

1 The Problem

We consider a feedforward neural network with a single hidden layer and activation function σ . It receives an input $x \in \mathbb{R}^n$ and produces a scalar output $\hat{y} \in \mathbb{R}$. The hidden layer has K units. The weights for the first and second layer are $W_1 \in \mathbb{R}^{K \times n}$ and $W_2 \in \mathbb{R}^{1 \times K}$, respectively, and the corresponding biases are $b_1 \in \mathbb{R}^K$ and $b_2 \in \mathbb{R}$.

$$\hat{y} = W_2 \sigma(W_1 x + b_1) + b_2. \quad (1)$$

Our data x are sampled from a mixture of two translation-invariant distributions in some family $\{p_\xi\}_\xi$ parameterized by a correlation length-scale ξ . That is, we sample $x \sim p_{\xi_1}$ with probability $\frac{1}{2}$ and $x \sim p_{\xi_2}$ otherwise. If x is sampled from p_{ξ_1} , then $y(x) = 1$; otherwise, $y(x) = 0$. We can train using either mean-squared error or cross-entropy loss, though we primarily consider the former.

Alessandro’s paper primarily considers the case where $W_2 = \frac{1}{K} \mathbf{1}^\top$ (take the mean of the hidden activations) is fixed and $\sigma(h) = \text{erf}(\frac{h}{\sqrt{2}})$. I have also tried $\sigma = \text{sigmoid}, \text{ReLU}$. For the former, the results are qualitatively identical, while for the latter we get localization if $\xi_1 > \xi_2$ and short-range oscillations otherwise. For $\sigma = \text{ReLU}$, one can further remove the bias terms b_1 and b_2 (though not for sigmoid).

We consider two types of datasets: the nonlinear Gaussian process (NLGP) and the single pulse (SP). We explain them in more detail later. They differ primarily in that the former has continuous support on \mathbb{R}^n , while the latter has discrete support on a subset of $\{0, 1\}^n$. The former also has a gain parameter that controls the degree of localization, while the latter does not.

We’ll present the results in reverse order, since it makes more sense logically. We start by attempting to analyze the ReLU model directly. This will force us to assume Gaussian data, which captures only half of what we’d like to describe. We will discuss some ideas about how to extend this to the non-Gaussian case, and perhaps also the SP dataset.

To address these analytical roadblocks, we’ll explore using a gated deep linear network (GDLN) to model the ReLU network. This will require some assumptions about the gating structure, which we test empirically. However, we’re not really sure how to properly set up the “neural race” and map the winner onto the ReLU case. We consider a few approaches, though we are not sure which is correct. We’ll conclude with some questions, concerns, and ideas, with specific focus on the discrete Fourier transform, uncertainty principle, and characteristic functions.

2 ReLU Analysis

To get an idea of what gating looks like early on during training, let’s try to analyze the gradient flow for a ReLU network. We won’t be able to solve it exactly, but we can get some intuition. We will have to assume the data are Gaussian to say something interesting after just a few steps. Of course, this is the case we are less interested in, since it’s the non-Gaussian data that shows localization.

Let w_i be the i -th row in W_1 . We make predictions with

$$\hat{y}(x) = \frac{1}{K} \sum_{k \in [K]} \text{ReLU}(\langle w_k, x \rangle). \quad (2)$$

We use MSE loss,

$$\mathcal{L} = \frac{1}{2} \mathbb{E}_{X,Y} \left[(\hat{y}(X) - Y)^2 \right]. \quad (3)$$

The corresponding gradient flow for w_1 is given by

$$\tau \frac{d}{dt} w_1 = - \mathbb{E}_{X,Y} \left[\left(\frac{1}{K} \sum_{k \in [K]} \text{ReLU}(\langle w_k, X \rangle) - Y \right) \frac{\partial}{\partial w_1} [\text{ReLU}(\langle w_1, X \rangle)] \right] \quad (4)$$

$$= \frac{1}{2} \mathbb{E}_{X, \langle w_1, X \rangle > 0 | Y=1} [X] - \frac{1}{K} \sum_{k \in [K]} \mathbb{E}_{X,Y, \langle w_1, X \rangle > 0, \langle w_k, X \rangle > 0} [\langle w_k, X \rangle X] \quad (5)$$

Gaussian Data

To compute these conditional expectations, we will have to assume the data are Gaussian. We begin by computing the first conditional expectation. Define the random variable $S = \langle w_1, X \rangle$.

$$\mathbb{E}_{X, S > 0 | Y=1} [X] = \mathbb{E}_{S | S > 0, Y=1} [\mathbb{E}_{X | S, Y=1} [X]] \mathbb{P}(S > 0 | Y = 1). \quad (6)$$

Let us consider X sampled from p_{ξ_1} (i.e. with label $Y = 1$) and write

$$X = AX + Sv, \quad (7)$$

where

$$v = \frac{1}{w_1^\top \Sigma_1 w_1} \Sigma_1 w_1, \quad (8)$$

$$A = I_n - v w_1^\top. \quad (9)$$

Equation (7) clearly holds. Our specific selection of v and A guarantees that AX and S have zero covariance. *Since X is Gaussian, this implies they are independent.* So, $X | S \sim \mathcal{N}(Sv, A\Sigma_1 A^\top)$. Note that $\mathbb{E}_{S | S > 0, Y=1} [S] = \left(\frac{2}{\pi} w_1^\top \Sigma_1 w_1 \right)^{\frac{1}{2}}$. Plugging this into equation (6),

$$\mathbb{E}_{X, S | S > 0, Y=1} [X] = \mathbb{E}_{S | S > 0, Y=1} [Sv] \mathbb{P}(S > 0 | Y = 1) \quad (10)$$

$$= \frac{1}{\sqrt{2\pi}} (w_1^\top \Sigma_1 w_1)^{-\frac{1}{2}} \Sigma_1 w_1. \quad (11)$$

Now, let us evaluate the second conditional expectation. First, we consider the case $k = 1$. Then, we only have one positivity constraint. We use S again, just as before. Let us also only consider X with label $Y = 1$.

$$\mathbb{E}_{X, S | S > 0, Y=1} [SX] = \mathbb{E}_{S | S > 0, Y=1} [S \mathbb{E}_{X | S, Y=1} [X]] \mathbb{P}(S > 0 | Y = 1) = \frac{1}{2} \mathbb{E}_{S | S > 0, Y=1} [S^2] v. \quad (12)$$

By symmetry of S about 0, $\mathbb{E}_{S | S > 0, Y=1} [S^2] = \mathbb{E}_{S | Y=1} [S^2] = w_1^\top \Sigma_1 w_1$. (This step does not require Gaussianity!) So,

$$\mathbb{E}_{X, \langle w_1, X \rangle > 0 | Y=1} [\langle w_1, X \rangle X] = \frac{1}{2} w_1^\top \Sigma_1 w_1 \left(\frac{1}{w_1^\top \Sigma_1 w_1} \Sigma_1 w_1 \right) = \frac{1}{2} \Sigma_1 w_1. \quad (13)$$

Now, let us consider $k > 1$. We will need to consider both positivity constraints. To do this, let us define $S = \langle w_1, X \rangle$ (as above) and $T = \langle w_k, X \rangle$. Again, let us focus on $Y = 1$.

$$\mathbb{E}_{X, S > 0, T > 0 | Y=1} [TX] = \mathbb{E}_{S > 0, T > 0 | Y=1} [T \mathbb{E}_{X | S, T} [X]] \mathbb{P}(S > 0, T > 0 | Y = 1). \quad (14)$$

Define

$$U = \begin{bmatrix} w_1^\top \\ w_k^\top \end{bmatrix}, \quad \text{and} \quad b = \begin{bmatrix} S \\ T \end{bmatrix}. \quad (15)$$

Then, we can write the inner expectation as

$$\mathbb{E}_{X|UX=b} [X]. \quad (16)$$

Using a similar trick as above,

$$X = AX + Cb, \quad (17)$$

where

$$C = \Sigma_1 U^\top (U^\top \Sigma_1 U)^{-1}, \quad (18)$$

$$A = I_n - CU. \quad (19)$$

Again, our specific selection of C and A guarantees that AX and Cb have zero covariance. Since X is Gaussian, this means we can write $\mathbb{E}_{X|UX=b} [X] = Cb$. So,

$$\mathbb{E}_{X,S>0,T>0|Y=1} [TX] = C \mathbb{E}_{S>0,T>0|Y=1} \left[\begin{bmatrix} S^2 \\ ST \end{bmatrix} \right] \mathbb{P}(S > 0, T > 0 | Y = 1). \quad (20)$$

We begin with the first term. Let us define $\rho_{ij} = w_i^\top \Sigma w_j^\top$ for $i, j = 1, k$.

$$\mathbb{E}_{S>0,T>0|Y=1} [S^2] \quad (21)$$

$$= \int_0^\infty \int_0^\infty s^2 \frac{1}{2\pi \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}(\rho_{11}s^2 - 2\rho_{1k}st + \rho_{kk}t^2)} dt ds \quad (22)$$

$$= \int_0^\infty s^2 \frac{1}{\sqrt{2\pi} \sqrt{\rho_{kk}}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}\left(-\frac{\rho_{1k}^2}{\rho_{kk}}s^2 + \rho_{11}s^2\right)} \int_0^\infty \frac{1}{\sqrt{2\pi} \left(\sqrt{\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}}\right)} e^{-\frac{1}{2\left(\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}\right)}\left(t - \frac{\rho_{1k}}{\rho_{kk}}s\right)^2} dt ds \quad (23)$$

$$= \int_0^\infty s^2 \frac{1}{\sqrt{2\pi} \sqrt{\rho_{kk}}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}\left(-\frac{\rho_{1k}^2}{\rho_{kk}}s^2 + \rho_{11}s^2\right)} \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\rho_{1k}}{\sqrt{2\rho_{kk}(\rho_{11}\rho_{kk} - \rho_{1k}^2)}}s\right)\right) ds \quad (24)$$

$$= \frac{\rho_{kk}}{2\pi} \left(\cos^{-1}\left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}}\right) + \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2} \right). \quad (25)$$

Now, the second term:

$$\mathbb{E}_{S>0, T>0|Y=1} [ST] \quad (26)$$

$$= \int_0^\infty \int_0^\infty st \frac{1}{2\pi \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}(\rho_{11}s^2 - 2\rho_{1k}st + \rho_{kk}t^2)} dt ds \quad (27)$$

$$= \int_0^\infty s \frac{1}{\sqrt{2\pi} \sqrt{\rho_{kk}}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)} \left(-\frac{\rho_{1k}^2}{\rho_{kk}} s^2 + \rho_{11}s^2 \right)} \underbrace{\int_0^\infty t \frac{1}{\sqrt{2\pi} \left(\sqrt{\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}} \right)} e^{-\frac{1}{2 \left(\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}} \right)} \left(t - \frac{\rho_{1k}}{\rho_{kk}} s \right)^2} dt}_{\text{mean of truncated normal with } \mu = \frac{\rho_{1k}}{\rho_{kk}} s, \sigma^2 = \frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}, a=0, b=\infty} ds \quad (28)$$

$$= \int_0^\infty s \frac{1}{\sqrt{2\pi} \sqrt{\rho_{kk}}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)} \left(-\frac{\rho_{1k}^2}{\rho_{kk}} s^2 + \rho_{11}s^2 \right)} \left[\frac{1}{2} \left(\frac{\rho_{1k}}{\rho_{kk}} s \right) \left(1 + \operatorname{erf} \left(\frac{\frac{\rho_{1k}}{\rho_{kk}} s}{\sqrt{2 \frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}}} \right) \right) + \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} \frac{\left(\frac{\rho_{1k}}{\rho_{kk}} s \right)^2}{\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}}} \left(\sqrt{\frac{\rho_{11}\rho_{kk} - \rho_{1k}^2}{\rho_{kk}}} \right) \right] ds \quad (29)$$

$$= \frac{\rho_{1k}}{2\pi} \left(\sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2} \right). \quad (30)$$

Thus,

$$C \mathbb{E}_{X, S>0, T>0|Y=1} [TX] \quad (31)$$

$$= C \left[\frac{\rho_{kk}}{2\pi} \left(\cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2} \right) \right. \quad (32)$$

$$\left. + \frac{\rho_{1k}}{2\pi} \left(\sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2} \right) \right] \quad (33)$$

Let's compute C

$$\Sigma_1 U^\top (U^\top \Sigma_1 U)^{-1} = \frac{1}{\rho_{11}\rho_{kk} - \rho_{1k}^2} \begin{bmatrix} \Sigma_1 w_1 & \Sigma_1 w_k \end{bmatrix} \begin{bmatrix} \rho_{kk} & -\rho_{1k} \\ -\rho_{1k} & \rho_{11} \end{bmatrix}. \quad (34)$$

Then,

$$C \begin{bmatrix} \rho_{kk} \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \\ \rho_{1k} \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \end{bmatrix} \quad (35)$$

$$= \frac{1}{\rho_{11}\rho_{kk} - \rho_{1k}^2} \begin{bmatrix} \Sigma_1 w_1 & \Sigma_1 w_k \end{bmatrix} \begin{bmatrix} \rho_{kk} & -\rho_{1k} \\ -\rho_{1k} & \rho_{11} \end{bmatrix} \begin{bmatrix} \rho_{kk} \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \\ \rho_{1k} \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \end{bmatrix} \quad (36)$$

$$= \frac{1}{\rho_{11}\rho_{kk} - \rho_{1k}^2} \begin{bmatrix} \Sigma_1 w_1 & \Sigma_1 w_k \end{bmatrix} \begin{bmatrix} \rho_{kk}^2 \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) - \rho_{1k}^2 \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \\ -\rho_{1k}\rho_{kk} \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \rho_{1k}\rho_{11} \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \end{bmatrix} \quad (37)$$

$$= \frac{1}{\rho_{11}\rho_{kk} - \rho_{1k}^2} \left[\left(\rho_{kk}^2 \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) - \rho_{1k}^2 \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \right) \Sigma_1 w_1 \right. \quad (38)$$

$$\left. + \left(-\rho_{1k}\rho_{kk} \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \rho_{1k}\rho_{11} \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \right) \Sigma_1 w_k \right]$$

And,

$$\frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \cdot \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2} C \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \cdot \frac{1}{\sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} \begin{bmatrix} \Sigma_1 w_1 & \Sigma_1 w_k \end{bmatrix} \begin{bmatrix} \rho_{kk} - \rho_{1k} \\ \rho_{11} - \rho_{1k} \end{bmatrix} \quad (39)$$

$$= \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \cdot \frac{1}{\sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} ((\rho_{kk} - \rho_{1k})\Sigma_1 w_1 + (\rho_{11} - \rho_{1k})\Sigma_1 w_k). \quad (40)$$

In toto,

$$\mathbb{E}_{X, \langle w_1, X \rangle > 0, \langle w_k, X \rangle > 0 | Y=1} [\langle w_k, X \rangle X] \quad (41)$$

$$= C \mathbb{E}_{X, S > 0, T > 0 | Y=1} [TX] \quad (42)$$

$$\begin{aligned} &= \frac{1}{2\pi} \cdot \frac{1}{\rho_{11}\rho_{kk} - \rho_{1k}^2} \left[\left(\rho_{kk}^2 \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) - \rho_{1k}^2 \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \right) \Sigma_1 w_1 \right. \\ &\quad \left. + \left(-\rho_{1k}\rho_{kk} \cos^{-1} \left(-\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) + \rho_{1k}\rho_{11} \sin^{-1} \left(\frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}} \right) \right) \Sigma_1 w_k \right] \\ &\quad + \frac{1}{2\pi} \cdot \frac{\rho_{1k}}{\rho_{11}\rho_{kk}} \cdot \frac{1}{\sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} ((\rho_{kk} - \rho_{1k})\Sigma_1 w_1 + (\rho_{11} - \rho_{1k})\Sigma_1 w_k). \end{aligned} \quad (43)$$

This is really messy! But, we can identify an important order parameter: $\gamma \equiv \frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}}$. This measures the cosine similarity between w_1 and w_k with respect to the inner product defined by Σ_1 . (Does the cosine identity still hold under this new inner product?)

$$\begin{aligned} (43) &= \frac{1}{2\pi(1-\gamma^2)\rho_{11}\rho_{kk}} \left[(\rho_{kk}^2 \cos^{-1}(-\gamma) - \rho_{1k}^2 \sin^{-1}(\gamma)) \Sigma_1 w_1 + (-\rho_{1k}\rho_{kk} \cos^{-1}(\gamma) + \rho_{1k}\rho_{11} \sin^{-1}(\gamma)) \Sigma_1 w_k \right] \\ &\quad + \frac{1}{2\pi\sqrt{1-\gamma^2}} \cdot \frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}^3} ((\rho_{kk} - \rho_{1k})\Sigma_1 w_1 + (\rho_{11} - \rho_{1k})\Sigma_1 w_k) \end{aligned} \quad (44)$$

$$\begin{aligned} &= \frac{1}{2\pi(1-\gamma^2)} \left[\left(\frac{\rho_{kk}}{\rho_{11}} \cos^{-1}(-\gamma) - \gamma^2 \sin^{-1}(\gamma) \right) \Sigma_1 w_1 + \left(-\frac{\rho_{1k}}{\rho_{11}} \cos^{-1}(\gamma) + \frac{\rho_{1k}}{\rho_{kk}} \sin^{-1}(\gamma) \right) \Sigma_1 w_k \right] \\ &\quad + \frac{1}{2\pi\gamma\sqrt{1-\gamma^2}\rho_{11}\rho_{kk}} ((\rho_{kk} - \rho_{1k})\Sigma_1 w_1 + (\rho_{11} - \rho_{1k})\Sigma_1 w_k) \end{aligned} \quad (45)$$

$$\begin{aligned} &= \left[\frac{1}{2\pi(1-\gamma^2)} \left(\frac{\rho_{kk}}{\rho_{11}} \cos^{-1}(-\gamma) - \gamma^2 \sin^{-1}(\gamma) \right) + \frac{\rho_{kk} - \rho_{1k}}{2\pi\gamma\sqrt{1-\gamma^2}\rho_{11}\rho_{kk}} \right] \Sigma_1 w_1 \\ &\quad + \left[\frac{1}{2\pi(1-\gamma^2)} \left(-\frac{\rho_{1k}}{\rho_{11}} \cos^{-1}(\gamma) + \frac{\rho_{1k}}{\rho_{kk}} \sin^{-1}(\gamma) \right) + \frac{\rho_{11} - \rho_{1k}}{2\pi\gamma\sqrt{1-\gamma^2}\rho_{11}\rho_{kk}} \right] \Sigma_1 w_k \end{aligned} \quad (46)$$

$$\begin{aligned} &= \left[\frac{1}{2\pi(1-\gamma^2)} \left(\frac{\rho_{kk}}{\rho_{11}} \cos^{-1}(-\gamma) - \gamma^2 \sin^{-1}(\gamma) \right) + \frac{\frac{1}{\gamma\rho_{11}} - \frac{1}{\sqrt{\rho_{11}\rho_{kk}}}}{2\pi\sqrt{1-\gamma^2}} \right] \Sigma_1 w_1 \\ &\quad + \left[-\frac{\gamma}{2\pi(1-\gamma^2)} \sqrt{\frac{\rho_{11}}{\rho_{kk}}} \left(\frac{\rho_{kk}}{\rho_{11}} \cos^{-1}(\gamma) - \sin^{-1}(\gamma) \right) + \frac{\frac{1}{\gamma\rho_{kk}} - \frac{1}{\sqrt{\rho_{11}\rho_{kk}}}}{2\pi\sqrt{1-\gamma^2}} \right] \Sigma_1 w_k \end{aligned} \quad (47)$$

$$\begin{aligned} &= \left[\frac{\gamma^2 - \frac{\rho_{kk}}{\rho_{11}}}{2\pi(1-\gamma^2)} \cos^{-1}(\gamma) + \frac{2\frac{\rho_{kk}}{\rho_{11}} - \gamma^2}{4(1-\gamma^2)} + \frac{\frac{1}{\gamma\rho_{11}} - \frac{1}{\sqrt{\rho_{11}\rho_{kk}}}}{2\pi\sqrt{1-\gamma^2}} \right] \Sigma_1 w_1 \\ &\quad + [??] \Sigma_1 w_k \end{aligned} \quad (48)$$

Hi!

Forgot to consider Σ_0 in gradient flow. Look for order parameters, I think it will be more meaningful.

Interpretation

This is really messy! How do we interpret what this is telling us? Let's consider two extreme cases, first where $\rho_{1k} = 0$, that is, w_1 and w_k are orthogonal w.r.t. to the inner product defined by Σ_1 . Then, we'll consider the case where they are almost identical.

Orthogonal w_1 and w_k We simplify equation (43),

$$\mathbb{E}_{X, \langle w_1, X \rangle > 0, \langle w_k, X \rangle > 0 | Y=1} [\langle w_k, X \rangle X] = \frac{1}{2\pi} \cdot \frac{1}{\rho_{11}\rho_{kk}} \left[\left(\rho_{kk}^2 \frac{\pi}{2} \right) \Sigma_1 w_1 \right] = \frac{\rho_{kk}}{4\rho_{11}} \Sigma_1 w_1. \quad (49)$$

So, the gradient flow is

$$\tau \frac{d}{dt} w_1 = -\mathbb{E}_{X,Y} \left[\left(\frac{1}{K} \sum_{k \in [K]} \text{ReLU}(\langle w_k, X \rangle) - Y \right) \frac{\partial}{\partial w_1} [\text{ReLU}(\langle w_1, X \rangle)] \right] \quad (50)$$

$$= \frac{1}{2\sqrt{2\pi}} (w_1^\top \Sigma_1 w_1)^{-\frac{1}{2}} \Sigma_1 w_1 - \frac{1}{2K} \left(\Sigma_1 w_1 + \sum_{k>1} \frac{\rho_{kk}}{2\rho_{11}} \Sigma_1 w_1 \right) \quad (51)$$

$$= \frac{1}{2} \left[\frac{1}{\sqrt{2\pi}} (w_1^\top \Sigma_1 w_1)^{-\frac{1}{2}} - \frac{1}{K} \left(1 + \sum_{k>1} \frac{w_k^\top \Sigma_1 w_k}{2w_1^\top \Sigma_1 w_1} \right) \right] \Sigma_1 w_1. \quad (52)$$

So, the update to w_1 only depends on w_k through its weighted norm. There is no additional “competition term” that persists, although of course we've assumed that w_1 and w_k are orthogonal.

Almost identical w_1 and w_k Now, let's consider the case where $\rho_{1k} = \sqrt{\rho_{11}\rho_{kk}} + \varepsilon^2$. Then, we have

$$\mathbb{E}_{X, \langle w_1, X \rangle > 0, \langle w_k, X \rangle > 0 | Y=1} [\langle w_k, X \rangle X] \quad (53)$$

$$\begin{aligned} &\approx \frac{\rho_{kk}^2}{2\pi} \cdot \frac{1}{\varepsilon^2} \left[(\cos^{-1}(-1) - \sin^{-1}(1)) \Sigma_1 w_1 + \rho_{kk}^2 (-\cos^{-1}(-1) + \sin^{-1}(1)) \Sigma_1 w_k \right] \\ &\quad + \frac{1}{2\pi} \cdot \frac{1}{\varepsilon} ((\rho_{kk} - \rho_{1k}) \Sigma_1 w_1 + (\rho_{11} - \rho_{1k}) \Sigma_1 w_k). \end{aligned} \quad (54)$$

Hi!

$$\mathbb{E}_{S>0, T>0|Y=1} [S^2] \quad (55)$$

$$= \int_0^\infty \int_0^\infty s^2 \frac{1}{2\pi \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}(\rho_{11}s^2 - 2\rho_{1k}st + \rho_{kk}t^2)} ds dt \quad (56)$$

$$= \int_0^\infty \frac{1}{\sqrt{2\pi}} \int_0^\infty s^2 \frac{1}{\sqrt{2\pi} \cdot \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} e^{-\frac{1}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}[(\rho_{11}s - \rho_{1k}t)^2 - \rho_{1k}^2t^2 + \rho_{kk}t^2]} ds dt \quad (57)$$

$$= \int_0^\infty \frac{1}{\sqrt{2\pi}\rho_{11}} e^{-\frac{(\rho_{kk} - \rho_{1k}^2)}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}t^2} \underbrace{\int_0^\infty s^2 \frac{\rho_{11}}{\sqrt{2\pi} \cdot \sqrt{\rho_{11}\rho_{kk} - \rho_{1k}^2}} e^{-\frac{\rho_{11}^2}{2(\rho_{11}\rho_{kk} - \rho_{1k}^2)}(s - \frac{\rho_{1k}}{\rho_{11}}t)^2} ds}_{\text{}} \quad (58)$$

To compute this expectation, we need to find the distribution of (S, T) . Let us write $\rho_{ij} = w_i^\top \Sigma_1 w_j$ and $\rho = \frac{\rho_{1k}}{\sqrt{\rho_{11}\rho_{kk}}}$. Then,

$$\begin{bmatrix} S' \\ T' \end{bmatrix} \equiv \begin{bmatrix} S/\sqrt{\rho_{11}} \\ T/\sqrt{\rho_{kk}} \end{bmatrix} | Y = 1 = \begin{bmatrix} w_1^\top / \sqrt{\rho_{11}} \\ w_k^\top / \sqrt{\rho_{kk}} \end{bmatrix} X | Y = 1 \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right). \quad (59)$$

Note that

$$\mathbb{E}_{S>0, T>0|Y=1} [(S')^2] = \mathbb{E}_{S'>0, T'>0|Y=1} [(S')^2] \quad (60)$$

$$= \mathbb{E}_{T'>0|Y=1} [\mathbb{E}_{S'>0|T', Y=1} [(S')^2] \mathbb{P}(S' > 0 | T', Y = 1)]. \quad (61)$$

Note that

$$S' | T' \sim \mathcal{N}(\rho T', 1 - \rho^2). \quad (62)$$

Then, the inner expectation is the second moment of a folded normal. So, we can use Wikipedia:

$$\mathbb{E}_{S'>0|T', Y=1} [(S')^2] = (\rho T')^2 + (1 - \rho^2) = 1 + [(T')^2 - 1]\rho^2. \quad (63)$$

The probability is

$$\mathbb{P}(S' > 0 | T', Y = 1) = \int_0^\infty \frac{1}{\sqrt{2\pi(1 - \rho^2)}} e^{-\frac{1}{2}\left(\frac{s' - \rho T'}{\sqrt{1 - \rho^2}}\right)^2} ds' \quad (64)$$

$$= \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}}\right) \right). \quad (65)$$

So, we have

$$\mathbb{E}_{S>0, T>0|Y=1} [(S')^2] \quad (66)$$

$$= \frac{1}{2} \mathbb{E}_{T'>0|Y=1} \left[(1 + [(T')^2 - 1]\rho^2) \left(1 + \operatorname{erf}\left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}}\right) \right) \right] \quad (67)$$

$$= \frac{1}{2} \mathbb{E}_{T'>0|Y=1} \left[(1 + [(T')^2 - 1]\rho^2) + (1 - \rho^2) \operatorname{erf}\left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}}\right) + \rho^2 (T')^2 \operatorname{erf}\left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}}\right) \right]. \quad (68)$$

We evaluate each of these terms separately. The first can be found using (half) the second moment of a folded normal:

$$\mathbb{E}_{T'>0|Y=1} [1 + [(T')^2 - 1]\rho^2] = 1 + \rho^2 (\mathbb{E}_{T'>0|Y=1} [(T')^2] - 1) = 1 + \rho^2 \left(\frac{1}{2}(1) - 1 \right) \quad (69)$$

$$= 1 - \frac{\rho^2}{2}. \quad (70)$$

The second is

$$(1 - \rho^2) \mathbb{E}_{T' > 0 | Y=1} \left[\operatorname{erf} \left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}} \right) \right] = (1 - \rho^2) \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(T')^2} \operatorname{erf} \left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}} \right) dT' \quad (71)$$

$$= \frac{1 - \rho^2}{2} \left(1 - \frac{2}{\pi} \tan^{-1} \left(\frac{\sqrt{1 - \rho^2}}{\rho} \right) \right). \quad (72)$$

The third is

$$\rho^2 \mathbb{E}_{T' > 0 | Y=1} \left[(T')^2 \operatorname{erf} \left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}} \right) \right] = \rho^2 \int_0^\infty \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(T')^2} (T')^2 \operatorname{erf} \left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}} \right) dT' \quad (73)$$

$$= \frac{\rho^2}{\sqrt{2\pi} \cdot 2\sqrt{\pi}} \left(2^{\frac{3}{2}} \tan^{-1} \left(\frac{\sqrt{1 - \rho^2}}{\rho} \right) - \frac{\rho / \sqrt{2(1 - \rho^2)}}{\frac{1}{2} \left(\frac{\rho^2}{1 - \rho^2} + \frac{1}{2} \right)} \right) \quad (74)$$

$$= \frac{\rho^2}{2\sqrt{2\pi}} \left(2^{\frac{3}{2}} \tan^{-1} \left(\frac{\sqrt{1 - \rho^2}}{\rho} \right) - \frac{2\sqrt{2}\rho\sqrt{1 - \rho^2}}{3\rho^2 + 1} \right). \quad (75)$$

The mean is (?):

$$\mathbb{E}_{S' > 0 | T', Y=1} [S'] = \sqrt{\frac{2}{\pi}} (1 - \rho^2) \cdot e^{-\frac{\rho^2 (T')^2}{2\sqrt{1 - \rho^2}}} + \rho T' \operatorname{erf} \left(\frac{\rho T'}{\sqrt{2(1 - \rho^2)}} \right). \quad (76)$$

Non-Gaussian Data

We needed Gaussianity to say that AX and S were independent. Our construction of v and A was chosen to make AX and S have zero covariance. For general data, this does not imply independence. So, we will have to do something else to compute the conditional expectations.

What do we know about X ? First, it is symmetric about 0, and it is translation invariant. As $g \rightarrow \infty$, $X_i(Z) \xrightarrow{d} \operatorname{sign}(Z_i)$ (perhaps even almost surely, and even if not in distribution, then certainly in probability). I cannot figure out how to compute this for anything larger than 2 dimensions, though I feel it should be possible.

3 Have You Tried Making It Linear?

Gating lets us decompose the ReLU post-activation in terms of the pre-activation's sign and magnitude.

$$\operatorname{ReLU}(\langle w_1(t), x \rangle) = g(t, x) \langle w_1(t), x \rangle \quad \text{where} \quad g(t, x) = \mathbb{1}(\langle w_1(t), x \rangle \geq 0). \quad (77)$$

We generally assume that g does not vary during learning, even though w_1 may. Later on, we'll try to analyze what happens when this does not hold.

To assess the validity of this assumption, we need to see how much $g(x)$, as defined above, changes during learning. Additionally, post-hoc, we can usually pick a somewhat sensible gating structure that mimics a specific run's behavior. But we'd like to be able to determine this gating upfront. We explore all this in the following subsections.

3.1 Sign Flipping

Note that g is invariant to the scale of w_1 . We’ve observed in previous experiments that w_1 appears to grow uniformly in size during much of its training. (There is, importantly, a phase where it goes from Gaussian to non-Gaussian, but the localization seems to be more likely to occur around its mode.) This suggests that $g(x)$ may be relatively constant during learning. If this is so, then it would be reasonable to try using a standard GDLN to model the ReLU network.

We will model the ReLU network as in equation (77), focusing on how g varies with time for each hidden neuron. We will look at the metrics

$$p(t) = \mathbb{P}_x(g(t, x) = g(t + \delta t, x)) \quad \text{for all } t \quad (78)$$

$$p_{\text{unif}} = \mathbb{P}_x(\{g(t, x) = g(t', x) \ \forall t, t'\}) \quad (79)$$

3.2 Predicting Loca(liza)tion

3.3 Evolving Gates?

4 Let’s Consider a Single Layer with Linear Activation...

4.1 Model

Our GDLN model is defined as follows:

$$\hat{y}(x) = \frac{1}{K} \left(\sum_{k \in [K]} g_k(x) w_k^\top \right) x, \quad (80)$$

where g_k are (node) gates, and $w_k \in \mathbb{R}^n$ are the rows of the first-layer weight matrix $W_1 \in \mathbb{R}^{K \times n}$. That is,

$$W_1 = \begin{pmatrix} w_1^\top \\ \vdots \\ w_K^\top \end{pmatrix} \quad (81)$$

4.2 Gradient Flow

Recalling the GDLN paper, the gradient flow for w_1 is given by

$$\tau \frac{d}{dt} w_1^\top = \frac{1}{K} \left[\Sigma^{yx}(p_1) - \sum_{k \in [K]} w_k^\top \Sigma^{xx}(p_1, p_k) \right], \quad (82)$$

where

$$\Sigma^{yx}(p_i) = \langle g_i y x^\top \rangle_{g, x, y} \quad (83)$$

$$\Sigma^{xx}(p_i, p_j) = \langle g_i g_j x x^\top \rangle_{g, x}. \quad (84)$$

4.3 General Case

Let us relabel $b_i = \Sigma^{yx}(p_i)^\top$ and $A_{ij} = \Sigma^{xx}(p_i, p_j)$. Note that A_{ij} is symmetric and $A_{ij} = A_{ji}$. Then, we can write the gradient flow for all weights as

$$K\tau \frac{d}{dt} \underbrace{\begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}}_{w \in \mathbb{R}^{Kn}} = \underbrace{\begin{bmatrix} b_1 \\ \vdots \\ b_K \end{bmatrix}}_{b \in \mathbb{R}^{Kn}} - \underbrace{\begin{bmatrix} A_{11} & \cdots & A_{1K} \\ \vdots & \ddots & \vdots \\ A_{K1} & \cdots & A_{KK} \end{bmatrix}}_{A \in \mathbb{R}^{Kn \times Kn}} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix}. \quad (85)$$

Observe that w is the vectorized form of our $K \times n$ first-layer weight matrix. Note also that A is a symmetric real matrix, so we can diagonalize it as $A = P\Lambda P^\top$, where the columns of P are the eigenvectors of A and the diagonal entries of Λ are the corresponding (nonnegative) eigenvalues. (It is symmetric because A is block symmetric with blocks A_{ij} , and the blocks are also symmetric.) To see this more clearly, let us write

$$\tilde{g}(x) = \begin{bmatrix} g_1(x) \\ \vdots \\ g_K(x) \end{bmatrix}. \quad (86)$$

Then,

$$A = \langle (\tilde{g} \otimes x)(\tilde{g} \otimes x)^\top \rangle_{g,x}, \quad (87)$$

which is clearly a symmetric matrix. (Interjection: Whatever distribution we have over g should satisfy that $\tilde{g} \sim \Pi\tilde{g}$, where Π is some permutation matrix on K elements. That is, the distribution should be invariant to the ordering of the gates, since this is what we want empirically.)

We can reparameterize in terms of $u = P^\top w$ and $c = P^\top b$.

$$K\tau \frac{d}{dt} u = -\Lambda u + c \implies u(t) = \Lambda^{-1} e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + \Lambda^{-1} c, \quad (88)$$

where C is a constant diagonal matrix that defines the initial condition. So,

$$w(t) = P\Lambda^{-1} e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + P\Lambda^{-1} c \quad (89)$$

$$= A^{-1} P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + A^{-1} b. \quad (90)$$

4.4 Winning Gating Structure

Can we read off the winning gating structure from the gradient flow? For simplicity, let us assume we are sampling uniformly from a finite set of G gates, $\{g\}$. Then, we can write equation (85) as

$$K\tau \frac{d}{dt} w = \frac{1}{G} \sum_{g \in \{g\}} [b_{x,y|g} - A_{x,y|g} w], \quad (91)$$

where the subscript on b and A indicates the conditioning on a specific gating structure g . Intuitively, a gating structure that minimizes the norm of A will shrink the slowest. This is somewhat equivalent to minimizing the eigenvalues of A , since they are all nonnegative. (What happens if an eigenvalue is zero?)

Let us consider a single block in A :

$$A_{ij} = \langle g_i g_j x x^\top \rangle_{g,x} = \mathbb{P}(g_i = 1, g_j = 1) \langle x x^\top \rangle_{x|g_i=g_j=1}. \quad (92)$$

Let us quickly ask: Does what we observe empirically match this intuition? We see that receptive fields come in pairs and tile the space.

TODO: empirically look at dominating eigenvalues for finite case!

4.4.1 Early Dynamics

For small t , but sufficiently large to see separation among different eigenvalues, can we predict the leading structure?

4.4.2 Limiting Behavior

If none of the eigenvalues are zero, then $w(\infty) = A^{-1}b$. If we write

$$\tilde{x} = \begin{bmatrix} g_1(x)x \\ \vdots \\ g_K(x)x \end{bmatrix} \in \mathbb{R}^{Kn}, \quad (93)$$

then $A = \langle \tilde{x}\tilde{x}^\top \rangle_{x,y,g}$ and $b = \langle \tilde{x}y \rangle_{x,y,g}$. Then, it is clear that this is the population solution to the OLS problem of regressing y on \tilde{x} , averaging across the distributions of the data *and* the gating architectures.

In this context, one might ask, which gating structure minimizes the MSE loss? The loss is

$$\mathcal{L}_{OLS} = \left\langle \left(\tilde{x}'^\top (\langle \tilde{x}\tilde{x}^\top \rangle_{x,y,g})^{-1} \langle \tilde{x}y \rangle_{x,y,g} - y' \right)^2 \right\rangle_{x',y',g'} \quad (94)$$

$$= \left\langle (\tilde{x}'^\top (\langle \tilde{x}\tilde{x}^\top \rangle_{x,y,g})^{-1} \langle \tilde{x}y \rangle_{x,y,g})^2 - 2(y' \tilde{x}'^\top (\langle \tilde{x}\tilde{x}^\top \rangle_{x,y,g})^{-1} \langle \tilde{x}y \rangle_{x,y,g}) + (y')^2 \right\rangle_{x',y',g'} \quad (95)$$

$$= \frac{1}{2} - \langle y\tilde{x}^\top \rangle_{x,y,g} (\langle \tilde{x}\tilde{x}^\top \rangle_{x,y,g})^{-1} \langle \tilde{x}y \rangle_{x,y,g} \quad (96)$$

$$= \frac{1}{2} - \frac{1}{2} \langle \tilde{x} \rangle_{x,g|y=1}^\top (\langle \tilde{x}\tilde{x}^\top \rangle_{x,g|y=1} + \langle \tilde{x}\tilde{x}^\top \rangle_{x,g|y=0})^{-1} \langle \tilde{x} \rangle_{x,g|y=1}. \quad (97)$$

In the final step, we assumed (WLOG) that the negative class is $y = 0$ and the positive class is $y = 1$. (Throughout, we also assume that the classes are balanced.) The question is: For fixed p_{ξ_1} and p_{ξ_2} , how do we choose the gates g_k to minimize equation (97)?

It may be useful to write this in terms of Kronecker products. Let

$$\tilde{g}(x) = \begin{bmatrix} g_1(x) \\ \vdots \\ g_K(x) \end{bmatrix} \in \{0, 1\}^K. \quad (98)$$

Then, minimizing equation (97) is equivalent to maximizing

$$\mathcal{L}^*(\tilde{g}) = \langle \tilde{g} \otimes x \rangle_{x,g|y=1}^\top (\langle (\tilde{g} \otimes x)(\tilde{g} \otimes x)^\top \rangle_{x,g|y=1} + \langle (\tilde{g} \otimes x)(\tilde{g} \otimes x)^\top \rangle_{x,g|y=0})^{-1} \langle \tilde{g} \otimes x \rangle_{x,g|y=1} \quad (99)$$

over $\tilde{g} : \text{supp}(p_{\xi_1}) \cup \text{supp}(p_{\xi_2}) \rightarrow \{0, 1\}^K$. *I will have to think more about this.*

After a bit more thinking... I think that the best precision matrix would be maximally diagonal (no clue if this is actually true! but maybe it holds empirically?). For Gaussian data (at least), this mean that the blocks are independent conditioned on all the other blocks. Gates that tile the space without overlap would achieve this (I think?). But tbh I haven't got the slightest clue!!

4.5 Exclusive Gates

Let us assume that the gates are exclusive, that is, only one gate is active at a time. Then, $\Sigma^{xx}(p, q) = 0$ for $p \neq q$.

Then A becomes block diagonal. We can write the gradient flow for w_1 as

$$K\tau \frac{d}{dt} w_1 = -A_{11} w_1 + b_1. \quad (100)$$

Note that $A_{11} = \Sigma^{xx}(p_1, p_1)$ is always symmetric (and real). So, we can diagonalize it as $A_{11} = P\Lambda P^\top$, where the columns of P are v_1, \dots, v_n and the diagonal entries of Λ are $\lambda_1, \dots, \lambda_n$. Let us introduce $u_1 = P^\top w_1$ and $c_1 = P^\top b_1$. Then,

$$K\tau \frac{d}{dt} u_1 = -\Lambda u_1 + c_1. \quad (101)$$

This ODE is solved by

$$u_1(t) = \Lambda^{-1} e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + \Lambda^{-1} c_1, \quad (102)$$

where C is a constant diagonal matrix that defines the initial condition. Then,

$$w_1(t) = P \left(\Lambda^{-1} e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + \Lambda^{-1} c_1 \right) \quad (103)$$

$$= P \Lambda^{-1} e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + P \Lambda^{-1} c_1 \quad (104)$$

$$= A_{11}^{-1} \left(P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + b_1 \right). \quad (105)$$

Recalling A_{11} and b_1 ,

$$w_1(t) = (\Sigma^{xx}(p_1, p_1))^{-1} \left(P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + \Sigma^{yx}(p_1)^\top \right). \quad (106)$$

So,

$$w_1(\infty) = (\Sigma^{xx}(p_1, p_1))^{-1} \left(P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + \Sigma^{yx}(p_1)^\top \right). \quad (107)$$

As with above, this is the population solution to OLS, $(X^\top X)^{-1} X^\top y = (\langle xx^\top \rangle)^{-1} (\langle xy \rangle)$.

So, each weight matrix converges to the OLS solution on the subset of the data determined by its gate.

4.6 Redundant Gates

What if $g_1 = g_2$? Then, $b_1 = b_2$ and $A_{11} = A_{12} = A_{22}$. So,

$$K\tau \frac{d}{dt} w_1 = -A_{11}(w_1 + w_2) + b_1, \quad (108)$$

$$K\tau \frac{d}{dt} w_2 = -A_{11}(w_1 + w_2) + b_1. \quad (109)$$

Clearly, then, $w_1 - w_2$ is a constant vector. Moreover, $\frac{1}{2}(w_1 + w_2)$ evolves according to

$$\frac{K\tau}{2} \frac{d}{dt} (w_1 + w_2) = -A_{11}(w_1 + w_2) + b_1. \quad (110)$$

Writing $2\Delta = w_1 - w_2$ and $w_1 + w_2 = 2(w_1 - \Delta)$, we have

$$K\tau \frac{d}{dt} w_1 = K\tau \frac{d}{dt} (w_1 - \Delta) = -2A_{11}(w_1 - \Delta) + b_1. \quad (111)$$

We can plug this into our solution from the previous section to get

$$w_1(t) = \frac{1}{2} A_{11}^{-1} \left(P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + b_1 \right). \quad (112)$$

So,

$$w_2(t) = \frac{1}{2} A_{11}^{-1} \left(P e^{-\frac{t}{K\tau} \Lambda + C} \mathbf{1} + b_1 \right) - (w_1(0) - w_2(0)). \quad (113)$$

5 Theory-driven Experiments

5.1 General Case

5.2 Single Gate

5.3 Redundant Gates

5.4 Simple ReLU Network

Let's consider what happens in a single step of a network with ReLU activation. We make predictions with

$$\hat{y}(x) = \frac{1}{K} \sum_{k \in [K]} \text{ReLU}(\langle w_k, x \rangle). \quad (114)$$

We consider MSE loss,

$$\mathcal{L} = \frac{1}{2} \mathbb{E}_{X,Y} \left[(\hat{y}(X) - Y)^2 \right]. \quad (115)$$

The gradient flow for w_1 is given by

$$\tau \frac{d}{dt} w_1 = - \mathbb{E}_{X,Y} \left[(\text{ReLU}(\langle w_1(t), X \rangle) - Y) \frac{\partial}{\partial w_1} [\text{ReLU}(\langle w_1(t), X \rangle)] \right] \quad (116)$$

$$= \underbrace{\left(\mathbb{E}_{X,Y|\langle w_1(t), X \rangle > 0} [YX] - \mathbb{E}_{X,Y|\langle w_1(t), X \rangle > 0} [\langle w_1(t), X \rangle X] \right)}_{\equiv (I)} \mathbb{P}(\langle w_1(t), X \rangle > 0). \quad (117)$$

Recall that X is a mixture of $X | Y = 1$ and $X | Y = 0$. So, we will compute these expectations separately. Let us write $S = \langle w_1(t), X \rangle$. Then, we can use the law of total expectation to write

$$\mathbb{E}_{X|Y=1, \langle w_1(t), X \rangle > 0} [f(X)] = \mathbb{E}_{S>0|Y=1} [\mathbb{E}_{X|S,y=1} [f(X)]] . \quad (118)$$

So, we need to find the distribution of X conditioned on $S \equiv \langle w_1(t), X \rangle = s$. In general, this is very challenging. We will split this into two terms, one of which disappears when X is Gaussian. Let Σ be the covariance of x (recall it has mean 0). We write

$$X = AX + Sv, \quad (119)$$

where

$$v = \frac{1}{w_1(t)^\top \Sigma w_1(t)} \Sigma w_1(t), \quad (120)$$

$$A = I_n - v w_1(t)^\top. \quad (121)$$

Thus, $\mathbb{E}_{X|S=s}[X] = \mathbb{E}_{X|S=s}[AX] + sv$. Our choice of A and v implies that AX and S have zero *covariance* (see [this post](#)). When X is Gaussian, this implies that AX and S are independent, so we'd have $\mathbb{E}_{X|S=s}[X] = A \mathbb{E}[X] + sv = sv$.

With this representation,

$$\mathbb{E}_{X|Y=1, S>0} [YX] = \mathbb{E}_{S>0|Y=1} [\mathbb{E}_{X|S,Y=1} [AX] + Sv] \quad (122)$$

$$= \mathbb{E}_{X|Y=1, S>0} [AX] + \frac{\mathbb{E}_{S>0|Y=1} [S]}{w_1(t)^\top \Sigma w_1(t)} \Sigma w_1(t), \quad (123)$$

$$\mathbb{E}_{X|Y=1, S>0} [SX] = \mathbb{E}_{S>0|Y=1} [S \mathbb{E}_{X|S,Y=1} [AX] + S^2 v] \quad (124)$$

$$= \mathbb{E}_{X|Y=1, S>0} [SAX] + \mathbb{E}_{S>0|Y=1} [S^2] v \quad (125)$$

$$= \mathbb{E}_{X|Y=1, S>0} [SAX] + w_1(t)^\top \Sigma w_1(t) v, \quad (126)$$

$$= \mathbb{E}_{X|Y=1, S>0} [SAX] + \Sigma w_1(t). \quad (127)$$

Note that if X were Gaussian, then the first terms in equations (123) and (127) would be zero.

Now, we evaluate (I) and (II).

$$(I) = \mathbb{E}_{X|S>0} [SAX] + [\mathbb{P}(Y = 1 | S > 0)\Sigma_1 + \mathbb{P}(Y = 0 | S > 0)\Sigma_0] w_1(t), \quad (128)$$

$$(II) = \mathbb{E}_{X|S>0} [AX] + \mathbb{P}(Y = 1 | S > 0) \frac{\mathbb{E}_{S>0|Y=1} [S]}{w_1(t)^\top \Sigma_1 w_1(t)} \Sigma_1 w_1(t) \quad (129)$$

Then,

$$[(I) - (II)] \mathbb{P}(S > 0) \quad (130)$$

$$= -\mathbb{E}_{X|S>0} [(S-1)AX] - \left[\mathbb{P}(Y = 1 | S > 0) \left(\frac{\mathbb{E}_{S>0|Y=1} [S]}{w_1(t)^\top \Sigma_1 w_1(t)} - 1 \right) \Sigma_1 + \mathbb{P}(Y = 0 | S > 0) \frac{\mathbb{E}_{S>0|Y=0} [S]}{w_1(t)^\top \Sigma_0 w_1(t)} \Sigma_0 \right] w_1(t). \quad (131)$$

By symmetry, $\mathbb{P}(Y = 1 | S > 0) = \frac{1}{2}$ and $PR(S > 0) = \frac{1}{2}$. Then,

$$4\tau \frac{d}{dt} w_1 = -\mathbb{E}_{X|S>0} [(S-1)AX] - \left[\left(\frac{\mathbb{E}_{S>0|Y=1} [S]}{w_1(t)^\top \Sigma_1 w_1(t)} - 1 \right) \Sigma_1 + \frac{\mathbb{E}_{S>0|Y=0} [S]}{w_1(t)^\top \Sigma_0 w_1(t)} \Sigma_0 \right] w_1(t) \quad (132)$$

Also recall that Σ_1 and Σ_0 both diagonalize in the discrete Fourier basis, which we denote with P , and their corresponding diagonal matrices of eigenvalues Λ_1 and Λ_0 . Write $u_1 = P^\top w_1$.

$$4\tau \frac{d}{dt} u_1 = -\mathbb{E}_{X|S>0} [(S-1)P^\top AX] - \left[\left(\frac{\mathbb{E}_{S>0|Y=1} [S]}{u_1(t)^\top \Lambda_1 u_1(t)} - 1 \right) \Lambda_1 + \frac{\mathbb{E}_{S>0|Y=0} [S]}{u_1(t)^\top \Lambda_0 u_1(t)} \Lambda_0 \right] u_1(t). \quad (133)$$

Let us expand the first term for $Y = 1$. Define $\Xi = P^\top X$.

$$\mathbb{E}_{X|Y=1, S>0} [(S-1)P^\top AX] = \mathbb{E}_{X|Y=1, S>0} \left[(S-1)P^\top \left(I_n - \frac{\Sigma_1 w_1(t) w_1(t)^\top}{w_1(t)^\top \Sigma_1 w_1(t)} \right) X \right] \quad (134)$$

$$= \mathbb{E}_{\Xi|Y=1, \langle u_1(t), \Xi \rangle > 0} \left[(\langle u_1(t), \Xi \rangle - 1) \left(I_n - \frac{\Lambda_1 u_1(t) u_1(t)^\top}{u_1(t)^\top \Lambda_1 u_1(t)} \right) \Xi \right]. \quad (135)$$

Now, we must stop and ask: what is Ξ ? Recall that P is the discrete Fourier basis. Importantly, it is actually the *real* part of the discrete Fourier basis since Σ_1 is symmetric. That is, with $\omega = e^{-\frac{2\pi i}{n}}$,

$$P_{:,j} = \Re \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ \omega^j \\ \omega^{2j} \\ \vdots \\ \omega^{(n-1)j} \end{bmatrix} = \frac{1}{\sqrt{n}} \begin{bmatrix} 1 \\ \cos(\frac{2\pi}{n}j) \\ \cos(\frac{2\pi}{n}2j) \\ \vdots \\ \cos(\frac{2\pi}{n}(n-1)j) \end{bmatrix}. \quad (136)$$

Ξ is the discrete Fourier transform (DFT) of X .

TRY COMPUTING DENSITY OF NON-GAUSSIAN USING TRICK; DON'T THINK FOURIER APPROACH WILL BE VERY HELPFUL TBH.

5.4.1 Gaussian X

Now, let us assume that X is Gaussian. Then, $\mathbb{E}_{X|S>0} [(S-1)P^\top AX] = 0$. Furthermore,

$$\mathbb{E}_{S>0|Y=1} [S] = \sqrt{\frac{2}{\pi}} (w_1(t)^\top \Sigma_1 w_1(t))^\frac{1}{2} = \sqrt{\frac{2}{\pi}} (u_1(t)^\top \Lambda_1 u_1(t))^\frac{1}{2}. \quad (137)$$

Then, the gradient flow becomes

$$4\tau \frac{d}{dt} u_1 = - \left[\left(\frac{1}{\sqrt{u_1(t)^\top \Lambda_1 u_1(t)}} - 1 \right) \Lambda_1 + \frac{1}{\sqrt{u_1(t)^\top \Lambda_0 u_1(t)}} \Lambda_0 \right] u_1(t). \quad (138)$$

So,

$$\mathbb{E}_{S>0|y=1} [\mathbb{E}_{x|S,y=1} [x]] = \mathbb{E}_{S>0|y=1} [sv_1] = \mathbb{E}_{S>0|y=1} [s] v_1, \quad (139)$$

$$\mathbb{E}_{S>0|y=1} [\mathbb{E}_{x|S,y=1} [xx^\top]] = \mathbb{E}_{S>0|y=1} [A_1 \Sigma_1 A_1^\top] = A_1 \Sigma_1 A_1^\top. \quad (140)$$

Note that $S \sim \mathcal{N}(0, w_1(t)^\top \Sigma w_1(t))$. So, $\mathbb{E}_{S>0|y=1} [s] = \left(\frac{2}{\pi} w_1(t)^\top \Sigma w_1(t)\right)^{\frac{1}{2}}$. In summary,

$$\mathbb{E}_{x|y=1, \langle w_1(t), x \rangle > 0} [x] = \frac{\left(\frac{2}{\pi} w_1(t)^\top \Sigma w_1(t)\right)^{\frac{1}{2}}}{w_1(t)^\top \Sigma w_1(t)} \Sigma w_1(t) = \sqrt{\frac{2}{\pi}} (w_1(t)^\top \Sigma w_1(t))^{-\frac{1}{2}} \Sigma w_1(t), \quad (141)$$

$$\mathbb{E}_{x|y=1, \langle w_1(t), x \rangle > 0} [xx^\top] = A_1 \Sigma_1 A_1^\top = \left(I_n - \frac{1}{w_1(t)^\top \Sigma_1 w_1(t)} \Sigma_1 w_1(t) w_1(t)^\top \right) \Sigma_1. \quad (142)$$

Rewriting the gradient flow,

$$\frac{\tau}{\mathbb{P}(\langle w_1(t), x \rangle > 0)} \frac{d}{dt} w_1 \quad (143)$$

$$= \left[\left(I_n - \frac{1}{w_1(t)^\top \Sigma_1 w_1(t)} \Sigma_1 w_1(t) w_1(t)^\top \right) \Sigma_1 + \left(I_n - \frac{1}{w_1(t)^\top \Sigma_0 w_1(t)} \Sigma_0 w_1(t) w_1(t)^\top \right) \Sigma_0 \right] w_1(t) - \left[\sqrt{\frac{2}{\pi}} (w_1(t)^\top \Sigma w_1(t))^{-\frac{1}{2}} \Sigma w_1(t) \right] \quad (144)$$

$$= -\sqrt{\frac{2}{\pi}} (w_1(t)^\top \Sigma_1 w_1(t))^{-\frac{1}{2}} \Sigma_1 w_1(t). \quad (145)$$

By symmetry, $\mathbb{P}(\langle w_1(t), x \rangle > 0) = \frac{1}{2}$.

$$\tau \frac{d}{dt} w_1 = -2 \sqrt{\frac{2}{\pi}} [w_1(t)^\top \Sigma_1 w_1(t)]^{-\frac{1}{2}} \Sigma_1 w_1(t). \quad (146)$$

Let us write $\Sigma_1 = P \Lambda P^\top$, where Λ is diagonal and P is orthogonal, and $u_1 = P^\top w_1$. Then,

$$\tau \frac{d}{dt} u_1 = -2 \sqrt{\frac{2}{\pi}} [u_1(t)^\top \Lambda u_1(t)]^{-\frac{1}{2}} \Lambda u_1(t). \quad (147)$$

So it appears that this shrinks u_1 to 0, and this is accelerated as u_1 gets small by the weighted norm. This is consistent with what we observe for the Gaussian case. But what happens when x is non-Gaussian. Analytically, it's hard to say exactly. But it seems like the input-input covariance terms will *not* cancel, and so we will have another term that hopefully does more than just shrink u_1 to 0. It is surprising that $w_1(t)$ pops out of the input-output term. This *does not happen in the gated linear network*, and this seems like a key difference. What does this say about gating early during training? I should double-check this to make sure it's right. If this weren't the case though, we would get a bias term that probably yields localization (or some nonzero structure).

What we see in equation (147) is that the largest eigenvalues are shrunk fastest. This corresponds to the longest-frequency signals disappearing quickly. So, we see that the Gaussian noise quickly turns into a short-range oscillation, which is then damped out to 0. This seems consistent with my ReLU simulations, but I need to make sure the results are perfectly comparable (same learning rate, using exact same data at each time step, etc.)