

1 Statement

The goal of this analysis is to understand the shared representations model.

We have words in a dictionary:

$$\mathcal{D} \triangleq \{\mathbf{x}_1, \dots, \mathbf{x}_V\} \subseteq \mathbb{R}^D,$$

and a set of L languages, indexed by $[L] = \{1, \dots, L\}$. For simplicity, we make the following assumption:

Assumption 1. *The dictionary \mathcal{D} is a set of one-hot vectors that span \mathbb{R}^D . So, $V = D$.*

We train a gated deep linear network (GDLN) to “translate” words from one language into another. The gated deep linear network is represented by a collection of L input and output weight matrices of dimensions $D \times H$ and $H \times D$, respectively, and a hidden layer of size $H \times H$. The forward pass of the GDLN is given by:

$$\hat{y}(\mathbf{x}) = \sum_{l' \in [L]} g_{l'}(\mathbf{x}) \mathbf{W}_o^{l'} \mathbf{W}_h \left(\sum_{l \in [L]} g_l(\mathbf{x}) \mathbf{W}_i^l \mathbf{x} \right) = \sum_{l, l' \in [L]} g_l(\mathbf{x}) g_{l'}(\mathbf{x}) \mathbf{W}_o^{l'} \mathbf{W}_h \mathbf{W}_i^l \mathbf{x}.$$

Here, the input \mathbf{x} contains a one-hot vector for the word we are translating, and an input-output language pair. To make this more explicit, we will now write inputs as (\mathbf{x}, p) , where \mathbf{x} is the one-hot word and $p = (k, k')$ is the language pair, where $k, k' \in [L]$ are the languages of the input and output, respectively. Thus, we write,

$$\hat{y}(\mathbf{x}, p) = \underbrace{\left[\sum_{l, l' \in [L]} g_l(p) g_{l'}(p) \mathbf{W}_o^{l'} \mathbf{W}_h \mathbf{W}_i^l \right]}_{\triangleq \mathbf{A}_p} \mathbf{x},$$

where \mathbf{A}_p is a matrix of size $V \times V$ that depends only on the input via the language pair p .

During training, we will sample input-output pairs p uniformly from a subset of the language pairs, which we denote $T \subseteq [L]^2$. We will equivalently use T to denote the cardinality of the set T , and \mathbf{T} to denote a binary matrix, where $\mathbf{T}_{k, k'} = 1$ if $(k, k') \in T$ and 0 otherwise. The gates operate on p as follows: $g_l(p) = \mathbb{1}(l = k)$ and $g_{l'}(p) = \mathbb{1}(l' = k')$.

We study *full-batch gradient flow*. For simplicity, we will assume that $y(\mathbf{x}_v, p) \equiv \mathbf{x}_{\pi_p(v)}$ where π_p is a permutation on $[L]$ corresponding to language pair p , i.e., the target is another one-hot vector in the dictionary \mathcal{D} ¹. In particular, we assume $\pi_p = \pi_{k'} \circ \pi_k^{-1}$, where π_k is a permutation on the input language k and $\pi_{k'}$ is a permutation on the output language k' . Intuitively, this maps each word in the input language to a word to an abstract representation, and then maps that abstract representation to a word in the output language. Since these are one-hot vectors (by assumption 1), we can also write $\mathbf{x}_{\pi_p(v)} = \mathbf{\Pi}_p \mathbf{x}_v = \mathbf{\Pi}_{k'} \mathbf{\Pi}_k^\top \mathbf{x}_v$, where $\mathbf{\Pi}_k$ is the permutation matrix corresponding to π_k .

¹I believe this is how Andrew set it up.

With this structure, we can write the loss function as:

$$\begin{aligned}
\mathcal{L} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, p \sim T} \left[\frac{1}{2} \|\hat{y}(\mathbf{x}, p) - y(\mathbf{x}, p)\|^2 \right] \\
&= \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, p \sim T} \left[\|\mathbf{A}_p \mathbf{x} - \mathbf{\Pi}_p \mathbf{x}\|^2 \right] \\
&= \frac{1}{2} \mathbb{E}_{p \sim T} \left[\frac{1}{D} \sum_{i \in [D]} \left[\|(\mathbf{A}_p - \mathbf{\Pi}_p) \mathbf{e}_i\|^2 \right] \right] \\
&= \frac{1}{2D} \mathbb{E}_{p \sim T} \left[\|\mathbf{A}_p - \mathbf{\Pi}_p\|_F^2 \right].
\end{aligned}$$

Expanding \mathbf{A}_p and summing over p gives:

$$\mathcal{L} = \frac{1}{2DT} \sum_{(k, k') \in T} \left\| \sum_{l, l' \in [L]} g_l(p) g_{l'}(p) \mathbf{W}_o^{l'} \mathbf{W}_h \mathbf{W}_i^l - \mathbf{\Pi}_{k'} \mathbf{\Pi}_k^\top \right\|_F^2.$$

We use block matrix notation to clean this up; define:

$$\begin{aligned}
\mathbf{W}_o &= \begin{pmatrix} \mathbf{W}_o^1 \\ \vdots \\ \mathbf{W}_o^L \end{pmatrix} \in \mathbb{R}^{LH \times D}, \\
\mathbf{W}_i &= \begin{pmatrix} \mathbf{W}_i^1 & \dots & \mathbf{W}_i^L \end{pmatrix} \in \mathbb{R}^{D \times LH}, \\
\mathbf{\Pi} &= \begin{pmatrix} \mathbf{\Pi}_1 \\ \vdots \\ \mathbf{\Pi}_L \end{pmatrix} \in \mathbb{R}^{LD \times D},
\end{aligned}$$

Then, using \mathbf{J}_L to denote the $L \times L$ matrix of all ones, we can write:

$$\mathcal{L} = \frac{1}{2DT} \left\| (\mathbf{W}_o \mathbf{W}_h \mathbf{W}_i - \mathbf{\Pi} \mathbf{\Pi}^\top) \odot (\mathbf{T} \otimes \mathbf{J}_L) \right\|_F^2,$$

where \otimes is the Kronecker product and \odot is the Hadamard product.

I believe Andrew said we could make the following, additional, assumption:

Assumption 2. *The permutation matrices $\mathbf{\Pi}_k$ can all be taken to be the identity matrix, i.e. $\mathbf{\Pi}_k = \mathbf{I}_D$ for all $k \in [L]$.*

Under assumption 2, the loss function simplifies to:

$$\mathcal{L} = \frac{1}{2DT} \left\| (\mathbf{W}_o \mathbf{W}_h \mathbf{W}_i - \mathbf{J}_L \otimes \mathbf{I}_L) \odot (\mathbf{T} \otimes \mathbf{J}_L) \right\|_F^2. \tag{1}$$

This is the loss function we simulate below. It is worth noting that we would have obtained the same loss function if we had used $\mathbf{x} \sim \mathcal{N}(0, \frac{1}{D} \mathbf{I}_D)$ instead of one-hot vectors.

2 Simulations

To validate and gain insight our model, we simulate the loss function in equation (1). We sample \mathbf{T} by fixing the diagonal to be 1 and randomly making off-diagonal entries 1 until we achieve a desired fraction of 1s;

all other entries are set to 0^2 . We use JAX to run these experiments, giving us much greater control over randomness and initialization. In particular, for each seed, when increasing the number of 1s by one over a previous experiment, we do not resample the entire matrix \mathbf{T} , but rather we only sample one new entry. Additionally, when change the initialization scale σ , we do not resample the entire matrix $\mathbf{W}_o\mathbf{W}_h\mathbf{W}_i$, but rather we only rescale the entries of those matrices.

Figure 1 shows the results of our simulations. We see that generalization typically emerges after roughly 30% of the pairs have been trained on. However, this is not always the case, and we do see spikes in test loss for some experiments after this point (my initial guess is that this is numerical, but I could also believe that these simulations are highly sensitive to the structure of \mathbf{T} , though I believe Andrew said he did not see this). Additionally, the generalization threshold appears to *decrease* as the initialization scale σ decreases, which suggests that, as $\sigma \rightarrow 0$, the generalization threshold may go to zero as well, a counterintuitive result. Moreover, we see considerable variation across random seeds.

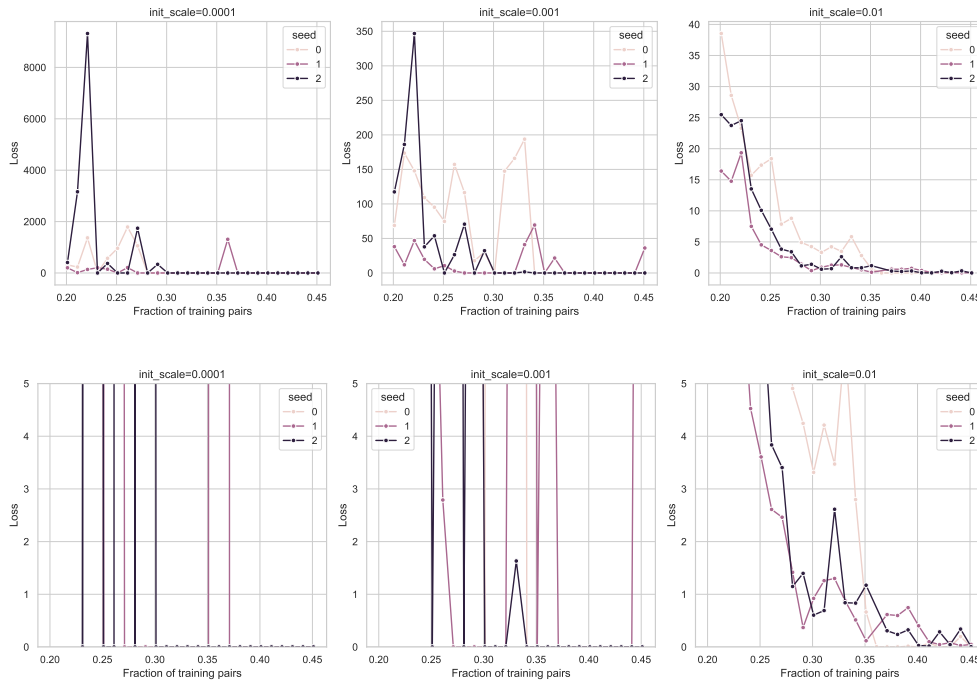


Figure 1: (a) MSE loss on *untrained* language pairs after gradient flow achieving a value of 10^{-6} on the loss \mathcal{L} from equation (1) for various σ and three random seeds. (b) Same setup as in (a), but the loss is truncated at $\mathcal{L} = 5$ to show spikes later in training. The model is a GDLN with $L = 10$, $H = 32$, and $D = 32$. Each point represents a single training run.

²I don't see any inherent reason to make the diagonal all 1s; this is something I would change, but I had sunk too much compute into these sims by the time I realized this.

3 A Simple Case

The block structure of the loss suggests we may be able to gain insight from studying a simpler, non-block, analogue:

$$\mathcal{L} = \frac{1}{2} \|(\mathbf{u} \mathbf{v}^\top - \mathbf{J}_n) \odot \mathbf{T}\|_F^2. \quad (2)$$

Here, $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ and $s \in \mathbb{R}$ are learnable parameters and $\mathbf{T} \in \{0, 1\}^{n \times n}$ is a binary matrix. We would obtain a similar expression if we diagonalized the matrices in the loss function in equation (1).

The gradients of the loss function are given by:

$$\nabla_{\mathbf{u}} \mathcal{L} = s \mathbf{E} \mathbf{v}, \quad \nabla_s \mathcal{L} = \langle \mathbf{E}, \mathbf{u} \mathbf{v}^\top \rangle_F, \quad \nabla_{\mathbf{v}} \mathcal{L} = s \mathbf{E}^\top \mathbf{u},$$

where $\mathbf{E} \triangleq (\mathbf{u} \mathbf{v}^\top - \mathbf{J}_L) \odot \mathbf{T}$ is the “error” matrix. This implies that gradient flow (under unit learning rate) is:

$$\frac{d\mathbf{u}}{dt} = -s \mathbf{E} \mathbf{v}, \quad \frac{ds}{dt} = -\langle \mathbf{E}, \mathbf{u} \mathbf{v}^\top \rangle_F, \quad \frac{d\mathbf{v}}{dt} = -s \mathbf{E}^\top \mathbf{u}.$$

Consider the behavior of the model when σ , the initialization scale, is small. Then, we have $\mathbf{E} \approx -\mathbf{J}_L \odot \mathbf{T} = \mathbf{T}$. So,

$$\frac{d\mathbf{u}}{dt} \approx s \mathbf{T} \mathbf{v}, \quad \frac{ds}{dt} \approx \langle \mathbf{T}, \mathbf{u} \mathbf{v}^\top \rangle_F, \quad \frac{d\mathbf{v}}{dt} \approx s \mathbf{T}^\top \mathbf{u}.$$

Consider the SVD of $\mathbf{T} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$. We write

$$\frac{d\mathbf{u}}{dt} \approx s \mathbf{T} \mathbf{v} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top \mathbf{v} \implies \frac{d\mathbf{U}^\top \mathbf{u}}{dt} \approx \mathbf{\Sigma} \mathbf{V}^\top \mathbf{v}.$$

Similarly,

$$\frac{d\mathbf{v}}{dt} \approx s \mathbf{T}^\top \mathbf{u} = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^\top \mathbf{u} \implies \frac{d\mathbf{V}^\top \mathbf{v}}{dt} \approx \mathbf{\Sigma} \mathbf{U}^\top \mathbf{u}.$$

Defining $\tilde{\mathbf{u}} \triangleq \mathbf{U}^\top \mathbf{u}$ and $\tilde{\mathbf{v}} \triangleq \mathbf{V}^\top \mathbf{v}$, we have:

$$\begin{aligned} \frac{d\tilde{\mathbf{u}}}{dt} &\approx s \mathbf{\Sigma} \tilde{\mathbf{v}}, \\ \frac{ds}{dt} &\approx \langle \mathbf{T}, \mathbf{u} \mathbf{v}^\top \rangle_F = \langle \mathbf{\Sigma}, \tilde{\mathbf{u}} \tilde{\mathbf{v}}^\top \rangle_F, \\ \frac{d\tilde{\mathbf{v}}}{dt} &\approx s \mathbf{\Sigma} \tilde{\mathbf{u}}. \end{aligned}$$

Focusing just on $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{v}}$, we have:

$$\begin{aligned} \tilde{\mathbf{u}}(t) &= \cosh(\mathbf{\Sigma} S(t)) \tilde{\mathbf{u}}(0) + \sinh(\mathbf{\Sigma} S(t)) \tilde{\mathbf{v}}(0), \\ \tilde{\mathbf{v}}(t) &= \sinh(\mathbf{\Sigma} S(t)) \tilde{\mathbf{u}}(0) + \cosh(\mathbf{\Sigma} S(t)) \tilde{\mathbf{v}}(0), \end{aligned}$$

where $S(t) \triangleq \int_0^t s(t') dt'$ and functions f are applied to matrices \mathbf{X} via their singular values, i.e. $f(\mathbf{X}) = \mathbf{U} f(\mathbf{\Sigma}) \mathbf{V}^\top$; so, in this case, we only apply it along the diagonal of $\mathbf{\Sigma}$. Thus, assuming $S(t) \rightarrow \infty$ as $t \rightarrow \infty$ (which clearly must be the case), we see that $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{v}}$ will align, since $\cosh(x) \sim \sinh(x) \sim e^x$ for large x .³

³I expect this happens around time $t \sim -\log(\sigma)$ as $\sigma \rightarrow 0$.

Moreover, it is known that for \mathbf{T} with entries i.i.d. $\text{Bern}(p)$, there is an outlier singular value on the order of pn , while the rest follow the Marchenko-Pastur law, which concentrate in the interval $[0, 2\sqrt{p(1-p)}]$. Assuming we have a \mathbf{T} of this form, the largest singular value will dominate the dynamics of the system *by a lot*. In fact, as σ is sufficiently small, we will always see that the largest singular value dominates because the hyperbolic functions grow exponentially. Therefore, $\tilde{\mathbf{u}} \approx c_u \mathbf{e}_1$, $\tilde{\mathbf{v}} \approx c_v \mathbf{e}_1$, which implies $\mathbf{u} \approx c_u \mathbf{u}_*$ and $\mathbf{v} \approx c_v \mathbf{v}_*$, where \mathbf{u}_* , \mathbf{v}_* are the first columns of \mathbf{U} , \mathbf{V} , respectively, i.e. the first left and right singular vectors of \mathbf{T} .

This means that $\mathbf{u}\mathbf{v}^\top$ emerges to be the best rank-1 approximation of \mathbf{T} , in Frobenius norm, up to a scaling factor. Note that if $\mathbf{u} = \mathbf{v} = \mathbf{1}$, the vector of all 1s, then we immediately obtain perfect generalization. This would happen, for instance, if \mathbf{T} is a non-identity circulant matrix.

(So, does this mean any such circulant \mathbf{T} *should* generalize perfectly, asymptotically as $\sigma \rightarrow 0$? This seems surprising ...)

I tried to test this observation, and it generally seems to be true, but I can't make the initialization scale small enough to get it to always work. It seems like the gap between the largest singular value and the rest needs to be sufficiently large, but mathematically, I don't see why this should be the case. As the circulant width increases, this gap appears to grow, which supports why we only see this for wider circulant matrices.