
How Cross-Entropy Shapes Representation Geometry: A Spectral Study on Cycle Graphs

Anonymous Authors¹

Abstract

Why do embeddings trained on graph transition matrices remember geometry? On cycles, we show that this phenomenon is neither surprising nor mysterious: it admits a transparent spectral explanation in terms of the global optima and implicit bias of softmax cross-entropy. We make this mechanism explicit by characterizing the optimal embedding geometry of tied and untied parameterizations trained on cycle graphs, under both sparse targets and dense label-smoothed targets. For sparse targets, tied embeddings converge to a finite rank-2 solution whose embeddings recover the cycle structure, whereas untied embeddings diverge along a max-margin direction whose rank scales with the number of nodes and whose spectrum differs from that of the target. With label smoothing, untied embeddings recover the target structure exactly up to scaling, while tied embeddings approximate its positive semidefinite component. Together, our theoretical and numerical results show that geometric memory on cycles is a consequence of how cross-entropy transforms the spectral structure of the target distribution into representation geometry through its implicit optimization bias.

1 Introduction

Classical word embedding models such as word2vec and GloVe learn representations from word-context statistics (Mikolov et al., 2013; Pennington et al., 2014), and an extensive line of work relates these embeddings to implicit matrix factorization, PMI-type statistics, and the spectral structure in the data (Levy & Goldberg, 2014; Gittens et al., 2017; Allen & Hospedales, 2019; Qiu et al., 2018). More

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the FoGen Workshop at ICML 2026. Do not distribute.

recently, a growing body of works studies how different objectives form the geometry of embeddings in supervised and contrastive training settings (Papayan et al., 2020; Zhu et al., 2021; Jiang et al., 2023; Li et al., 2023; Bangachev et al., 2025). This geometric perspective has also become important in language models, where the representation structure is studied to probe what models store, how semantic or relational information is organized, and what geometry is induced by next-token prediction (Saxe et al., 2019; Elhage et al., 2022; Gurnee & Tegmark, 2023; Park et al., 2023; Jiang et al., 2024; Lee et al., 2025; Li et al., 2025; Zhao et al., 2024; Zhao & Thramoulidis, 2025).

Recently, Noroozizadeh et al. (2025) empirically showed that when transformers with tied embeddings are trained on pairs (u, v) of neighboring nodes from a graph G , the learned node embeddings reflect global structure in the graph rather than merely encoding local adjacency information, a phenomenon they termed as *geometric memory*. They further observe that the dominant directions of this learned embedding geometry, as seen through 3D PCA or UMAP (McInnes et al., 2018) visualizations, are closely aligned with those learned by a one-layer *node2vec* model (Grover & Leskovec, 2016). In this model, the node embeddings $\mathbf{W} \in \mathbb{R}^{n \times d}$ are trained directly to minimize the cross-entropy (CE) between the *sparse* random-walk transition probabilities of G and the softmax probabilities induced by the logits $\mathbf{W}\mathbf{W}^\top$ (see Sec. 2).

A particularly striking feature of the *node2vec* model pointed in Noroozizadeh et al. (2025) is that, even without an explicit dimensional bottleneck or regularization, it empirically converges to low-rank embeddings. Since the tied parameterization constrains the logit matrix $\mathbf{W}\mathbf{W}^\top$ to be positive semidefinite (PSD), one should already expect its spectrum to have lower rank than that of the graph adjacency matrix, whose eigenvalues can be both positive and negative. However, the observed rank collapse is much stronger than what the PSD constraint alone would suggest. For example, on the cycle graph in Fig. 1, with n nodes, the learned tied solution always has rank 2, whereas the positive spectral part of the adjacency matrix has rank $\Theta(n)$. By contrast, if one relaxes the PSD constraint and parameterizes the logits using untied embeddings as $\mathbf{W}\mathbf{H}$, equivalently

the unconstrained features model in (Zhao et al., 2024), the learned geometry changes substantially: the resulting solution is no longer low-rank. Both the rank of the full matrix and the rank of its PSD component grow as $\Theta(n)$.

On the other hand, as shown in the bottom row of Fig. 1, if we smooth the transition probabilities of the graph, thereby removing sparsity while preserving the graph structure, the strict low-rank bias under tied embeddings disappears. In this case, the eigenspectrum of the learned embeddings becomes a much closer approximation to the positive component of the target probabilities.

Motivated by these observations, we focus on the cycle graph and analytically study the optimal embedding geometry learned under the CE loss, to characterize, theoretically or numerically, how the learned embeddings are affected by the sparsity of the target labels and the tying of the embedding layers.

After introducing the notation and problem in Sec. 2, along with preliminary discussions in Secs. 3 and 4, we characterize the optimal embeddings of the loss function for sparse and dense target probabilities in Sec. 5. Our analysis shows that, while label smoothing primarily rescales the eigenspectrum of the target label matrix, it can substantially alter the eigenspectrum of the learned embeddings in both tied and untied settings. Specifically: (1) With untied embeddings and label-smoothed targets, the learned logits recover the target spectrum up to a scaling factor. In contrast, with sparse targets, although the optimal solution still interpolates the target probabilities, as predicted by the implicit bias of the loss, it converges to a markedly different spectrum. (2) With tied embeddings, label smoothing again drives the learned spectrum closer to the PSD component of the target spectrum, and this approximation improves as the smoothing strength increases. Without label smoothing, however, the tied solution exhibits a fixed low-rank spectrum, retaining only the two dominant eigendirections of the adjacency matrix.

2 Setup

Let $G = (V, E)$ be a graph with n nodes and $\mathbf{A} \in \{0, 1\}^{n \times n}$ be the adjacency matrix with $\mathbf{D} = \text{diag}(\mathbf{A})$ and define $\mathbf{R} = \mathbf{D}^{-1}\mathbf{A}$ to be the random walk transition matrix, with \mathbf{R}_{ji} being the probability of moving from node i to node j . Consider an embedding model $f(\mathbf{W}, \mathbf{H}) := \mathbf{W}\mathbf{H}$, where $\mathbf{H} \in \mathbb{R}^{d \times n}$ are the input embeddings and $\mathbf{W} \in \mathbb{R}^{n \times d}$ are the output embeddings. We train the model on the bi-graph edges $(v_i, v_j) \in E$ by minimizing the CE loss

$$\min_{\mathbf{W}, \mathbf{H}} \mathcal{L}_{\text{ce}}(\mathbf{W}, \mathbf{H}) := \sum -\mathbf{Y}_{ij} \log(\mathbb{S}_i(\mathbf{W}\mathbf{h}_j)), \quad (1)$$

where $\mathbb{S}_i(z) := \left(\frac{e^{z_i}}{\sum_{i'} e^{z_{i'}}}\right)$ is the softmax and \mathbf{Y} is the label matrix, given either by \mathbf{R} itself or by its label-smoothed version $(1 - \epsilon)\mathbf{R} + \frac{\epsilon}{n}\mathbf{1}\mathbf{1}^\top$. We consider two cases, 1) **untied** embeddings, where \mathbf{W} and \mathbf{H} are free independent embeddings, and 2) **tied** embeddings, where the input and output embeddings of each node are constrained to be equal, i.e., $\mathbf{W}^\top = \mathbf{H}$.

When the embeddings are untied, the model is mathematically equivalent to word2vec (Mikolov et al., 2013) or the unconstrained features model (UFM) (Zhao et al., 2024), but is trained on the edges of a graph rather than on natural-language words and contexts. With weight-tying, this reduces to the one-layer *Node2vec* model (Grover & Leskovec, 2016) trained on random walks of length two and is the proxy model considered by Noroozizadeh et al. (2025).

We are interested in understanding the properties of the optimal embeddings learned by the CE loss in each case. We focus on the case where there is no dimension bottleneck, i.e., $d \geq n$. We use \square_t and \square_u to refer to variables corresponding to the tied and untied settings, respectively.

In this paper, we perform a case study on the cycle graph, and in the rest of the paper, we use \mathbf{A} to refer to the adjacency matrix of this graph specifically.

Definition 2.1 (Cycle graph). Let C_n denote the cycle graph on n nodes with adjacency matrix $\mathbf{A} \in \{0, 1\}^{n \times n}$ defined by $\mathbf{A}_{i+1,i} = \mathbf{A}_{i-1,i} = 1$, $\mathbf{A}_{j,i} = 0$ otherwise, where indices are taken modulo n .

Experimental setup As we will discuss in the next section, for the analysis, we consider the convex relaxation of (1) in the logit space. For the experiments, however, we directly train the bilinear model $\mathbf{W}\mathbf{H}$ (untied) or $\mathbf{W}\mathbf{W}^\top$ (tied) with CE loss (1). In all experiments, we set $n = 30$ and train the model long enough, with gradient descent, to be close to the global optimizer. For visualization, we project the embeddings \mathbf{W} and \mathbf{H} to three dimensions using PCA. Similar visualizations obtained with UMAP are provided in the appendix. Since these are only low-rank visualizations that can hide spectral information, we also present the full spectrums of the optimal solutions.

3 Preliminaries

For the analysis, we study the tight convex relaxation of the problem in the logit space $\mathbf{L} = \mathbf{W}\mathbf{H}$. Specifically,

$$\mathbf{L}_u^* \in \arg \min_{\mathbf{L}} \sum -\mathbf{Y}_{ij} \log(\mathbb{S}_i(\ell_j)) + \mu \|\mathbf{L}\|_*, \quad (2a)$$

$$\mathbf{L}_t^* \in \arg \min_{\mathbf{L} \succeq 0} \sum -\mathbf{Y}_{ij} \log(\mathbb{S}_i(\ell_j)). \quad (2b)$$

Note that the CE loss is strictly convex in every direction other than $\mathbf{1}\mathbf{1}^\top$. Thus, ignoring the constant shift, which

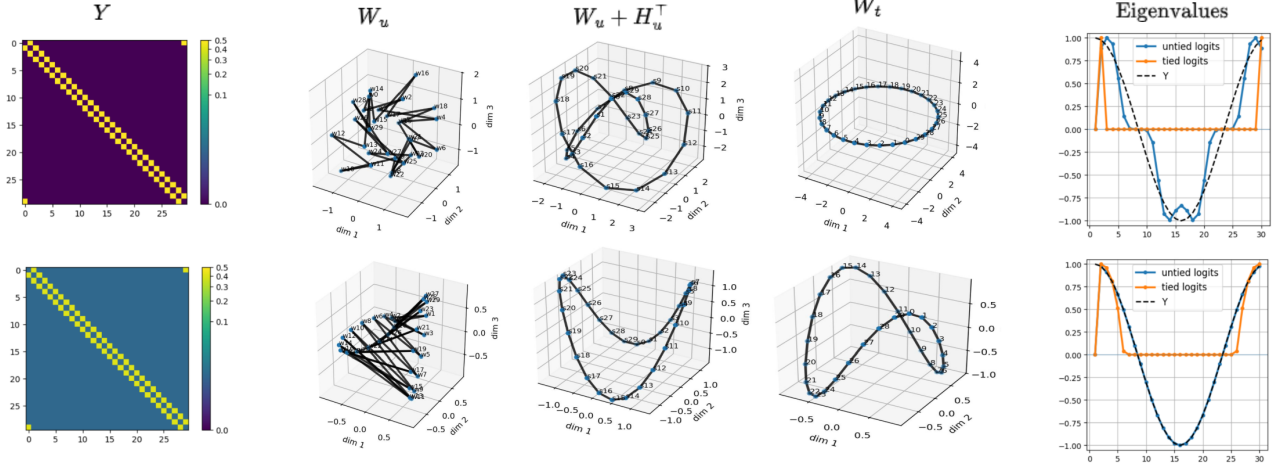


Figure 1. Bilinear embedding model \mathbf{WH} trained on cycle graph by cross-entropy loss, with (bottom) and without (top) label-smoothing for graph size $n = 30$. (Left) The target label matrix \mathbf{Y} used for training. (Middle) The PCA of embeddings. ($\mathbf{W}_u, \mathbf{H}_u$) refers to the minimizers of (1) and \mathbf{W}_t is the minimizer of (1) with weight-tying, i.e., constrained to the space $\mathbf{W} = \mathbf{H}^\top$. (Right) The normalized eigen values of the (blue) untied logits $\mathbf{W}_u \mathbf{H}_u$, (orange) tied logits $\mathbf{W}_t \mathbf{W}_t^\top$, (black) target label \mathbf{Y} . See Secs. 2-5 for details.

does not change the value of the softmax, the optimal logits are unique. In the untied case (2a), we add a (small) regularization μ to the loss, which corresponds to L2-regularization of the parameters by $\mu(\|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2)$ in (1). We consider this since, without regularization, the optimal embeddings \mathbf{W}, \mathbf{H} , i.e., the factorization of the optimal logit matrix, are not identifiable: if $\mathbf{L} = \mathbf{WH}$, then also $\mathbf{L} = \mathbf{W}'\mathbf{H}'$ for $\mathbf{W}' = \mathbf{W}\mathbf{Q}, \mathbf{H}' = \mathbf{Q}^{-1}\mathbf{H}$ for any invertible \mathbf{Q} . The regularization will make the optimal embeddings unique (up to some rotation). Then, we have a one-to-one correspondence between the geometry of the embeddings and the optimal logits as in Thrapoulidis et al. (2022); Zhao et al. (2024)

Lemma 3.1 (Convex relaxation). *Assume $d \geq n$. Let $\mathbf{L}_u^* = \mathbf{U}_u \Sigma_u \mathbf{V}_u^\top$ be the SVD of the optimal logits of (2a) and $\mathbf{L}_t^* = \mathbf{U}_t \Sigma_t \mathbf{U}_t^\top$ be the eigenvalue decomposition of the optimal logits of (2b). Also, let \mathbf{Q} be any rotation matrix, i.e., $\mathbf{Q}\mathbf{Q}^\top = \mathbf{I}$.*

1. (Untied) The optimal embeddings $(\mathbf{W}_u^*, \mathbf{H}_u^*)$ of (L2-regularized) (1) are unique up to a rotation.

$$\mathbf{W}_u^* = \mathbf{U}_u \Sigma_u^{1/2} \mathbf{Q}, \quad \mathbf{H}_u = \mathbf{Q}^\top \Sigma_u^{1/2} \mathbf{V}_u^\top,$$

2. (Tied) The optimal embeddings \mathbf{W}_t^* of (1), constrained to $\mathbf{W} = \mathbf{H}^\top$, are unique up to a rotation.

$$\mathbf{W}_t^* = \mathbf{U}_t \Sigma_t^{1/2} \mathbf{Q}.$$

Note that the above lemma implies the *geometry* of the embeddings, characterized by their Gram matrix is unique, e.g., $\mathbf{W}_u^* \mathbf{W}_u^{*\top} = \mathbf{U}_u^* \Sigma_u^* \mathbf{U}_u^{*\top}$.

As mentioned in the previous section, we consider the circle graph throughout this paper. We make this choice since it yields a circulant symmetric target matrix \mathbf{Y} .

Lemma 3.2. *Assume \mathbf{Y} is circulant symmetric. Then, the optimal logits in (2) are also circulant and symmetric.*

The lemma follows from a symmetry argument detailed in A.1. Since all circulant and symmetric matrices are diagonalizable by the Discrete Fourier Transform (DFT) matrix, we can focus on simply characterizing the optimal eigenvalues of the logit matrix $\lambda(\mathbf{L}^*)$ to compare the global optimizers in different cases. Also this together with Lem. 3.1 implies that characterizing $\lambda(\mathbf{L})$ also specifies which directions, along the Fourier vectors of different frequencies, are dominant in the embedding space of \mathbf{W} and \mathbf{H} . Thus, in Sec. 5, we only focus on finding the optimal eigen-spectrum $\lambda(\mathbf{L}^*)$ in each scenario, comparing them with the spectrum of the Cycle graph.

Fact 1. *For the cycle graph \mathbf{A} in Defn. 2.1, the eigenvalues are $\lambda_m(\mathbf{Y}) = \cos(2\pi m/n)$.*

In our experiments (Figs. 1-2), we also show the PCA of $\mathbf{W} + \mathbf{H}^\top$ in the untied case. This is motivated by some classic works on word embeddings, which consider the sum of the input and output embeddings as the final embedding for downstream tasks (e.g., Pennington et al., 2014). In this case, the sum of the embeddings retains only the positive part of the optimal spectrum.

Lemma 3.3. *Let $\mathbf{L}_{n \times n}$ be a symmetric logit matrix, and let $\mathbf{L} = \mathbf{U}\Sigma\mathbf{V}^\top$ be an SVD. Define $\mathbf{W} = \mathbf{U}\sqrt{\Sigma}\mathbf{R}$ and $\mathbf{H} = \mathbf{R}^\top\sqrt{\Sigma}\mathbf{V}^\top$, for some rotation matrix. Then, with $\mathbf{Z} = \mathbf{W}^\top + \mathbf{H}$, we have $\mathbf{Z}^\top\mathbf{Z} = \mathbf{L}_+$, where \mathbf{L}_+ is the PSD part of \mathbf{L} .*

4 Warm-up: Approximating CE with MSE

The CE objective is generally harder to analyze directly because of the nonlinear softmax normalization. A common simplification is to replace CE by its second-order Taylor approximation around zero logits, which yields a quadratic proxy. Although our analysis in this paper does not rely on this approximation, it is useful as a baseline for understanding what it would predict in our setting. Applying this approximation to (1) gives

$$\arg \min \|\mathbf{W}\mathbf{H} - (\mathbb{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)\mathbf{Y}\|_F^2, \quad (3)$$

i.e., minimizing the MSE between the logits $\mathbf{L} = \mathbf{W}\mathbf{H}$ and the centered target matrix $\bar{\mathbf{Y}} = (\mathbb{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top)\mathbf{Y}$. Let $\bar{\mathbf{Y}} = \bar{\mathbf{Y}}_+ + \bar{\mathbf{Y}}_-$ denote the decomposition of $\bar{\mathbf{Y}}$ into its positive and negative spectral parts, with $\bar{\mathbf{Y}}_+ \succeq 0$ and $\bar{\mathbf{Y}}_- \preceq 0$. In the untied case, with no rank bottleneck, the minimizer of (3) is $\mathbf{L}_u^* = \mathbf{W}_u\mathbf{H}_u = \bar{\mathbf{Y}}$. In the tied case, the logits are constrained to be PSD, $\mathbf{L} = \mathbf{W}\mathbf{W}^\top \succeq 0$, and the minimizer is the PSD projection $\mathbf{L}_t^* = \mathbf{W}_t\mathbf{H}_t = \bar{\mathbf{Y}}_+$. Thus, under the MSE proxy, the untied logits recover the full target spectrum, while the tied logits retain only its positive part. Fig. 2 illustrates this prediction for the cycle graph.

This quadratic picture is useful, but it does not capture the behavior of the CE objective in the sparse setting shown in Fig. 1. For sparse targets, the CE optimum can be far from the origin, and in the untied case the logits diverge in norm. As a result, the second-order approximation around zero can give a qualitatively incorrect prediction of the learned spectrum. In the next section, we therefore analyze the CE objective directly. We will see that the MSE picture becomes more relevant only after label smoothing makes the target dense and the finite CE optimum moves closer to the quadratic regime.

Remark 4.1. Second-order proxies have been used in prior work to make embedding dynamics analytically tractable. For example, Karkada et al. (2025) use such an approximation to study word2vec-like models, and Karkada et al. (2026) uses this proxy to connect symmetry in data statistics to symmetry in learned embeddings of LLMs trained with next-token-prediction. Our results suggest that this MSE picture is most informative in regimes where the target statistics are dense.

5 CE Loss

Here, we characterize the global solutions of (2) for two scenarios: training on the sparse transition matrix $\mathbf{Y} = \mathbf{R} = \frac{1}{2}\mathbf{A}$ in Sec. 5.1, and the dense transition matrix $\mathbf{Y}_\epsilon = (1 - \epsilon)\mathbf{R} + \frac{\epsilon}{n}\mathbf{1}\mathbf{1}^\top$, the result of label smoothing, in Sec. 5.2.

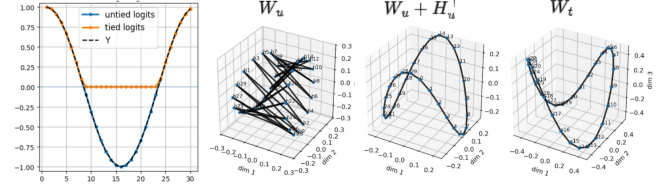


Figure 2. Embeddings trained by MSE loss in (3) on cycle graph of size $n = 30$. (Left) Eigenvalues of the optimal logits and the targets \mathbf{Y} . (Right) PCA of the learned embedding in the untied ($\mathbf{W}_u, \mathbf{H}_u$) and tied \mathbf{W}_t case.

5.1 Sparse Targets

Since the graph adjacency matrix is sparse, there is no finite \mathbf{L} that can reach the entropy lower-bound of the CE loss, i.e., $\sum \mathbf{Y}_{ij} \log(\mathbf{Y}_{ij})$. This is the case as reaching the entropy requires $\mathbb{S}(\mathbf{L}) = \mathbf{R}$, and entries of $\mathbb{S}(\mathbf{L})$ are strictly positive for finite \mathbf{L} . On the other hand, any \mathbf{L} that satisfies the interpolating conditions below, identifies the asymptotic direction of an optimizer, as $\mathbb{S}(\alpha\mathbf{L}) = \mathbf{R}$, $\alpha \rightarrow \infty$.

$$\begin{aligned} \mathbf{L}_{ij} - \mathbf{L}_{i'j} &= 0 & \text{if } \mathbf{A}_{ij} = \mathbf{A}_{i'j} = 1, \\ \mathbf{L}_{ij} - \mathbf{L}_{\ell j} &> 0 & \text{if } \mathbf{A}_{ij} = 1, \mathbf{A}_{\ell j} = 0. \end{aligned} \quad (4)$$

When there are no constraints on \mathbf{L} , as in the untied case of (2a), there are clearly infinite \mathbf{L} 's that satisfy the interpolating conditions (4). However, in the limit of vanishing regularization, inductive bias results point at a specific max-margin solution (Zhao et al., 2024).

Theorem 5.1 (Zhao et al. (2024, Thm. 1)). *As $\mu \rightarrow 0$ in (2a), \mathbf{L}_u^* increases in norm but converges in direction to the following max-margin solution. That is $\mathbf{L}_u^* / \|\mathbf{L}_u^*\|_* \rightarrow \hat{\mathbf{L}} / \|\hat{\mathbf{L}}\|_*$, where*

$$\begin{aligned} \hat{\mathbf{L}} &= \arg \min_{\mathbf{L}} \|\mathbf{L}\|_* \quad \text{s.t.} \\ \mathbf{L}_{ij} - \mathbf{L}_{i'j} &= 0 & \text{if } \mathbf{A}_{ij} = \mathbf{A}_{i'j} = 1, \\ \mathbf{L}_{ij} - \mathbf{L}_{\ell j} &\geq 1 & \text{if } \mathbf{A}_{ij} = 1, \mathbf{A}_{\ell j} = 0. \end{aligned} \quad (5)$$

We cannot solve (5) in closed form. However, we can efficiently solve it numerically for different graph sizes n . First, we note that similar to Lem. 3.2 by symmetry, we can argue that the optimal solution to (5) is circulant and symmetric. So, let the first row of \mathbf{L} be (c_0, \dots, c_{n-1}) , with $c_k = c_{n-k}$ and $L_{ij} = c_{j-i}$. Since DFT matrix specifies the eigenvectors of a circulant symmetric matrix, we can write the Fourier eigenvalues of $\hat{\mathbf{L}}$ as $\lambda_m = \sum_{k=0}^{n-1} c_k e^{-i(2\pi mk/n)}$, $m = 0, \dots, n-1$. Then $\lambda_m \in \mathbb{R}$ and $\lambda_{n-m} = \lambda_m$, and $\|\mathbf{L}\|_* = \sum_{m=0}^{n-1} |\lambda_m|$. On the other hand, using the inverse DFT, the nontrivial constraints $c_1 - c_k \geq 1$ for all $k \notin \{1, n-1\}$ becomes $\frac{1}{n} \sum_{m=0}^{n-1} \lambda_m \left(\cos \frac{2\pi m}{n} - \cos \frac{2\pi mk}{n} \right) \geq 1$, as $c_k = \frac{1}{n} \sum_{m=0}^{n-1} \lambda_m e^{2\pi imk/n}$. So the nuclear norm minimization

in (5) reduces to

$$\min_{\lambda \in \mathbb{R}^n} \sum_{m=0}^{n-1} |\lambda_m| \quad \text{s.t.} \quad \forall k \notin \{1, n-1\} \quad (6)$$

$$\frac{1}{n} \sum_{m=0}^{n-1} \lambda_m \left(\cos \frac{2\pi m}{n} - \cos \frac{2\pi mk}{n} \right) \geq 1,$$

where λ_m 's are the eigenvalues of the optimal $\hat{\mathbf{L}}$. We cannot solve this linear program in closed form, but we can solve it numerically, as shown in Fig. 3-(a) for various values of n .

While there are infinite \mathbf{L} 's that satisfy (4), it turns out none satisfy the PSD constraint $\mathbf{L} \succeq 0$ in the tied case.

Proposition 5.2. *Let \mathbf{A} be the adjacency matrix of a cycle graph of size $n \geq 3$. Any logits matrix \mathbf{L} that satisfies (4) has a negative eigenvalue.*

In fact, by solving the KKT conditions in (2b), we get a closed-form for the unique optimal solution \mathbf{L}_t^* , which is finite and does not achieve the entropy lower-bound.

Theorem 5.3. *Let $G = C_n$ be the cycle graph with $n \geq 5$, and let $\mathbf{Y} = \mathbf{A}/2$. Problem (2b) has a unique optimal solution (up to a constant shift) of the form,*

$$\mathbf{L}_{t,ij}^* = t^* \cos(2\pi(i-j)/n) + a, \quad (7)$$

for any $a \geq 0$ and some $t^* \geq 0$. Consequently, ignoring the constant shift a , \mathbf{L}_t^* has rank 2, with eigenvalues $\lambda_1 = \lambda_{n-1} = \text{const}$, and $\lambda_m = 0$ for $m \in \{0, 2, 3, \dots, n-2\}$, and the optimal solution of (1) constrained on tied weights, up to some rotation, is realized by,

$$\mathbf{h}_i^* = \mathbf{w}_i^* = \sqrt{t^*} (\cos(2\pi i/n), \sin(2\pi i/n))^\top.$$

In short, the proof translates the KKT conditions to conditions on the eigenvalues $\ell_m := \lambda_m(\mathbf{L})$, $s_m := \lambda_m(\mathbf{S}(\mathbf{L}))$, and $y_m := \lambda_m(\mathbf{Y})$. Specifically, $\ell_m \geq 0$ for feasibility, $s_m \geq \ell_m$ for dual feasibility and stationarity condition, and $\ell_m(s_m - y_m) = 0$ from complementary slackness. It is the complementary slackness together with the fact that s_m is the Fourier transform of a probability vector that enforces $\ell_m = 0$ for $m \neq 1, n-1$ since s_m cannot match $y_m = \cos(2\pi m/n)$. We defer the detailed proof to A.4.

While the tied model $f(\mathbf{W}) = \mathbf{W}\mathbf{W}^\top$ does not have the capacity to interpolate the labels, Eq. (7) shows the optimal embeddings form a circle in the embedding space, which is also aligned with our observations in Fig. 1. Notably, the embeddings themselves, even without the PCA dimension reduction, form this circular shape (up to some rotation), which has a nice reflection of the graph, and the rank of the optimal solution in this setting is constant irrespective of the size of the graph n . This is in contrast to the untied solution found by (6), which has a relatively dense spectrum with its rank growing as $\Theta(n)$ as shown in Fig. 3.

Remark 5.4. The eigenspectrum in the untied case has large values around four different frequency bands, rather than concentrating dominantly on a pair of modes. A three-dimensional PCA projection can therefore capture only part of this structure, which makes the resulting embedding visualization appear less organized than the tied solution in Fig. 1. For the untied model, visualizing $\mathbf{W} + \mathbf{H}^\top$ can look more structured because its Gram matrix retains only the positive spectral part of \mathbf{L}_u^* . Similar behavior appears with UMAP. This highlights that low-dimensional projections can obscure the symmetries in the embedding space, making it hard to make interpretations of the geometry.

Next, we will see how the optimal spectrum changes in each case if we remove the sparsity of the label matrix by applying label-smoothing before the loss computation.

5.2 Dense Targets with Label-Smoothing

Here, we consider label-smoothed target $\mathbf{Y}_\epsilon = (1 - \epsilon)\mathbf{R} + \frac{\epsilon}{n}\mathbf{1}\mathbf{1}^\top$ for training. This will remove all the zero entries in the target matrix, which in turn relaxes the need for parameters growing to infinity for minimizing the loss function. In this case, we can find the optimal finite solution by solving a linear system of equations.

Theorem 5.5 (Zhao et al. (2024, Thm. 2)). *Let \mathbf{Y}_ϵ be the target label in (2a) for some $\epsilon \in (0, 1)$. The optimal solution \mathbf{L}_u^* is circulant with first row $\mathbf{c} = (c_0, c_1, \dots, c_{n-1})$, where \mathbf{c} solves, $\forall j \neq 1, n-1$,*

$$c_1 - c_j = \log \left(\frac{n + (2-n)\epsilon}{2\epsilon} \right) =: \Delta$$

$$c_1 = c_{n-1}$$

In other words, $\mathbf{c} = \Delta(e_1 + e_{n-1} - \frac{2}{n}\mathbf{1}) + a\mathbf{1}$.

The eigenvalues (ignoring the DC, which is affected by the constant shift a) of the corresponding \mathbf{L}_u^* then become

$$\lambda_{m \neq 0}(\mathbf{L}_u^\epsilon) = \sum_{j=0}^{n-1} c_j \cos(2\pi m j/2) = 2\Delta \cos\left(\frac{2\pi m}{n}\right),$$

which is the same as the eigenvalues of the graph adjacency matrix \mathbf{A} after normalization (black line in Fig. 3-(c)).

In the tied case (2b) with dense targets \mathbf{Y}_ϵ , we cannot solve the solution in closed form, but we can simplify the objective in terms of the eigenvalues. Per Lem. 3.2, the solution is circulant symmetric, which simplifies the problem as

$$\min_{\mathbf{c}} \left(\log \sum_{i=0}^{n-1} e^{c_i} - \frac{1-\epsilon}{2}(c_1 + c_{n-1}) - \frac{\epsilon}{n} \sum_{i=1}^{n-1} c_i \right)$$

Rewrite $c_i(\lambda) = \sum_{m=1}^{n-1} \lambda_m \cos(mi\theta)$. Then the loss can

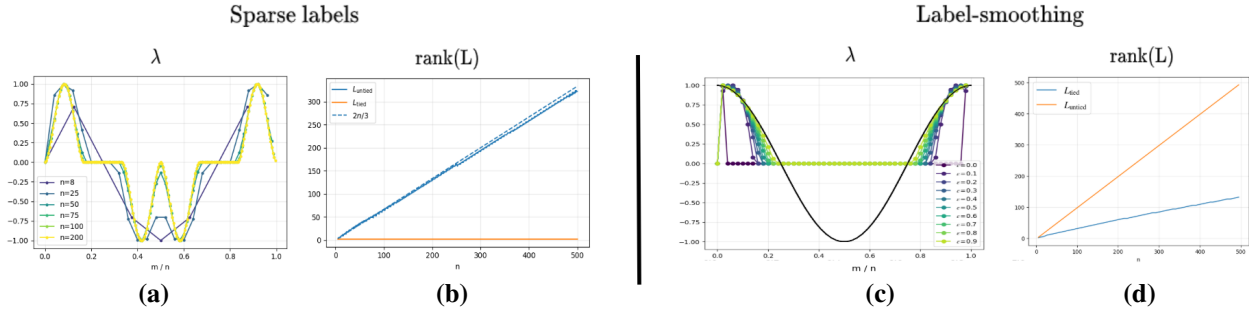


Figure 3. Properties of the global optimizers of (2) trained on (Left) sparse random walk matrix \mathbf{R} of the cycle graph and (Right) \mathbf{R} after label-smoothing with $\epsilon \in (0, 1)$. (a) Normalized eigenvalues λ_m of the optimal untied logit \mathbf{L}_u for various values of graph size n , found by solving the equivalent max-margin problem (6). (b) Rank of the untied logits and tied logits across n . (c) Normalized eigenvalues of the optimal tied logits \mathbf{L}_t for $n = 30$ and different values of ϵ , along with the normalized eigenvalues of the untied logits, which match the eigenvalues of the smoothed label matrix \mathbf{Y}_ϵ after normalization (black line). Eigenvalues in the tied case are found by solving the optimization in (8). (d) Same as (b) but with smoothed labels ($\epsilon = 0.2$).

be written directly in terms of the eigenvalues of \mathbf{L} .

$$\log \sum_{i=0}^{n-1} \exp \left(\sum_{m=1}^{n-1} \lambda_m \cos(mi\theta) \right) - (1 - \epsilon) \left(\sum_{m=1}^{n-1} \lambda_m \cos(m\theta) \right) \quad (8)$$

We can now numerically, optimizing this over $\lambda_m \geq 0$, to enforce PSD condition. We visualize the optimal spectrum for various values of ϵ and $n = 30$ in Fig. 3-(c). For $\epsilon > 0$ the (normalized) optimal spectrum approximates the positive part of the untied solution in a way that progressively becomes closer as ϵ grows and the target moves away from sparsity.

6 Discussion and Limitations

Through this case study, we highlighted how small tweaks in the training can significantly change the geometry learned by CE. However, we do not claim that one of these geometries is universally preferable. Different geometries preserve different spectral summaries of the same training distribution, and their usefulness should depend on the downstream application. Formalizing this task dependence and the generalization performance is an important next step. Also, our analysis relies on the nice properties of the circulant matrices. We leave extending this characterization to other data structures for future work.

References

Allen, C. and Hospedales, T. Analogies explained: Towards understanding word embeddings. In *International Conference on Machine Learning*, pp. 223–231. PMLR, 2019.

Bangachev, K., Bresler, G., Noman, I., and Polyanskiy, Y.

Global minimizers of sigmoid contrastive loss. *arXiv preprint arXiv:2509.18552*, 2025.

Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.

Gittens, A., Achlioptas, D., and Mahoney, M. W. Skip-gram- zipf+ uniform= vector additivity. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 69–76, 2017.

Grover, A. and Leskovec, J. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 855–864, 2016.

Gurnee, W. and Tegmark, M. Language models represent space and time. *arXiv preprint arXiv:2310.02207*, 2023.

Jiang, J., Zhou, J., Wang, P., Qu, Q., Mixon, D., You, C., and Zhu, Z. Generalized neural collapse for a large number of classes. *arXiv preprint arXiv:2310.05351*, 2023.

Jiang, Y., Rajendran, G., Ravikumar, P., Aragam, B., and Veitch, V. On the origins of linear representations in large language models. *arXiv preprint arXiv:2403.03867*, 2024.

Karkada, D., Simon, J. B., Bahri, Y., and DeWeese, M. R. Closed-form training dynamics reveal learned features and linear structure in word2vec-like models. *arXiv preprint arXiv:2502.09863*, 2025.

Karkada, D., Korchinski, D. J., Nava, A., Wyart, M., and Bahri, Y. Symmetry in language statistics shapes the geometry of model representations. *arXiv preprint arXiv:2602.15029*, 2026.

- 330 Lee, A., Weber, M., Viégas, F., and Wattenberg, M. Shared
 331 global and local geometry of language model embeddings.
 332 *arXiv preprint arXiv:2503.21073*, 2025.
 333
- 334 Levy, O. and Goldberg, Y. Neural word embedding as
 335 implicit matrix factorization. *Advances in neural infor-*
 336 *mation processing systems*, 27, 2014.
- 337 Li, M. Z., Agrawal, K. K., Ghosh, A., Teru, K. K., Santoro,
 338 A., Lajoie, G., and Richards, B. A. Tracing the represen-
 339 tation geometry of language models from pretraining to
 340 post-training. *arXiv preprint arXiv:2509.23024*, 2025.
 341
- 342 Li, P., Li, X., Wang, Y., and Qu, Q. Neural collapse in multi-
 343 label learning with pick-all-label loss. *arXiv preprint*
 344 *arXiv:2310.15903*, 2023.
- 345 McInnes, L., Healy, J., and Melville, J. Umap: Uniform
 346 manifold approximation and projection for dimension
 347 reduction. *arXiv preprint arXiv:1802.03426*, 2018.
 348
- 349 Mikolov, T., Chen, K., Corrado, G., and Dean, J. Efficient
 350 estimation of word representations in vector space. *arXiv*
 351 *preprint arXiv:1301.3781*, 2013.
 352
- 353 Noroozizadeh, S., Nagarajan, V., Rosenfeld, E., and Kumar,
 354 S. Deep sequence models tend to memorize geometri-
 355 cally; it is unclear why. *arXiv preprint arXiv:2510.26745*,
 356 2025.
- 357 Papyan, V., Han, X., and Donoho, D. L. Prevalence of
 358 neural collapse during the terminal phase of deep learn-
 359 ing training. *Proceedings of the National Academy of*
 360 *Sciences*, 117(40):24652–24663, 2020.
 361
- 362 Park, K., Choe, Y. J., and Veitch, V. The linear represen-
 363 tation hypothesis and the geometry of large language
 364 models. *arXiv preprint arXiv:2311.03658*, 2023.
 365
- 366 Pennington, J., Socher, R., and Manning, C. D. Glove:
 367 Global vectors for word representation. In *Proceedings*
 368 *of the 2014 conference on empirical methods in natural*
 369 *language processing (EMNLP)*, pp. 1532–1543, 2014.
- 370 Qiu, J., Dong, Y., Ma, H., Li, J., Wang, K., and Tang, J.
 371 Network embedding as matrix factorization: Unifying
 372 deepwalk, line, pte, and node2vec. In *Proceedings of the*
 373 *eleventh ACM international conference on web search*
 374 *and data mining*, pp. 459–467, 2018.
 375
- 376 Saxe, A. M., McClelland, J. L., and Ganguli, S. A mathe-
 377 matical theory of semantic development in deep neural
 378 networks. *Proceedings of the National Academy of Sci-*
 379 *ences*, 116(23):11537–11546, 2019.
- 380 Thrapoulidis, C., Kini, G. R., Vakilian, V., and Behnia, T.
 381 Imbalance trouble: Revisiting neural-collapse geometry.
 382 *Advances in Neural Information Processing Systems*, 35:
 383 27225–27238, 2022.
 384
- Zhao, Y. and Thrapoulidis, C. On the geometry of se-
 mantics in next-token prediction. In *The First Workshop*
on the Interplay of Model Behavior and Model Internals,
 2025.
- Zhao, Y., Behnia, T., Vakilian, V., and Thrapoulidis, C.
 Implicit geometry of next-token prediction: From lan-
 guage sparsity patterns to model representations. *arXiv*
preprint arXiv:2408.15417, 2024.
- Zhu, Z., Ding, T., Zhou, J., Li, X., You, C., Sulam, J., and
 Qu, Q. A geometric analysis of neural collapse with
 unconstrained features. *Advances in Neural Information*
Processing Systems, 34:29820–29834, 2021.

A Proofs

A.1 Proof of Prop. 5.2

Suppose for the sake of contradiction that such a positive semi-definite matrix $\mathbf{L} \succeq 0$ exists.

Because \mathbf{A} represents a simple cycle graph, it contains no self-loops, meaning $\mathbf{A}_{jj} = 0$ for all vertices j . For any adjacent vertices i and j , we have $\mathbf{A}_{ij} = 1$. Letting $\ell = j$, the strict inequality in (4) requires

$$\mathbf{L}_{ij} - \mathbf{L}_{jj} > 0 \implies \mathbf{L}_{ij} > \mathbf{L}_{jj}.$$

By the symmetry of the undirected graph, we also have $\mathbf{A}_{ji} = 1$. Letting $\ell = i$ in (4) similarly yields

$$\mathbf{L}_{ji} - \mathbf{L}_{ii} > 0 \implies \mathbf{L}_{ji} > \mathbf{L}_{ii}.$$

Summing these two inequalities gives

$$\mathbf{L}_{ij} + \mathbf{L}_{ji} > \mathbf{L}_{ii} + \mathbf{L}_{jj}. \quad (9)$$

Since \mathbf{L} is assumed to be positive semi-definite, we must have $x^\top \mathbf{L} x \geq 0$ for any vector $x \in \mathbb{R}^n$. Choosing $x = e_i - e_j$, where e_k denotes the k -th standard basis vector, we obtain

$$(e_i - e_j)^\top \mathbf{L} (e_i - e_j) = \mathbf{L}_{ii} + \mathbf{L}_{jj} - \mathbf{L}_{ij} - \mathbf{L}_{ji} \geq 0.$$

Rearranging this yields

$$\mathbf{L}_{ii} + \mathbf{L}_{jj} \geq \mathbf{L}_{ij} + \mathbf{L}_{ji},$$

which directly contradicts the strict inequality established in (9). Therefore, no such positive semi-definite matrix \mathbf{L} can exist.

A.2 Proof of Lem. 3.2

Consider a permutation matrix P_t where it either shifts the indices or mirrors and shifts them ($\pi(i) = i + t$ or $\pi(i) = t - i$). Then, $P\mathbf{A}P^\top = \mathbf{A}$. For $\mathbf{L}' = P\mathbf{L}P^\top$, we have $\mathbf{L}'_{ij} = \mathbf{L}_{\pi(i),\pi(j)}$ and thus, $\mathcal{L}_{\text{ce}}(\mathbf{L}, \mathbf{A}) = \mathcal{L}_{\text{ce}}(\mathbf{L}', \mathbf{A}') = \mathcal{L}_{\text{ce}}(\mathbf{L}', \mathbf{A})$. Now consider all such permutations of a feasible \mathbf{L} and define the convex combination

$$\bar{\mathbf{L}} = \frac{1}{2n} \sum_{t=1}^n (P_t \mathbf{L} P_t^\top + P_{-t} \mathbf{L} P_{-t}^\top) \quad (10)$$

Then, $\bar{\mathbf{L}}$ is circular and symmetric, and since CE is convex, we have,

$$\mathcal{L}_{\text{ce}}(\bar{\mathbf{L}}) \leq \frac{1}{2n} \sum_{t=1}^n (\mathcal{L}_{\text{ce}}(P_t \mathbf{L} P_t^\top) + \mathcal{L}_{\text{ce}}(P_{-t} \mathbf{L} P_{-t}^\top)) = \mathcal{L}_{\text{ce}}(\mathbf{L}). \quad (11)$$

Thus, if \mathbf{L} is optimal, there exists a circulant and symmetric $\bar{\mathbf{L}}$ that is optimal. However, the solution to CE is unique up to some $\mathbf{1}\mathbf{1}^\top$ shift, which means any optimal solution needs to be circulant and symmetric.

A.3 Proof of Lem. 3.3

Since \mathbf{L} is symmetric, it has an eigenvalue decomposition $\mathbf{L} = P\Lambda P^\top$. Then, $\mathbf{U} = P$, $\Sigma = |\Lambda|$, $\mathbf{V} = P\mathbf{S}$, where $\mathbf{S} = \text{sign}(\Lambda)$, gives the SVD. Then, $\mathbf{L} = \mathbf{W}^\top + \mathbf{H} = \mathbf{R}^\top \sqrt{\Sigma}(\mathbf{I} + \mathbf{S})P^\top$.

$$\mathbf{L}^\top \mathbf{L} = P(\mathbf{I} + \mathbf{S})\Sigma(\mathbf{I} + \mathbf{S})P^\top = 4P\Lambda_+P^\top.$$

A.4 Proof of Thm. 5.3

To derive the optimality conditions, we introduce a dual variable matrix $\Lambda \in \mathbb{S}^n$ for the positive semidefinite constraint and define the Lagrangian as $\mathcal{L}(L, \Lambda) = f(L) - \text{tr}(\Lambda L)$. A candidate primal-dual pair (L^*, Λ^*) is globally optimal if and only

if it satisfies the following four KKT conditions:

$$\text{Stationarity: } ((\mathbb{S}(\mathbf{L}^*) - \mathbf{Y}) + (\mathbb{S}(\mathbf{L}^*) - \mathbf{Y})^\top) = 2\Lambda^*$$

$$\text{Primal Feasibility: } \mathbf{L}^* \succeq 0$$

$$\text{Dual Feasibility: } \Lambda^* \succeq 0$$

$$\text{Complementary Slackness: } \Lambda^* \mathbf{L}^* = 0$$

Since \mathbf{Y} is symmetric and circular, by Lemma 3.2, the primal optimization can be restricted to be over circular matrices, i.e., \mathbf{L}^* is circular. Then $\mathbb{S}(\mathbf{L}^*)$ is also circular and symmetric and the KKT conditions can be written in terms of the eigenvalues of the matrices, call them ℓ_i , y_i and s_i for \mathbf{L}^* , \mathbf{Y} and $\mathbb{S}(\mathbf{L}^*)$ respectively:

$$\ell_i \geq 0, \quad s_i \geq y_i, \quad \ell_i(s_i - y_i) = 0$$

For the cycle graph, the transition matrix is $\mathbf{Y} = \mathbf{A}/2$, so its circulant eigenvalues are $y_m = \cos(m\theta)$ for $m = 0, \dots, n-1$. Let (c_0, \dots, c_{n-1}) be the first row of \mathbf{L}^* , and let (p_0, \dots, p_{n-1}) be the first row of $\mathbb{S}(\mathbf{L}^*)$, namely $p_k = \frac{e^{c_k}}{\sum_{r=0}^{n-1} e^{c_r}}$. Since $\mathbb{S}(\mathbf{L}^*)$ is also symmetric circulant, its eigenvalues are the Fourier coefficients of p , i.e., $s_m = \sum_{k=0}^{n-1} p_k e^{-imk\theta} = \sum_{k=0}^{n-1} p_k \cos(mk\theta)$. Hence the KKT conditions reduce to

$$\ell_m \geq 0, \quad s_m \geq \cos(m\theta), \quad \ell_m(s_m - \cos(m\theta)) = 0, \quad m = 0, \dots, n-1.$$

We now show that no Fourier mode beyond the first one can be active at the optimum. Fix $m \in \{2, \dots, \lfloor n/2 \rfloor\}$ and define $\beta_m := \frac{1 - \cos(m\theta)}{1 - \cos\theta}$. On the cycle grid $k = 0, \dots, n-1$, one has the pointwise inequality by Lemma A.1

$$\cos(mk\theta) \geq 1 - \beta_m + \beta_m \cos(k\theta),$$

with equality only for $k \in \{0, 1, n-1\}$. Averaging with respect to the probability vector p gives

$$s_m = \sum_{k=0}^{n-1} p_k \cos(mk\theta) \geq 1 - \beta_m + \beta_m \sum_{k=0}^{n-1} p_k \cos(k\theta) = 1 - \beta_m + \beta_m s_1.$$

Since $p_k > 0$ for every k and $n \geq 5$, the inequality is in fact strict. Assume now that $\ell_m > 0$ for some $m \geq 2$. Then complementary slackness gives $s_m = \cos(m\theta)$. Combining this with the previous inequality yields

$$\cos(m\theta) = s_m > 1 - \beta_m + \beta_m s_1.$$

Since $\cos(m\theta) = 1 - \beta_m + \beta_m \cos\theta$, this implies $s_1 < \cos\theta$, which contradicts dual feasibility at $m = 1$, namely $s_1 \geq \cos\theta$. Therefore $\ell_m = 0$ for every $m = 2, \dots, \lfloor n/2 \rfloor$, and by symmetry also $\ell_{n-m} = 0$.

Thus the only possibly nonzero Fourier modes of \mathbf{L}^* are the zero mode and the first harmonic. Since the zero mode only adds a constant shift, we obtain

$$\mathbf{L}_{ij}^* = t^* \cos\left(\frac{2\pi(i-j)}{n}\right) + a$$

for some $t^* \geq 0$ and constant a . It remains to identify t^* . The first row is $c_k = t^* \cos(k\theta) + a$, so the corresponding softmax probabilities are

$$p_k(t) = \frac{e^{t \cos(k\theta)}}{\sum_{r=0}^{n-1} e^{t \cos(r\theta)}},$$

where the constant shift a cancels. Hence

$$s_1(t) = \sum_{k=0}^{n-1} p_k(t) \cos(k\theta) = \frac{\sum_{k=0}^{n-1} \cos(k\theta) e^{t \cos(k\theta)}}{\sum_{k=0}^{n-1} e^{t \cos(k\theta)}}.$$

Since ℓ_1 cannot be zero, otherwise all nonconstant modes vanish and $s_1 = 0 < \cos \theta$, we must have $\ell_1 > 0$. Therefore, complementary slackness at $m = 1$ gives

$$s_1(t^*) = \cos \theta.$$

Finally, the map $t \mapsto s_1(t)$ is strictly increasing, because

$$s_1'(t) = \sum_{k=0}^{n-1} p_k(t) \cos^2(k\theta) - \left(\sum_{k=0}^{n-1} p_k(t) \cos(k\theta) \right)^2 = \text{Var}_{p(t)}(\cos(K\theta)) > 0.$$

Moreover, $s_1(0) = 0$ and $\lim_{t \rightarrow \infty} s_1(t) = 1$. Since $n \geq 5$, we have $0 < \cos \theta < 1$, so there exists a unique $t^* > 0$ such that $s_1(t^*) = \cos \theta$. This proves that every optimal solution is of the form

$$L_{ij}^* = t^* \cos\left(\frac{2\pi(i-j)}{n}\right) + a,$$

with constant $a \geq 0$ to satisfy the PSD constraint. Ignoring the DC component, the optimal embeddings $\mathbf{L}^* = \mathbf{W}^* \mathbf{W}^{*\top}$ of (1) with weight-tying become rank 2 and of the form $\mathbf{w}_i^* = \sqrt{t^*} (\cos(2\pi i/n), \sin(2\pi i/n))^\top$, forming a circle in 2D.

Lemma A.1. *Let $n \geq 5$ and $\theta = 2\pi/n$. Define $\Delta_m(k) := 1 - \cos(mk\theta)$, $\beta_m := \frac{1 - \cos(m\theta)}{1 - \cos(\theta)}$. Then, for every $m = 1, \dots, n-1$ and every $k = 0, \dots, n-1$, $\Delta_m(k) \leq \beta_m \Delta_1(k)$. Moreover, when $m = 2, \dots, \lfloor n/2 \rfloor$, equality holds only for $k \in \{0, 1, n-1\}$.*

Proof. The case $k = 0$ is immediate, so assume $k \neq 0$. Using $1 - \cos x = 2 \sin^2(x/2)$, we have

$$\frac{\Delta_m(k)}{\Delta_1(k)} = \left(\frac{\sin(mk\theta/2)}{\sin(k\theta/2)} \right)^2.$$

Let $a := \theta/2 = \pi/n$, $x := ka$, and $p := \min\{m, n-m\}$. Since $\sin((n-m)x) = (-1)^{k+1} \sin(mx)$, it is enough to show

$$\left| \frac{\sin(px)}{\sin x} \right| \leq \frac{\sin(pa)}{\sin a}.$$

By replacing k with $n-k$ if needed, we may assume $x \in [a, \pi/2]$. If $p = 1$, the claim is equality. Hence assume $p \geq 2$.

First suppose $x \leq \pi/p$. On $(0, \pi/p)$, the function $g_p(x) := \sin(px)/\sin x$ is decreasing, since

$$\frac{d}{dx} \log g_p(x) = p \cot(px) - \cot x < 0,$$

where the inequality follows from the fact that $t \mapsto t \cot t$ is strictly decreasing on $(0, \pi)$. Therefore,

$$\frac{\sin(px)}{\sin x} \leq \frac{\sin(pa)}{\sin a}.$$

Now suppose $x > \pi/p$. Since $x \in [\pi/p, \pi/2]$,

$$\left| \frac{\sin(px)}{\sin x} \right| \leq \frac{1}{\sin(\pi/p)}.$$

Also $pa \leq \pi/2$, and using $\sin u \geq 2u/\pi$ on $[0, \pi/2]$ and $\sin a \leq a$, we get

$$\sin(pa) \sin(\pi/p) \geq \frac{2pa}{\pi} \cdot \frac{2}{p} = \frac{4a}{\pi} \geq \sin a.$$

Thus,

$$\frac{1}{\sin(\pi/p)} \leq \frac{\sin(pa)}{\sin a},$$

which proves the desired bound.

Therefore

$$\Delta_m(k) \leq \frac{\sin^2(pa)}{\sin^2 a} \Delta_1(k) = \beta_m \Delta_1(k).$$

Rearranging gives the equivalent cosine form. For $n \geq 5$ and $m = 2, \dots, \lfloor n/2 \rfloor$, the above inequalities are strict unless $k = 0$ or $x = a$ up to the symmetry $x \mapsto \pi - x$, i.e., $k \in \{0, 1, n - 1\}$. \square

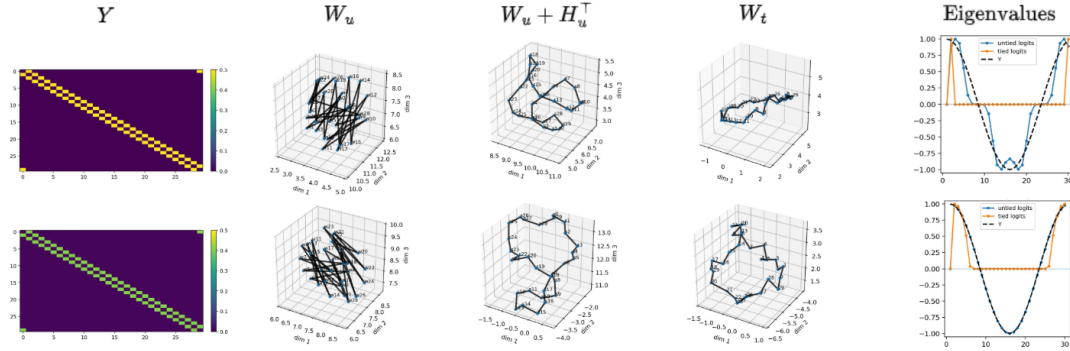


Figure 4. Same as Fig. 1 with UMAP instead of PCA for dimension reduction.