

Supplemental material

A. MinimalRNN Architecture

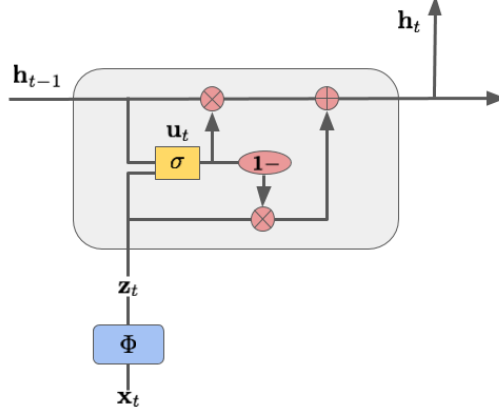


Figure Supp.1. Model architecture of minimalRNN.

B. Diagonal Recurrence Relation

Here we analyze the mean field dynamics of the minimalRNN. The minimalRNN features a hidden state $h^t \in \mathbb{R}^N$ and inputs x^t . The inputs are transformed via a fully-connected network $z^t = \Phi(x^t) \in \mathbb{R}^M$ before being fed into the network. The RNN cell is then described by the equations,

$$v_{i;a}^t = \sum_j W_{ij} h_{j;a}^{t-1} + \sum_j V_{ij} z_{j;a}^t + b_i \quad (16)$$

$$u_{i;a}^t = \sigma(v_{i;a}^t) \quad (17)$$

$$h_{i;a}^t = u_{i;a}^t h_{i;a}^{t-1} + (1 - u_{i;a}^t) z_{i;a}^t. \quad (18)$$

Here i denotes the (pre)-activation and a denotes an input to the network. Thus, $u_{i;a}^t$ acts as a gate on the t 'th step. We take $W_{ij} \sim \mathcal{N}(0, \sigma_w^2/N)$, $V_{ij} \sim \mathcal{N}(0, \sigma_v^2/M)$ and $b_i \sim \mathcal{N}(\mu_b, \sigma_b^2)$.

By the CTL we can make a mean field assumption that $v_{i;a}^t \sim \mathcal{N}(\mu_b, q_{ab}^t)$ where,

$$q_{ab}^t = \sigma_w^2 \mathbb{E}[h_{i;a}^{t-1} h_{i;b}^{t-1}] + \sigma_v^2 \mathbb{E}[z_{i;a}^t z_{i;b}^t] + \sigma_b^2 = \sigma_w^2 Q_{ab}^{t-1} + \sigma_v^2 R_{ab}^t + \sigma_b^2 \quad (19)$$

where we have defined $Q_{ab}^t = \mathbb{E}[h_{i;a}^t h_{i;b}^t]$ and $R_{ab}^t = \mathbb{E}[z_{i;a}^t z_{i;b}^t]$. We note that R_{ab}^t is fixed by the input, but it remains for us to work out Q_{ab}^t . We find that,

$$Q_{ab}^t = \mathbb{E}[h_{i;a}^t h_{i;b}^t] = \mathbb{E}[u_{i;a}^t h_{i;a}^{t-1} u_{i;b}^t h_{i;b}^{t-1}] + \mathbb{E}[u_{i;a}^t h_{i;a}^{t-1} (1 - u_{i;b}^t) z_{i;b}^t] \quad (20)$$

$$+ \mathbb{E}[(1 - u_{i;a}^t) z_{i;a}^t u_{i;b}^t h_{i;b}^{t-1}] + \mathbb{E}[(1 - u_{i;a}^t) z_{i;a}^t (1 - u_{i;b}^t) z_{i;b}^t] \quad (21)$$

$$\approx \mathbb{E}[u_{i;a}^t u_{i;b}^t] \mathbb{E}[h_{i;a}^{t-1} h_{i;b}^{t-1}] + \mathbb{E}[(1 - u_{i;a}^t)(1 - u_{i;b}^t)] \mathbb{E}[z_{i;a}^t z_{i;b}^t] \quad (22)$$

where we have assumed that the expectation factorizes so that $h_{i;a}^{t-1}$ and $u_{i;a}^t$ are approximately independent.

We choose to normalize the data so that $R_{aa}^t = R_{bb}^t = R$ independent of time. An immediate consequence of this normalization is that $Q_{aa}^t = Q_{bb}^t = Q^t$ and $q_{aa}^t = q_{bb}^t = q^t$. We then write $R_{ab}^t = R \Sigma^t$, $Q_{ab}^t = Q^t C^t$ and $q_{ab}^t = q^t c^t$ where Σ^t , C^t , and c^t are cosine similarities between the inputs, the hidden states, and the $v_{a,b}^t$ respectively. With this normalization, we can work out the mean-field recurrence relation characterizing the covariance matrix for the minimalRNN.

We begin by considering the diagonal recurrence relations. We find that the dynamics are described by the equation,

$$Q^t = Q^{t-1} \int \mathcal{D}z \sigma^2(\sqrt{q^t}z + \mu_b) + R \int \mathcal{D}z \left[1 - \sigma(\sqrt{q^t}z + \mu_b)\right]^2 \quad (23)$$

$$q^t = \sigma_w^2 Q^{t-1} + \sigma_v^2 R + \sigma_b^2 \quad (24)$$

As expected, the first and second integrands determine how much of the update of the random network is controlled by the norm of the hidden state and how much is determined by the norm of the input. Since $\sigma(z) = 1 - \sigma(-z)$ it follows that when $\mu_b = 0$ the first and second term will be equal and so,

$$Q^t = (Q^{t-1} + R) \int \mathcal{D}z \sigma^2(\sqrt{q^t}z). \quad (25)$$

In general, μ_b will therefore control the degree to which the hidden state of the random minimalRNN is updated based on the previous hidden state or based on the inputs with $\mu_b = 0$ implying parity between the two. This is reflected in eq. (25).

C. Existence of a Q^* Fixed Point

In the event that the norm of the inputs is time-independent, $R^t = R$ for all t , then the minimalRNN will have a fixed point provided there exists a Q^* that satisfies a transcendental equation, namely that

$$\mathcal{F}(Q^*) \equiv \frac{\int \mathcal{D}q^* z [1 - \sigma(z)]^2}{\int \mathcal{D}q^* z [1 - \sigma^2(z)]} - \frac{Q^*}{R} = 0. \quad (26)$$

It is easy to see that such a solution always exists. When $Q^* \rightarrow \infty$ the first term of $\mathcal{F}(Q^*)$ approaches 1 while the magnitude of the second increases without bound and so $\mathcal{F}(Q^*) < 0$. Conversely, when $Q^* \rightarrow 0$ the first term is positive while $Q^*/R \rightarrow 0$ and so $\mathcal{F}(Q^*) > 0$. The existence of a Q^* satisfying the transcendental equation then follows directly from the intermediate value theorem.

D. Q^* Dynamics

We can now investigate the dynamics of the norm of the hidden state in the vicinity of Q^* . To do this suppose that $Q^t = Q^* + \epsilon^t$ with $\epsilon \ll 1$. Our goal is then to expand eq.(23) about Q^* . First, we note that,

$$\sigma(\sqrt{q^t}z + \mu_b) = \sigma(\sqrt{q^* + \sigma_w^2 \epsilon^t}z + \mu_b) \quad (27)$$

$$\approx \sigma\left(\sqrt{q^*}z + \mu_b + \frac{1}{2\sqrt{q^*}}\sigma_w^2 \epsilon^t z\right) \quad (28)$$

$$\approx \sigma(\sqrt{q^*}z + \mu_b) + \frac{1}{2\sqrt{q^*}}\sigma_w^2 \epsilon^t z \sigma'(\sqrt{q^*}z + \mu_b) + \mathcal{O}((\epsilon^t)^2). \quad (29)$$

Letting $\zeta(z) = \sqrt{q^*}z + \mu_b$ this implies that,

$$Q^t = Q^{t-1} \int \mathcal{D}z \sigma^2(\sqrt{q^* + \sigma_w^2 \epsilon^{t-1}}z + \mu_b) + R \int \mathcal{D}z \left[1 - \sigma(\sqrt{q^* + \sigma_w^2 \epsilon^{t-1}}z + \mu_b)\right]^2 \quad (30)$$

$$Q^* + \epsilon^t = Q^* \int \mathcal{D}z \sigma^2(\zeta(z)) + R \int \mathcal{D}z [1 - \sigma(\zeta(z))]^2 + \epsilon^{t-1} \left[\int \mathcal{D}z \sigma^2(\zeta(z)) \right. \quad (31)$$

$$\left. + \frac{Q^*}{\sqrt{q^*}} \sigma_w^2 \int \mathcal{D}z z \sigma(\zeta(z)) \sigma'(\zeta(z)) - R \sqrt{q^*} \sigma_w^2 \int \mathcal{D}z z (1 - \sigma(\zeta(z))) \sigma'(\zeta(z)) \right] \quad (32)$$

$$\epsilon^t = \epsilon^{t-1} \left[\int \mathcal{D}z \sigma^2(\zeta(z)) + \frac{Q^*}{\sqrt{q^*}} \sigma_w^2 \int \mathcal{D}z z \sigma(\zeta(z)) \sigma'(\zeta(z)) - R \sqrt{q^*} \sigma_w^2 \int \mathcal{D}z z (1 - \sigma(\zeta(z))) \sigma'(\zeta(z)) \right] \quad (33)$$

$$= \epsilon^{t-1} \int \mathcal{D}z \left[\sigma^2(\zeta(z)) + \frac{\sigma_w^2}{\sqrt{q^*}} \{ (Q^* + R) \sigma(\zeta(z)) - R \} z \sigma'(\zeta(z)) \right] \quad (34)$$

$$= \epsilon^{t-1} \int \mathcal{D}z \left[\sigma^2(\zeta(z)) + \frac{\sigma_w^2}{\sqrt{q^*}} \{ (Q^* + R) \sigma(\zeta(z)) - R \} z \sigma(\zeta(z)) (1 - \sigma(\zeta(z))) \right] \quad (35)$$

$$= \epsilon^{t-1} \int \mathcal{D}z \left[\sigma^2(\zeta(z)) + \frac{\sigma_w^2}{\sqrt{q^*}} \{ (Q^* + R) \sigma(\zeta(z)) - R \} z \sigma(\zeta(z)) (1 - \sigma(\zeta(z))) \right] \quad (36)$$

It follows that $q^t \rightarrow q^*$ as,

$$|q^t - q^*| \sim e^{-t/\xi_Q} \quad (37)$$

with

$$\xi_Q^{-1} = -\log \left(\int \mathcal{D}z \left[\sigma^2(\zeta(z)) + \frac{\sigma_w^2}{\sqrt{q^*}} \{ (Q^* + R)\sigma(\zeta(z)) - R \} \sigma'(\zeta(z)) \right] \right) \quad (38)$$

as expected.

E. Off-Diagonal Recurrence Relation

We now turn our attention to the off-diagonal term. From eq. (9) it follows that,

$$Q^t C^t = Q^{t-1} C^{t-1} \int \mathcal{D}z_1 \mathcal{D}z_2 \sigma(u_1^t) \sigma(u_2^t) + R \Sigma^t \int \mathcal{D}z_1 \mathcal{D}z_2 (1 - \sigma(u_1^t))(1 - \sigma(u_2^t)) \quad (39)$$

$$q^t c^t = \sigma_w^2 Q^{t-1} C^{t-1} + \sigma_v^2 R \Sigma^t + \sigma_b^2 \quad (40)$$

where

$$u_1^t = \sqrt{q^t} z_1 + \mu_b \quad \text{and} \quad u_2^t = \sqrt{q^t} \left(c^t z_1 + \sqrt{1 - (c^t)^2} z_2 \right) + \mu_b. \quad (41)$$

By expanding eq (39) as $c^t = c^* + \epsilon^t$ we find $\epsilon^{t+1} = \chi_{c^*} \epsilon^t$ where,

$$\chi_{c^*} = \int \mathcal{D}z_1 \mathcal{D}z_2 \sigma(u_1) \sigma(u_2) + q^* (c^* + J_-) \int \mathcal{D}z_1 \mathcal{D}z_2 \sigma'(u_1) \sigma'(u_2). \quad (42)$$

We note that when $c^* = 1$ it follows that $\chi_{c^*} = \chi_1$.

F. Additional Hyperparameter Ranges

We tune the learning hyper-parameters in the following ranges for all the models:

- learning rate: {0.1, 0.2, 0.3, 0.5, 1, 2}
- max-epoch: {4, 7, 11}
- decay: {0.5, 0.65, 0.8}
- dropout: {0.0, 0.2, 0.3, 0.5}