{\bar{X}}^{\parallel} + \vec{\bar{X}}^{\perp}$ , where

$$\vec{\bar{X}}_j^{\parallel}(\omega) = \alpha(\omega) \vec{\bar{x}}_j^1 \quad (21)$$

Then (20) is

$$(1+i\omega) \alpha(\omega) |\vec{\bar{x}}^1|^2 = \sum_{ijk} W_{ij} W_{ik} \phi(\bar{x}_j) \phi'(\bar{x}_k)$$

$$\times \left[ \omega_k^{FB} \frac{F(\omega)}{1+i\omega} + \alpha(\omega) \vec{\bar{x}}_k^1 + \vec{\bar{X}}_k^{\perp}(\omega) \right]$$

(6) Taking  $\langle \dots \rangle$  of this gives

$$(1+i\omega)\alpha(\omega)N\sigma^2 = g^2 \sum_k \langle \phi(\vec{x}_k) \phi'(\vec{x}_k) \times [\omega_k^{FB} \frac{F(\omega)}{1+i\omega} + \alpha(\omega) \vec{x}_k^1] \rangle \quad (22)$$

[The last term vanished "because  $\vec{x}^\perp$  is indep. of  $\vec{x}^1$ , and thus, on average, perp. to any function of  $\vec{x}^1$ ."]

Following (13)-(16), replace  $\vec{x}^1$  by random variables:

$$\vec{x}_i = wA + \sigma y$$

where  $w = \rho(\omega)$  and  $y \sim N(0, 1)$ .

$$\begin{aligned} \Rightarrow (1+i\omega)\alpha(\omega) &= \frac{g^2}{\sigma^2} \int Dw \int Dy \phi(wA + \sigma y) \phi'(wA + \sigma y) \\ &\quad \times [w \frac{F(\omega)}{1+i\omega} + \alpha(\omega) \sigma y] \\ &\equiv \frac{\beta_0 F(\omega)}{1+i\omega} + \beta_1 \alpha(\omega) \end{aligned} \quad (27)$$

$$\Rightarrow \alpha(\omega) = \frac{\beta_0 F(\omega)}{(1+i\omega)(1+i\omega - \beta_1)} \quad (28)$$

To get  $\tilde{\omega}^{\text{out}}$ , note that  $\tilde{\omega}^{\text{out}} \cdot \phi(\vec{x}) = A$  is solved by

$$\tilde{\omega}^{\text{out}} = \frac{A}{|\phi(\vec{x})|^2} \phi(\vec{x}) \stackrel{(16)}{\simeq} \frac{A}{N \langle \phi(\vec{x})^2 \rangle} \phi(\vec{x}) = \frac{g^2}{N \sigma^2} A \phi(\vec{x}) \quad (29)$$

⑦ Using this with (11) gives

$$\begin{aligned}
 G^{OL}(\omega) &= \frac{1}{F(\omega)} \sum_i \langle \omega^{\text{out}} - \phi'(\bar{x}_i) \Sigma_i(\omega) \rangle \\
 &= \frac{g^2 A}{N g^2 F(\omega)} \sum_i \langle \phi(\bar{x}_i) \phi'(\bar{x}_i) \Sigma_i(\omega) \rangle \\
 &= \frac{g^2 A}{N g^2 F(\omega)} \langle \phi(\bar{x}) \phi'(\bar{x}) \Sigma(\omega) \rangle
 \end{aligned}$$

$$\stackrel{(12)}{=} \frac{A}{N g^2 F(\omega)} N g^2 (1+i\omega) \alpha(\omega)$$

$$= \frac{A}{F(\omega)} \alpha(\omega) (1+i\omega)$$

$$\stackrel{(13)}{=} \frac{A \beta_0}{1+i\omega - \beta_1}$$

Then, from (12),

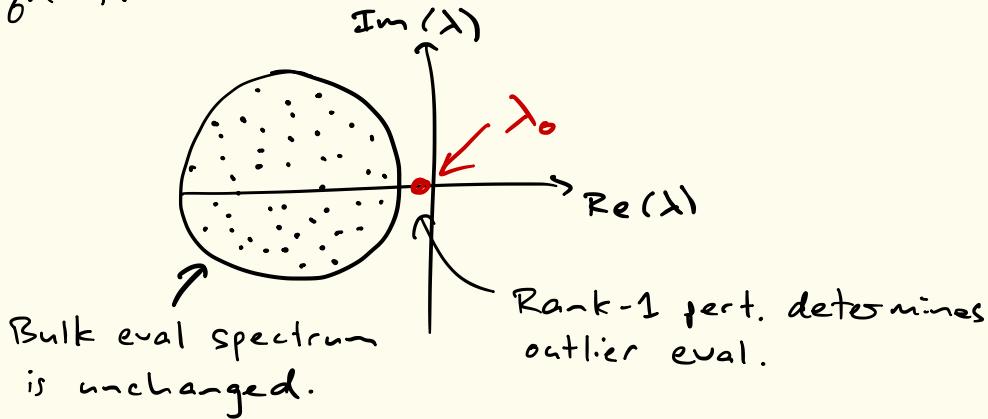
$$\begin{aligned}
 G^{CL}(\omega) &= \frac{1}{\frac{1+i\omega - \beta_1}{A \beta_0} - 1} \\
 &= \frac{A \beta_0}{i\omega - A \beta_0 - \beta_1 + 1}
 \end{aligned}$$
(25)

(26)

(8) This expression has a single pole at

$$\lambda_o = A\beta_0 + \beta_1 - 1 \quad (27)$$

This is the outlier eigenvalue of the effective CL connectivity matrix  $Q$  from Eq. (5).



Recipe:

- 1) Given  $g$  and  $A$ , calculate  $\sigma$  self-consistently from (16).
- 2) With  $g$ ,  $A$ , and  $\sigma$ , calculate  $\beta_0$  and  $\beta_1$  (p. 6).
- 3) Calculate  $\lambda_o$  from (27).

$$\left\{ \begin{array}{l} \lambda_o < 0 \Rightarrow \text{stable} \\ \lambda_o > 0 \Rightarrow \text{unstable} \end{array} \right.$$

(4)

Generalization to  $M$  outputs is possible.  
Resonance frequencies (large  $\zeta_{\text{cc}}(\omega \approx \omega_0)$ )  
can occur in this case.

Relation to Mastrogiuseppe + Ostojic (2019)

M+O: Given overlaps between  $\bar{\omega}^{\text{in}}$ ,  $\omega^{\text{out}}$ , and  $\bar{\omega}^{\text{FB}}$ , as well as target A, use DMFT to solve self-consistently for  $\omega^{\text{out}}$  and  $\Delta \equiv \text{Var}(x_i)$ . Then

- 1) Check stability of the  $z=1$  fixed point by computing the outlier eigenvalue with DMFT (M+O, 2018).
- 2) Check for other fixed points and test their stability.

Notes:

- i) Since DMFT doesn't require linearization, the fixed-point analysis is global, unlike Rikind's.
- ii) DMFT, unlike Rikind, assumes the recurrent input is random noise.