

How Cross-Entropy Learns Data Modes: Emergence and Implicit Bias in the Unconstrained Features Model

author names withheld

Under Review for the Workshop on High-dimensional Learning Dynamics, 2026

Abstract

A classical result for linear networks trained with mean-squared-error loss is that gradient flow learns the singular modes of the data sequentially, in order of importance. The precise mechanics of such sequential emergence under cross-entropy loss remain largely unknown. We study this question in a minimal nonconvex setting: a two-layer linear network with orthogonal inputs and step-imbalanced classes, equivalent to the unconstrained feature model used in neural collapse analyses. We derive a closed-form expression for the full regularization path, which exhibits sequential mode emergence, but with novel behavior: active singular values diverge, only normalized logits converge and can overshoot the limiting geometry. We then show that a related sequential picture holds for gradient flow under appropriate spectrally aligned initialization. Our analysis relies on a novel imbalance-adapted Hadamard basis in which softmax preserves a diagonal-plus-rank-one structure.

1. Introduction

The training dynamics of deep linear networks under mean-squared-error (MSE) loss have shaped much of the deep learning theory [2–4, 19, 33, 53]. A particularly influential example is the spectral-dynamics framework of Saxe et al. [43, 45], which showed that, under gradient flow, dominant data modes are learned sequentially and emerge in the network representations. This *sequential-learning* phenomenon has since informed numerous theoretical accounts [6, 24, 37, 55, 59, 60].

By contrast, the analogous picture for classification with cross-entropy (CE) loss remains largely incomplete. We still lack a precise characterization of whether representations exhibit comparable emergence, how such emergence unfolds, and which classifier CE selects at convergence. These questions are difficult even in simple models because of the softmax nonlinearity.

Motivated by this gap, we investigate CE convergence and dynamics in a two-layer linear neural network with orthogonal inputs under a step-imbalanced classification: K classes partitioned into equal numbers of majority and minority classes, with imbalance ratio $R > 1$. The choice of data and architecture is deliberately minimal. Step-imbalanced data is the minimal label-symmetry distortion that produces nontrivial mode structure: majority, majority-minority, and minority modes [60]. A two-layer linear architecture, unlike a single linear model, allows this structure to emerge sequentially. Orthogonality plays a double role: it simplifies the analysis, matching the original setting of Saxe et al. [44], and makes the model coincide with the unconstrained features model (UFM) used in neural collapse analyses, giving this otherwise minimal model a principled connection to deep networks, whose last-layer geometry the UFM has been shown to predict robustly. Our contributions are:

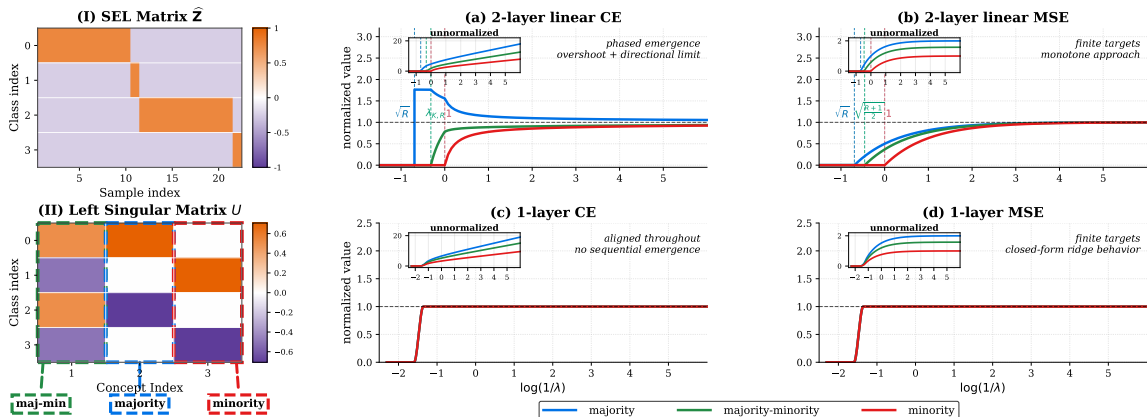


Figure 1: We study sequential learning in a deliberately minimal setting: step-imbalanced classification with $K/2$ majority and $K/2$ minority classes and imbalance ratio $R > 1$. (I) The SEL matrix $\hat{Z} = (I - \frac{1}{K}\mathbf{1}\mathbf{1}^\top)\mathbf{Y}$, a centering of the one-hot encoding matrix, has three distinct singular modes: majority–minority, majority, and minority. (II) These are reflected in the left singular matrix U [60]. Panels (a)–(d) show the evolution of logit singular values along the ridge-regularization path ($\lambda > 0$) for two-layer CE, two-layer MSE, one-layer CE, and one-layer MSE networks with orthogonal inputs. (a) We analytically solve the CE regularization path: modes activate in phases, but activation thresholds are more intricate than in MSE, unnormalized singular values diverge, and normalized coordinates converge only directionally and nonmonotonically, with possible overshoot. (b) The classical MSE picture [44, 46]: finite targets approached monotonically along sigmoidal curves. (c, d) One-layer CE and MSE paths show no sequential emergence. All coordinates are normalized so that the limiting profile equals one; insets show the unnormalized profiles.

- **Approximate softmax diagonalization.** We identify spectral coordinates adapted to the imbalanced label structure in which the softmax operator preserves a diagonal-plus-rank-one form. This approximate diagonalization reduces the matrix-valued CE problem to a lower-dimensional one.

- **Regularization path.** Using this reduction, we cast the ridge-regularized UFM as a four-dimensional optimization which we solve explicitly for *every* regularization strength λ . The solution exhibits a phase structure: as λ decreases, majority activates, then majority–minority, and finally the minority. Activation thresholds and active-mode magnitudes reduce to scalar nonlinear equations, which we give in closed form. From this we prove the vanishing-regularization limit, obtain a logarithmic convergence rate, and expose qualitative gaps with MSE sequential-learning; see Fig. 1.

- **Gradient-flow dynamics.** In the Appendix A, we perform an analogous study of gradient-flow under a spectrally aligned initialization compatible with the approximate diagonalization.

Closely related work. Garrod et al. [16] recently showed that the spectral dynamics framework for MSE dynamics [44], can be pushed through to CE despite the softmax nonlinearity. We extend this: (a) *Beyond pure diagonalization.* We show that an approximate diagonalization suffices to reduce the dynamics to low dimension. We expect this relaxation to be broadly useful, since prior exact diagonalization is restrictive outside symmetric settings. (b) *Phased mode emergence;* We give the first explicit analysis of sequential-emergence under CE since the balanced regime in [16] has no distinct singular modes. (c) *Regularization path.* We apply the spectral framework to the regularization path. This both connects the regularization path to gradient flow and yields the first analytic characterization of a CE regularization path beyond the closed-form MSE case.

Notation. Define $[N] := \{0, \dots, N - 1\}$ for $N \in \mathbb{N}$. Lowercase/uppercase boldface denote vectors/matrices; \otimes and \odot denote Kronecker and Hadamard products. e_i denotes the i -th standard basis vector; to reduce clutter, its dimension is kept implicit when clear from context. We reserve $\mathbf{f}_+ := (1, 0)^\top$ and $\mathbf{f}_- := (0, 1)^\top$ for basis in \mathbb{R}^2 . Throughout, subscripts $+/-$ flag quantities

associated with majorities/minorities. For matrices, $\bar{\mathbf{A}}$ denotes L2-normalization. For vectors, $\bar{\mathbf{a}}$ denotes L1-normalization. $\|\mathbf{A}\|_*$ is the nuclear norm. Inequalities on vectors apply entrywise. $\mathbb{S}(\cdot)$ denotes the softmax map. We denote the Sylvester–Hadamard matrix of order $K = 2^m, m \in \mathbb{N}$ by $\Phi_K \in \{\pm 1\}^{K \times K}$. With $\Phi_2 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$, it satisfies the recursion $\Phi_K = \Phi_M \otimes \Phi_2$. We write $\bar{\Phi}_K \in \{\pm 1\}^{K \times (K-1)}$ for the same matrix with its first (all-ones) column removed.

2. Problem setup and background

Data. We study imbalanced multiclass classification over $K = 2^m, m \in \mathbb{N}$ classes with n input-label pairs. Let $R \in \mathbb{N}$ denoting the *imbalance ratio*. Data comprises M majority classes with Rn' examples each and M minority classes with n' examples each, where $M := K/2$. Without restricting generality of the findings we set $n' = 1$ henceforth; thus $n = (R + 1)M$ total examples. We order classes so that even indices $(0, 2, \dots)$ correspond to majorities and odd indices $(1, 3, \dots)$ to minorities. The one-hot encoding matrix then takes the form

$$\mathbf{Y} := \mathbf{I}_M \otimes \begin{pmatrix} \mathbf{1}_R^\top & 0 \\ 0 & \mathbf{1}_R^\top \end{pmatrix} \in \mathbb{R}^{K \times n}. \quad (1)$$

When $R = 1$, data are balanced and $\mathbf{Y} = \mathbf{I}_K$. For later reference, we also define the *simplex-encoding label (SEL) matrix* $\hat{\mathbf{Z}} := (\mathbf{I}_K - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^\top) \mathbf{Y}$, i.e., \mathbf{Y} with each column centered.

Unconstrained Features Model (UFM). The UFM reduces DNN optimization to a bilinear problem over the classifier $\mathbf{W} \in \mathbb{R}^{K \times d}$ and a freely optimized embedding matrix $\mathbf{H} \in \mathbb{R}^{d \times n}$ (column \mathbf{h}_i is the last-layer representation of each example) [13, 35]. Throughout, we let $d \geq K$. Concretely, dropping the DNN’s parameterization of the embedding map, the UFM optimizes:

$$\mathcal{L}(\mathbf{W}\mathbf{H}) := \sum_{i \in [n]} -\log \left(\mathbf{e}_{y_i}^\top \mathbb{S}(\mathbf{W}\mathbf{h}_i) \right) = \sum_{i \in [n]} \log \left(1 + \sum_{c \neq y_i} e^{-(\mathbf{e}_{y_i} - \mathbf{e}_c)^\top \mathbf{W}\mathbf{h}_i} \right). \quad (\text{UFM})$$

Equivalently, (UFM) can be viewed as CE optimization of a two-layer linear network with orthogonal inputs $\mathbf{e}_i \in \mathbb{R}^n$, weight matrices \mathbf{H} and \mathbf{W} , and K outputs. The (UFM) loss has no finite minimizers: any direction pair $(\bar{\mathbf{W}}, \bar{\mathbf{H}})$ with strictly positive margin drives the loss to zero as parameter norm diverges. We study two complementary ways of making this precise: ridge regularization, which defines finite minimizers for each regularization strength, and unregularized gradient flow, which specifies an optimization trajectory. We focus on the former here and defer the latter to App. A.

Regularization path (RP). The ridge-regularized UFM,

$$(\mathbf{W}_\lambda, \mathbf{H}_\lambda) \in \arg \min_{\mathbf{W}, \mathbf{H}} \mathcal{L}(\mathbf{W}\mathbf{H}) + \lambda/2 \|\mathbf{W}\|_F^2 + \lambda/2 \|\mathbf{H}\|_F^2. \quad (\text{UFM}_\lambda)$$

models DNN training with weight decay [12, 15, 21, 22, 48, 63]. For $\lambda > 0$ it has finite minimizers, which diverge as $\lambda \rightarrow 0^+$, connecting the path to gradient flow [28, 41]. Two questions naturally arise: *convergence* (as $\lambda \rightarrow 0$) and *trajectory* (at finite λ). A key property that makes the RP more accessible than GF is that its convex SDP relaxation in logit space $\mathbf{Z} = \mathbf{W}\mathbf{H}$ is tight:

$$\mathbf{Z}_\lambda = \arg \min_{\mathbf{Z}} \mathcal{L}(\mathbf{Z}) + \lambda \|\mathbf{Z}\|_*. \quad (\text{Z-UFM}_\lambda)$$

Concretely, for all $\lambda > 0$, there exists orthogonal matrix $\mathbf{R} (\mathbf{R}^\top \mathbf{R} = \mathbf{I}_K)$ such that $\mathbf{W}_\lambda = \mathbf{U} \sqrt{\Sigma} \mathbf{R}^\top$ and $\mathbf{H}_\lambda = \mathbf{R} \sqrt{\Sigma} \mathbf{V}^\top$ where $\mathbf{Z}_\lambda = \mathbf{U} \Sigma \mathbf{V}^\top$ is the SVD of the solution to (Z-UFM $_\lambda$). In the balanced case ($R = 1$), it is known for all $\lambda < \lambda_{\max}$ (above which the solution is trivially zero) that $\mathbf{Z}_\lambda \propto \hat{\mathbf{Z}}$; thus, $(\mathbf{W}_\lambda, \mathbf{H}_\lambda)$ forms a simplex ETF: $\mathbf{W}_\lambda \mathbf{W}_\lambda^\top = \mathbf{H}_\lambda^\top \mathbf{H}_\lambda \propto \hat{\mathbf{Z}}$ [63]. But for $R > 1$, the picture changes. The solution varies nontrivially with λ , and although the normalized logits $\bar{\mathbf{Z}}_\lambda$ are known to become proportional to the SEL matrix $\hat{\mathbf{Z}}$ as $\lambda \rightarrow 0$, this proportionality fails at finite λ [22, 52]. Here, we ask: *precisely how do the singular modes of \mathbf{Z} enter and evolve as λ decreases?*

3. Technical Tool: Approximate Softmax Diagonalization

Both our RP and GF analyses rest on a technical observation about how the softmax operator interacts with the singular subspace of the SEL matrix $\widehat{\mathbf{Z}}$. To formalize this, recall $K = 2^m$, $m \in \mathbb{N}$ and let $\mathbf{U} \in \mathbb{R}^{K \times (K-1)}$ and $\mathbf{V} \in \mathbb{R}^{n \times (K-1)}$ be the following matrices, parameterized by the Hadamard basis $\overline{\Phi}_M$, whose columns form orthonormal bases for, the left and right singular subspaces of $\widehat{\mathbf{Z}}$:

$$\mathbf{U} := \left[\frac{1}{\sqrt{2M}} \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \frac{1}{\sqrt{M}} \overline{\Phi}_M \otimes \mathbf{I}_2 \right], \quad \mathbf{V} := \left[\frac{1}{\sqrt{M(R+1)}} \mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ -1 \end{pmatrix} \quad \frac{1}{\sqrt{M}} \overline{\Phi}_M \otimes \begin{pmatrix} \frac{1}{\sqrt{R}} & \frac{0}{\sqrt{R}} \\ 0 & 1 \end{pmatrix} \right].$$

It can then be checked that $\widehat{\mathbf{Z}} = \mathbf{U} \text{diag}(\boldsymbol{\sigma}) \mathbf{V}^\top$ for $\boldsymbol{\sigma} = [\sqrt{\frac{R+1}{2}}, \sqrt{R} \mathbf{1} \otimes \mathbf{f}_+ + \mathbf{1} \otimes \mathbf{f}_-]$. Thus, for $R > 1$, $\widehat{\mathbf{Z}}$ has *three* distinct ordered singular modes: a **majority mode** at value \sqrt{R} , a **majority-minority mode** at value $\sqrt{(R+1)/2}$, and, a **minority mode** at value 1.

We require a rank-1 extension of the basis. Set $\widetilde{\mathbf{V}} := [\widetilde{\mathbf{v}}, \mathbf{V}]$, with $\widetilde{\mathbf{v}} := \frac{1}{\sqrt{MR(R+1)}} \mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ R \end{pmatrix}$, and introduce the operator $\text{diag}_{+1}(\cdot) : \mathbb{R} \times \mathbb{R} \times \mathbb{R}^{M-1} \times \mathbb{R}^{M-1} \rightarrow \mathbb{R}^{(K-1) \times K}$ defined by

$$\text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) := [\alpha \mathbf{e}_0 \quad \mathbf{A}], \quad \mathbf{A} := \text{diag}(a, \mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-).$$

This packages a diagonal matrix \mathbf{A} together with a rank-1 column $\alpha \mathbf{e}_0 \in \mathbb{R}^{K-1}$.

Theorem 1 (Softmax Diagonalization up to Rank-1) For $\alpha, a \in \mathbb{R}$, $\mathbf{a}_\pm \in \mathbb{R}^{M-1}$ it holds

$$\mathbf{Y} - \mathbb{S} \left(\mathbf{U} \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \widetilde{\mathbf{V}}^\top \right) = \mathbf{U} \text{diag}_{+1}(\beta, b, \mathbf{b}_+, \mathbf{b}_-) \widetilde{\mathbf{V}}^\top. \quad (2)$$

where the map $(\beta, b, \mathbf{b}_+, \mathbf{b}_-) := \mathcal{F}((\alpha, a, \mathbf{a}_+, \mathbf{a}_-))$ is given in closed-form in Eq. (13) in App. D.1.

Since \mathbf{Y} itself decomposes as $\mathbf{Y} = \mathbf{U} \text{diag}_{+1}([0, \boldsymbol{\sigma}]) \widetilde{\mathbf{V}}^\top$, Thm. 1 states that softmax preserves this decomposition. This is a consequence of our specific Hadamard parameterization of \mathbf{U} , $\widetilde{\mathbf{V}}$, which is the key novelty of the result. The Hadamard property we rely on is *phase-shift under permutations*: $\Pi_{i \oplus j} \overline{\Phi}_M \mathbf{e}_i = (-1)^{i^\top j} \overline{\Phi}_M \mathbf{e}_i$; Π_i are XOR (\oplus) permutations and i, j binary representations.

4. Regularization path

We analytically solve (UFM $_\lambda$) for every $\lambda > 0$. From Sec. 2, it suffices to solve the convex relaxation (Z-UFM $_\lambda$) in logit space. While this is convex and can be solved numerically, there is no a priori reason to expect an analytic solution. We show that one nevertheless exists in two steps.

First, we use Thm. 1 to reduce the SDP to a strictly convex problem in only four scalars.

Theorem 2 (SDP-Vectorization) For $\lambda > 0$, the unique minimizer of (Z-UFM $_\lambda$) decomposes as $\mathbf{Z}_\lambda = \mathbf{U} \text{diag}_{+1}(\alpha_\lambda, a_\lambda, \mathbf{a}_{+,\lambda}, \mathbf{a}_{-,\lambda}) \widetilde{\mathbf{V}}^\top$, where $\alpha_\lambda, a_\lambda, \mathbf{a}_{+,\lambda}, \mathbf{a}_{-,\lambda}$ are unique minimizers of

$$\min_{\alpha, a, \mathbf{a}_\pm} \mathcal{L}(\mathbf{U} \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \widetilde{\mathbf{V}}^\top) + \lambda(\sqrt{\alpha^2 + a^2} + \|\mathbf{a}_+\|_1 + \|\mathbf{a}_-\|_1). \quad (3)$$

Moreover, $\mathbf{a}_{\pm,\lambda} = a_{\pm,\lambda} \mathbf{1}_{M-1}$. Hence (Z-UFM $_\lambda$) reduces to a four-dimensional optimization.

The crux of the proof is showing that the global minimizer of (Z-UFM $_\lambda$) lies on the $(\mathbf{U}, \widetilde{\mathbf{V}})$ subspace for all $\lambda > 0$, and this is where Thm. 1 enters. The optimality condition for (Z-UFM $_\lambda$) is $\mathbf{Y} - \mathbb{S}(\mathbf{Z}) \in \lambda \partial \|\mathbf{Z}\|_*$. For candidate $\mathbf{Z} = \mathbf{U} \text{diag}_{+1}(\mathbf{a}) \widetilde{\mathbf{V}}^\top$, Thm. 1 rewrites the LHS as

$U \text{diag}_{+1}(\mathcal{F}(\mathbf{a})) \tilde{\mathbf{V}}^\top$, so the condition becomes $U \text{diag}_{+1}(\mathcal{F}(\mathbf{a})) \tilde{\mathbf{V}}^\top \in \lambda \partial \|U \text{diag}_{+1}(\mathbf{a}) \tilde{\mathbf{V}}^\top\|_*$, which by orthogonality of U , $\tilde{\mathbf{V}}$ is equivalent to the optimality condition of (3).

Second, we solve the optimality conditions of (3) explicitly. To state the result, we reparameterize (α, a) in terms of $(\gamma, \delta) := \sqrt{\frac{2}{R+1}}(a + \frac{1}{\sqrt{R}}\alpha, a - \alpha\sqrt{R})$. For $Q(x, \theta) := \log(x^{M-1}((1+\theta)x + \theta))$, $\theta_+ := M[\sqrt{R}\lambda^{-1} - 1]_+$, and $\theta_- := M[\lambda^{-1} - 1]_+$, define $\gamma_\lambda := Q(s_\lambda, \theta_+)$, $\delta_\lambda := Q(r_\lambda, \theta_-)$. Finally, let r_λ solve a scalar nonlinear equation and s_λ be a closed-form function of it, both given in App. E.3.

Theorem 3 (Analytic solution) *Let $\mathbf{a}_\lambda = (\alpha_\lambda, a_\lambda, a_{+,\lambda}, a_{-,\lambda})$ the minimizer of (3) and $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ the unique solution of $Q(2\lambda/\sqrt{2\lambda^2 - 1} - 1, \theta_+) = 0$. Then, \mathbf{a}_λ satisfies:*

- (i) For $\lambda \geq \sqrt{R}$, $\mathbf{a}_\lambda = (0, 0, 0, 0)$.
 - (ii) For $\lambda \in [\lambda_{K,R}, \sqrt{R})$, $\mathbf{a}_\lambda = (0, 0, a_{+,\lambda}, 0)$, where $\gamma_\lambda = \delta_\lambda = 0$ and $a_{+,\lambda} = -M\sqrt{R} \log s_\lambda$.
 - (iii) For $\lambda \in [1, \lambda_{K,R})$, $\mathbf{a}_\lambda = (\alpha_\lambda, a_\lambda, a_{+,\lambda}, 0)$, where $a_{+,\lambda} = \sqrt{R}(\gamma_\lambda - M \log s_\lambda)$.
 - (iv) For $\lambda < 1$, $\mathbf{a}_\lambda = (\alpha_\lambda, a_\lambda, a_{+,\lambda}, a_{-,\lambda})$, where $a_{+,\lambda} = \sqrt{R}(\gamma_\lambda - M \log s_\lambda)$, $a_{-,\lambda} = \delta_\lambda - M \log s_\lambda$.
- Thus, as λ decreases, the majority mode ($a_{+,\lambda}$) appears first, then the majority–minority mode $(\alpha_\lambda, a_\lambda)$, and finally the minority mode ($a_{-,\lambda}$).

Emergence thresholds. The theorem identifies the λ s at which modes first become nonzero. Majority and minority activate at $\lambda_+ := \sqrt{R}$, $\lambda_- := 1$, which coincide with the largest, smallest nonzero singular values of the SEL matrix, and do not vary with K . As expected, increased imbalance R widens the ratio λ_+/λ_- . The intermediate threshold $\lambda_{K,R}$, at which majority–minority mode emerges, is decreasing in K and increasing in R (but slower than λ_+). Although it has no closed form, it solves a scalar nonlinear equation. App. E.4 and Fig. 2 for details and visualization.

Convergence. More importantly, the theorem identifies the exact evolution of the parameters within each phase in terms of a single scalar nonlinear equation. In App. E.2, we show that $a_{\pm,\lambda} \geq 0$ for all λ , with strict positivity once activation occurs. Moreover, the active coordinates are decreasing in λ diverge as $\lambda \rightarrow 0$. However, their normalized direction converges to the SEL singular direction.

Theorem 4 (Convergence to SEL) *With \mathbf{a}_λ as in Thm. 3, as $\lambda \rightarrow 0$, the normalized spectral modes $\bar{\mathbf{a}}_\lambda$ converge to $\bar{\mathbf{a}}_\infty$ with $\bar{\mathbf{a}}_\infty = (0, \sqrt{(R+1)/2}, \sqrt{R}, 1)$. Consequently, the logit singular values $\sigma_\lambda := (\sqrt{a_\lambda^2 + \alpha_\lambda^2}, a_{+,\lambda}, a_{-,\lambda})$ converge in direction to the SEL singular values $\sigma = (\sqrt{(R+1)/2}, \sqrt{R}, 1)$. Convergence is at rate $O(\log(\lambda^{-1}))$.*

Unlike prior vanishing-regularization convergence results (e.g. [28, 50–52]), which to our knowledge all rely on a non-constructive proof by contradiction following [40, 41], the result above is explicit (it follows directly from the phase-wise formulas of Thm. 3), and gives an explicit convergence rate.

Mode emergence without finite targets. Our fine-grained analysis highlights a key difference between the CE and the classical MSE sequential-learning picture. Under MSE, the target matrix has finite singular values; each mode follows a sigmoidal trajectory, leaving the initialization scale and growing from below monotonically toward its finite target value [17, 44]. In CE, by contrast, the active singular values have no finite targets: they diverge as $\lambda \rightarrow 0$, and only their normalized direction converges. Moreover, earlier activation does *not* mean monotone growth from below toward the limiting normalized value: Majority activates first but approaches its normalized SEL value from above; majority-minority mode activates next and approaches from below; and minority activates last, approaches from below and can approach faster than the majority-minority mode depending on K and R . This contrast with the MSE is illustrated in Fig. 1 and Fig. 3.

References

- [1] George Andriopoulos, Zixuan Dong, Li Guo, Zifan Zhao, and Keith Ross. The prevalence of neural collapse in neural multivariate regression. *arXiv preprint arXiv:2409.04180*, 2024.
- [2] Sanjeev Arora, Nadav Cohen, Noah Golowich, and Wei Hu. A convergence analysis of gradient descent for deep linear neural networks. *arXiv preprint arXiv:1810.02281*, 2018.
- [3] Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 244–253. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/arora18a.html>.
- [4] Francis Bach. Closed-form dynamics beyond quadratics, mar 2026. URL <https://francisbach.com/closed-form-dynamics/>.
- [5] Kiril Bangachev, Guy Bresler, Iliyas Noman, and Yury Polyanskiy. Global minimizers of sigmoid contrastive loss. *arXiv preprint arXiv:2509.18552*, 2025.
- [6] Ioannis Bantzis, James B Simon, and Arthur Jacot. Saddle-to-saddle dynamics in deep relu networks: Low-rank bias in the first saddle escape. *arXiv preprint arXiv:2505.21722*, 2025.
- [7] Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017. doi: 10.1137/1.9781611974997.
- [8] Tina Behnia, Ganesh Ramachandra Kini, Vala Vakilian, and Christos Thrampoulidis. On the implicit geometry of cross-entropy parameterizations for label-imbalanced data. In *International Conference on Artificial Intelligence and Statistics*, pages 10815–10838. PMLR, 2023.
- [9] Rajendra Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
- [10] Nadav Cohen and Noam Razin. Lecture notes on linear neural networks: A tale of optimization and generalization in deep learning. *arXiv preprint arXiv:2408.13767*, 2024.
- [11] Hien Dang, Tan Nguyen, Tho Tran, Hung Tran, and Nhat Ho. Neural collapse in deep linear network: From balanced to imbalanced data. *arXiv preprint arXiv:2301.00437*, 2023.
- [12] Hien Dang, Tho Tran, Tan Nguyen, and Nhat Ho. Neural collapse for cross-entropy class-imbalanced learning with unconstrained relu feature model. *arXiv preprint arXiv:2401.02058*, 2024.
- [13] Cong Fang, Hangfeng He, Qi Long, and Weijie J Su. Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training. *Proceedings of the National Academy of Sciences*, 118(43), 2021.
- [14] Quinn LeBlanc Fisher, Haoming Meng, and Vardan Papyan. Pushing boundaries: Mixup’s influence on neural collapse. In *The Twelfth International Conference on Learning Representations*, 2024.

- [15] Connall Garrod and Jonathan P Keating. The persistence of neural collapse despite low-rank bias: An analytic perspective through unconstrained features. *arXiv preprint arXiv:2410.23169*, 2024.
- [16] Connall Garrod, Jonathan P Keating, and Christos Thrampoulidis. Diagonalizing the softmax: Hadamard initialization for tractable cross-entropy dynamics. *arXiv preprint arXiv:2512.04006*, 2025.
- [17] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/f39ae9ff3a81f499230c4126e01f421b-Paper.pdf.
- [18] Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- [19] Suriya Gunasekar, Blake E. Woodworth, Srinadh Bhojanapalli, Behnam Neyshabur, and Nathan Srebro. Implicit regularization in matrix factorization. *2018 Information Theory and Applications Workshop (ITA)*, pages 1–10, 2017. URL <https://api.semanticscholar.org/CorpusID:3909231>.
- [20] Li Guo, George Andriopoulos, Zifan Zhao, Shuyang Ling, Zixuan Dong, and Keith Ross. Cross entropy versus label smoothing: A neural collapse perspective. *arXiv preprint arXiv:2402.03979*, 2024.
- [21] XY Han, Vardan Papyan, and David L Donoho. Neural collapse under mse loss: Proximity to and dynamics on the central path. *arXiv preprint arXiv:2106.02073*, 2021.
- [22] Wanli Hong and Shuyang Ling. Neural collapse for unconstrained feature model under cross-entropy loss with imbalanced data. *arXiv preprint arXiv:2309.09725*, 2023.
- [23] Wanli Hong and Shuyang Ling. Beyond unconstrained features: Neural collapse for shallow neural networks with general data. *arXiv preprint arXiv:2409.01832*, 2024.
- [24] Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv preprint arXiv:2106.15933*, 2021.
- [25] Arthur Jacot, Peter Sůkeník, Zihan Wang, and Marco Mondelli. Wide neural networks trained with weight decay provably exhibit neural collapse. *arXiv preprint arXiv:2410.04887*, 2024.
- [26] Wenlong Ji, Yiping Lu, Yiliang Zhang, Zhun Deng, and Weijie J Su. An unconstrained layer-peeled perspective on neural collapse. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=WZ3yjh8coDg>.
- [27] Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. *arXiv preprint arXiv:1810.02032*, 2018.

- [28] Ziwei Ji, Miroslav Dudík, Robert E Schapire, and Matus Telgarsky. Gradient descent follows the regularization path for general losses. In *Conference on Learning Theory*, pages 2109–2136. PMLR, 2020.
- [29] Jiachen Jiang, Jinxin Zhou, Peng Wang, Qing Qu, Dustin Mixon, Chong You, and Zihui Zhu. Generalized neural collapse for a large number of classes. *arXiv preprint arXiv:2310.05351*, 2023.
- [30] Ganesh Ramachandra Kini, Orestis Paraskevas, Samet Oymak, and Christos Thrampoulidis. Label-imbalanced and group-sensitive classification under overparameterization. *Advances in Neural Information Processing Systems*, 34:18970–18983, 2021.
- [31] Ganesh Ramachandra Kini, Vala Vakilian, Tina Behnia, Jaidev Gill, and Christos Thrampoulidis. Symmetric neural-collapse representations with supervised contrastive loss: The impact of relu and batching. In *The Twelfth International Conference on Learning Representations*, 2024.
- [32] Pengyu Li, Xiao Li, Yutong Wang, and Qing Qu. Neural collapse in multi-label learning with pick-all-label loss. *arXiv preprint arXiv:2310.15903*, 2023.
- [33] Zhiyuan Li, Yuping Luo, and Kaifeng Lyu. Towards resolving the implicit bias of gradient descent for matrix factorization: Greedy low-rank learning. *arXiv preprint arXiv:2012.09839*, 2020.
- [34] Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2020.
- [35] Dustin G Mixon, Hans Parshall, and Jianzong Pi. Neural collapse with unconstrained features. *arXiv preprint arXiv:2011.11619*, 2020.
- [36] Jean-Jacques Moreau. Fonctions convexes duales et points proximaux dans un espace hilbertien. *CR Acad. Sci. Paris Sér. A Math*, 255:2897–2899, 1962.
- [37] Tiberiu Musat. The geometry of grokking: Norm minimization on the zero-loss manifold. *arXiv preprint arXiv:2511.01938*, 2025.
- [38] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020.
- [39] Mee Young Park and Trevor Hastie. L 1-regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 69(4): 659–677, 2007.
- [40] Saharon Rosset, Ji Zhu, and Trevor Hastie. Margin maximizing loss functions. In *NIPS*, pages 1237–1244, 2003.
- [41] Saharon Rosset, Ji Zhu, and Trevor J. Hastie. Margin maximizing loss functions. In *NIPS*, 2003.
- [42] Saharon Rosset, Ji Zhu, and Trevor Hastie. Boosting as a regularized path to a maximum margin classifier. *Journal of Machine Learning Research*, 5(Aug):941–973, 2004.

- [43] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *CoRR*, abs/1312.6120, 2013. URL <https://api.semanticscholar.org/CorpusID:17272965>.
- [44] Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [45] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019. doi: 10.1073/pnas.1820226116. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1820226116>.
- [46] Andrew M Saxe, James L McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019.
- [47] Peter Sůkeník, Marco Mondelli, and Christoph H Lampert. Deep neural collapse is provably optimal for the deep unconstrained features model. In *Advances in Neural Information Processing Systems*, volume 36, pages 52991–53024. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/a60c43ba078b723d3d517d28c50ded4c-Paper-Conference.pdf.
- [48] Peter Sůkeník, Marco Mondelli, and Christoph Lampert. Neural collapse versus low-rank bias: Is deep neural collapse really optimal? *arXiv preprint arXiv:2405.14468*, 2024.
- [49] Peter Sůkeník, Christoph H Lampert, and Marco Mondelli. Neural collapse is globally optimal in deep regularized resnets and transformers. *arXiv preprint arXiv:2505.15239*, 2025.
- [50] Davoud Ataee Tarzanagh, Yingcong Li, Christos Thrampoulidis, and Samet Oymak. Transformers as support vector machines, 2023.
- [51] Davoud Ataee Tarzanagh, Yingcong Li, Xuechen Zhang, and Samet Oymak. Max-margin token selection in attention mechanism, 2023.
- [52] Christos Thrampoulidis, Ganesh R Kini, Vala Vakilian, and Tina Behnia. Imbalance trouble: Revisiting neural-collapse geometry. *arXiv preprint arXiv:2208.05512*, 2022.
- [53] Zhenfeng Tu, Santiago Aranguri, and Arthur Jacot. Mixed dynamics in linear networks: Unifying the lazy and active regimes. In *Advances in Neural Information Processing Systems*, volume 37, pages 106059–106104. Curran Associates, Inc., 2024.
- [54] Gal Vardi, Ohad Shamir, and Nathan Srebro. On margin maximization in linear and relu networks. *arXiv preprint arXiv:2110.02732*, 2021.
- [55] Bhavya Vasudeva, Puneesh Deora, Yize Zhao, Vatsal Sharan, and Christos Thrampoulidis. How muon’s spectral design benefits generalization: A study on imbalanced data. *arXiv preprint arXiv:2510.22980*, 2025.
- [56] Diyuan Wu and Marco Mondelli. Neural collapse beyond the unconstrained features model: Landscape, dynamics, and generalization in the mean-field regime. *arXiv preprint arXiv:2501.19104*, 2025.

- [57] Robert Wu and Vardan Papyan. Linguistic collapse: Neural collapse in (large) language models. *arXiv preprint arXiv:2405.17767*, 2024.
- [58] Yedi Zhang, Andrew M. Saxe, and Peter E. Latham. Saddle-to-saddle dynamics explains a simplicity bias across neural network architectures. *ArXiv*, abs/2512.20607, 2025. URL <https://api.semanticscholar.org/CorpusID:284133040>.
- [59] Yedi Zhang, Aaditya K Singh, Peter E Latham, and Andrew Saxe. Training dynamics of in-context learning in linear attention. *arXiv preprint arXiv:2501.16265*, 2025.
- [60] Yize Zhao and Christos Thrampoulidis. On the geometry of semantics in next-token prediction. *arXiv preprint arXiv:2505.08348*, 2025.
- [61] Yize Zhao, Tina Behnia, Vala Vakilian, and Christos Thrampoulidis. Implicit geometry of next-token prediction: From language sparsity patterns to model representations. *arXiv preprint arXiv:2408.15417*, 2024.
- [62] Jinxin Zhou, Chong You, Xiao Li, Kangning Liu, Sheng Liu, Qing Qu, and Zhihui Zhu. Are all losses created equal: A neural collapse perspective. *arXiv preprint arXiv:2210.02192*, 2022.
- [63] Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. *Advances in Neural Information Processing Systems*, 34, 2021.

Contents

1	Introduction	1
2	Problem setup and background	3
3	Technical Tool: Approximate Softmax Diagonalization	4
4	Regularization path	4
A	Gradient-flow Path	12
A.1	Background	12
A.2	Results	12
A.2.1	Logit dynamics under balancedness	13
A.2.2	Emergence	13
A.2.3	Convergence	14
B	Related Work and Limitations	15
B.1	Detailed Related Work	15
B.2	Limitations	18
C	Notation and basic facts	18
D	Theorem 1 Proof and Discussions	20
D.1	The Mapping of Theorem 1 and Proof Outline	20
D.2	Proof of Theorem 1	21
D.3	Simplified form of expressions	26
E	Proofs for Regularization-path	27
E.1	Proof of Theorem 2 and basic properties of the solution	27
E.1.1	Reduction to four variables	31
E.2	Proof of Theorem 3: Analytic solution	34
E.3	Analytic Solution of Theorem 3	38
E.4	Behavior of the threshold $\lambda_{K,R}$	41
E.5	Asymptotic behavior of RP solution and Proof of Theorem 4	42
E.6	Alternative proof of regularization-path limit	43
F	Proofs for Gradient-flow path	45
F.1	Proof of Theorem 5	45
F.2	Proof of Theorem 9: Majority and Minority tend to Uniformity	46
F.3	Proof of Theorem 7	47
F.3.1	Auxiliary lemmas	48
F.3.2	Proof of Proposition 41	50
F.3.3	Proof of Proposition 42	51
F.4	Lyapunov function counterexample	53

G Numerical Results	54
G.1 Regularization Path Experiments	54
G.1.1 Regularization Path Comparison for Different K and R	55
G.2 Gradient-flow experiments	58
G.2.1 Setup	58
G.2.2 Results	59
G.2.3 Preservation of block-constant spectral coordinates	60

Appendix A. Gradient-flow Path

A.1. Background

For the continuous gradient flow (GF):

$$\dot{\mathbf{W}}_t = (\mathbf{Y} - \mathbb{S}(\mathbf{W}_t \mathbf{H}_t)) \mathbf{H}_t^\top, \quad \dot{\mathbf{H}}_t = \mathbf{W}_t^\top (\mathbf{Y} - \mathbb{S}(\mathbf{W}_t \mathbf{H}_t)), \quad (4)$$

we ask: (i) *Convergence*: where do the normalized iterates $(\overline{\mathbf{W}}_t, \overline{\mathbf{H}}_t)$ converge as $t \rightarrow \infty$? (ii) *Trajectory*: what structural properties does the optimization path exhibit at finite times? Although simplified, (UFM) remains non-convex, so even convergence of the GF iterates is non-trivial. While it is known that $(\overline{\mathbf{W}}_t, \overline{\mathbf{H}}_t)$ must converge to a KKT point of a max-margin problem [26, 34]

$$\min_{\mathbf{W}, \mathbf{H}} \|\mathbf{W}\|_F^2 + \|\mathbf{H}\|_F^2 \quad \text{sub. to} \quad (\mathbf{e}_{y_i} - \mathbf{e}_c)^\top \mathbf{W} \mathbf{h}_i \geq 1, \quad \forall i, c \neq y_i, \quad (5)$$

such a characterization is *not* explicit. There are many KKT points, each with different structural properties, and, in general, convergence can be either at global or local minima and even at saddles of the max-margin problem in Eq. (5) [54].

In the binary case ($K = 2$), [27] shows that the dynamics asymptotically satisfy the balancedness condition $\mathbf{W}_t^\top \mathbf{W}_t = \mathbf{H}_t \mathbf{H}_t^\top$, and [54] establishes convergence to a global optimum of (5); but $K = 2$ is both restrictive and trivial in our setting, since the only KKT point is the global minimizer. In the multiclass case, Garrod et al. [16] recently introduced a Hadamard spectral initialization under which the dynamics remain on a fixed subspace, leaving only the singular values to evolve and making analysis tractable. Their construction, however, applies only to $R = 1$, and the balanced setting has no distinct modes, so mode emergence cannot be studied there. We extend this framework to the imbalanced setting using the diagonal-plus-rank-one softmax structure of Sec. 3.

A.2. Results

To analyze the GF dynamics (4), we follow the program of [16, 17, 44], and use a spectral initialization. Concretely, as dictated by Thm. 1, we initialize logits on the $\mathbf{U}, \tilde{\mathbf{V}}$ subspace. This initialization is motivated by three observations: (a) The conjectured limit of GF in [52], namely the SEL matrix $\tilde{\mathbf{Z}}$, admits this decomposition; (b) The entire regularization path lies exactly on this subspace, as shown in Sec. 4; and (c) Allowing us to leverage Theorem 1 ensures that under this initialization, the logits remain in this decomposition at *all* times, turning the matrix into a vector ODE .

Theorem 5 (Invariance of spectral subspace) *Initialize the GF updates in Eq. (4) such that $\mathbf{W}_0 = \mathbf{U} \text{diag}(w_0, \mathbf{w}_{+,0}, \mathbf{w}_{-,0}) \mathbf{R}^\top$ and $\mathbf{H}_0 = \mathbf{R} \text{diag}_{+1}(\alpha_{H,0}, h_0, \mathbf{h}_{+,0}, \mathbf{h}_{-,0}) \tilde{\mathbf{V}}^\top$ for a partial orthogonal matrix $\mathbf{R} \in \mathbb{R}^{d \times (K-1)}$. Then, for all times t , it holds that*

$$\mathbf{W}_t = \mathbf{U} \text{diag}(w_t, \mathbf{w}_{+,t}, \mathbf{w}_{-,t}) \mathbf{R}^\top \quad \text{and} \quad \mathbf{H}_t = \mathbf{R} \text{diag}_{+1}(\alpha_{H,t}, h_t, \mathbf{h}_{+,t}, \mathbf{h}_{-,t}) \tilde{\mathbf{V}}^\top, \quad (6a)$$

where the diagonal entries evolve as (suppressing t -dependence henceforth when clear from context):

$$\dot{w} = \beta\alpha_H + bh, \quad \dot{\alpha}_H = \beta w, \quad \dot{h} = bw, \quad \dot{\mathbf{w}}_{\pm} = \mathbf{b}_{\pm} \odot \mathbf{h}_{\pm}, \quad \dot{\mathbf{h}}_{\pm} = \mathbf{b}_{\pm} \odot \mathbf{w}_{\pm}, \quad (7)$$

for the mapping $[\beta, b, \mathbf{b}_+, \mathbf{b}_-] = \mathcal{F}([\alpha_H w, wh, \mathbf{w}_+ \odot \mathbf{h}_+, \mathbf{w}_- \odot \mathbf{h}_-])$ as specified in Theorem 1.

A.2.1. LOGIT DYNAMICS UNDER BALANCEDNESS

As with the regularization-path, we now focus on the logit matrix $\mathbf{Z}_t = \mathbf{W}_t \mathbf{H}_t$, which under the parameterization in Eq. (6) is fully described by the spectral coordinates $\alpha = \alpha_H w$, $a = wh$, and $\mathbf{a}_{\pm} = \mathbf{w}_{\pm} \odot \mathbf{h}_{\pm}$. We focus on their evolution. To simplify their dynamics, we adopt the standard balancedness condition at initialization, $\mathbf{W}_0^{\top} \mathbf{W}_0 = \mathbf{H}_0 \mathbf{H}_0^{\top}$ [10, 16, 17, 43, 45], which is preserved by the dynamics; we leave the general case to future work. It will also be convenient to track the maj-min mode in the (γ, δ) parameterization of Eq. (14), motivated by its appearance in the regularization-path.

Corollary 6 *Suppose the spectral initialization of Theorem 5 additionally satisfies the balancedness conditions $w_0 = (\alpha_{H,0}^2 + h_0^2)^{1/2}$ and $\mathbf{w}_{\pm,0} = \mathbf{h}_{\pm,0}$. These conditions are preserved by the dynamics, and the logit matrix takes the form $\mathbf{Z}_t = \mathbf{U} \text{diag}_{+1}(\alpha_t, a_t, \mathbf{a}_{+,t}, \mathbf{a}_{-,t}) \tilde{\mathbf{V}}^{\top}$ for all $t \geq 0$. With (γ_t, δ_t) defined as in Eq. (14) (a one-to-one transformation of (α_t, a_t)), the logit parameters evolve as*

$$\dot{\gamma} = \frac{(2R\gamma^2 + \delta^2) X_+ + \gamma\delta X_-}{\sqrt{2(R\gamma^2 + \delta^2)}}, \quad \dot{\delta} = \frac{(R\gamma^2 + 2\delta^2) X_- + R\gamma\delta X_+}{\sqrt{2(R\gamma^2 + \delta^2)}}, \quad \dot{\mathbf{a}}_{\pm} = 2 \mathbf{a}_{\pm} \odot \mathbf{b}_{\pm}. \quad (8)$$

We henceforth focus on the ODE in (8) with (strictly) positive initialization for all parameters. Because $\mathbf{b}_{\pm} \geq \mathbf{0}$ and $X_{\pm} \geq 0$, this directly yields that all the parameters are increasing in time.

A.2.2. EMERGENCE

In analogy with the regularization path, one expects the modes to appear sequentially, with the majority mode emerging first. This is indeed the case, although the statement is necessarily different: along GF the modes are not exactly zero before activation (indeed, the origin is a saddle point!), but rather start at a small scale $O(\varepsilon)$. We therefore compare the times at which the different modes first reach a fixed scale $\eta > \varepsilon$. The analysis is also more delicate than for the regularization path, since the reduced ODE does not admit an explicit solution. We give some insights below.

Theorem 7 (Majority first) *Fix $R > 1$. Consider the reduced logit dynamics (8) from positive initialization of size $O(\varepsilon)$, with at least one majority coordinate initialized at order ε . Then, for every sufficiently small fixed threshold $\eta > 0$, and all sufficiently small $\varepsilon \in (0, \eta)$, the majority mode \mathbf{a}_+ reaches scale η before both the minority mode \mathbf{a}_- and the maj-min mode $a_0 := \sqrt{(R\gamma^2 + \delta^2)}/2$.*

See App. F.3 for complete statement and proof. We now give the main intuition. Note that the majority and minority evolution $\dot{\mathbf{a}}_{\pm} = 2 \mathbf{a}_{\pm} \odot \mathbf{b}_{\pm}$ is coupled through γ, δ , since \mathbf{b}_{\pm} depends on \mathbf{x}_{\pm} , which in turn depends on γ, δ , and γ, δ are coupled in (8). For emergence, however, it suffices to control \mathbf{b}_{\pm} uniformly in the small-coordinate regime. A key observation is that by setting $\mathbf{y}_+ := (M\sqrt{R})^{-1} \bar{\Phi} \mathbf{a}_+$ and $\mathbf{y}_- := M^{-1} \bar{\Phi} \mathbf{a}_-$, it holds $|\frac{X_{\pm}}{R} \bar{\Phi}^{\top} e^{\mathbf{x}_{\pm}}| \leq \frac{1}{\mathbf{1}^{\top} e^{\mathbf{y}_{\pm}}} |\bar{\Phi}^{\top} e^{\mathbf{y}_{\pm}}|$; thus the dependence on γ, δ drops out of this bound. Hence, when $\mathbf{a}_{\pm} = O(\varepsilon)$, $\mathbf{b}_{\pm} = \sqrt{R} \mathbf{1} + O(\varepsilon)$ and

$\mathbf{b}_- = \mathbf{1} + O(\varepsilon)$, implying that, near the origin, $\dot{\mathbf{a}}_+ \approx 2\sqrt{R}\mathbf{a}_+$ while $\dot{\mathbf{a}}_- \approx 2\mathbf{a}_-$, so the majority mode grows faster than the minority mode for $R > 1$. It remains to compare with the majority–minority mode $a_0 := \sqrt{(R\gamma^2 + \delta^2)}/2$. From (8), $\dot{a}_0 = R\gamma X_+ + \delta X_-$. Thus, in the $O(\varepsilon)$ -regime, where $X_\pm = 1 + O(\varepsilon)$, it holds $\dot{a}_0 \approx R\gamma + \delta$. Unlike the majority/minority equations, this does not close in terms of a_0 alone without controlling the relative sizes of γ and δ . A simple Cauchy–Schwarz bound gives $\dot{a}_0 \approx R\gamma + \delta \leq 2\sqrt{(R+1)/2} \cdot a_0$. Thus a_0 grows at rate at most $2\sqrt{(R+1)/2}$, strictly below the majority rate $2\sqrt{R}$ for $R > 1$. This upper bound is enough to prove Thm. 7, but it is not tight: the actual growth of a_0 depends on the evolving ratio of γ and δ . Thus, the majority–minority mode is not governed by a closed scalar growth equation even in the linearized regime. This is the ODE analogue of the regularization-path picture: the majority and minority activation thresholds are fixed to \sqrt{R} and 1, while $\lambda_{K,R}$ depends nontrivially on K and R . The GF setting is even subtler. The crude lower bound $\dot{a}_0 \approx R\gamma + \delta \geq \sqrt{2}a_0$ does not imply that the majority–minority mode grows faster than the minority mode near the origin; in this sense, it does not provide a GF analogue of the regularization-path fact $\lambda_{K,R} > 1$. Moreover, even a sharper local lower bound would not by itself prove that the majority–minority mode reaches scale η before the minority mode, since the majority mode will have already left the linearized regime by the relevant time.

A.2.3. CONVERGENCE

To further relate GF equations to the RP, it is useful to view (8) as a preconditioned gradient flow over the loss $\mathcal{E}(\mathbf{a}) := \mathcal{L}(\mathbf{U} \text{diag}_{+1}(\mathbf{a}) \tilde{\mathbf{V}}^\top)$ for $\mathbf{a} = (\alpha, a, \mathbf{a}_+, \mathbf{a}_-)$. Indeed, one readily checks that $\dot{\mathbf{a}} = -\mathcal{P}(\mathbf{a})\nabla\mathcal{E}(\mathbf{a})$, where $\mathcal{P}(\mathbf{a})$ is block diagonal: the (α, a) -block is $\mathcal{P}_{\alpha a}(\alpha, a) = \frac{1}{\sqrt{\alpha^2 + a^2}} \begin{pmatrix} a^2 + 2\alpha^2 & \alpha a \\ \alpha a & 2a^2 + \alpha^2 \end{pmatrix}$ and the $\mathbf{a}_+, \mathbf{a}_-$ blocks are $2 \text{diag}(\mathbf{a}_+)$ and $2 \text{diag}(\mathbf{a}_-)$. Since we work in the positive regime, the regularizer $\|\mathbf{a}\|_{+1}$ in Eq. (3) is differentiable and a direct calculation verifies $\nabla\|\mathbf{a}\|_{+1} = 2\mathcal{P}(\mathbf{a})^{-1}\mathbf{a}$. Thus the RP optimality condition $\nabla\mathcal{E}(\mathbf{a}_\lambda) + \lambda\nabla\|\mathbf{a}_\lambda\|_{+1} = 0$ is equivalent to $-\mathcal{P}(\mathbf{a}_\lambda)\nabla\mathcal{E}(\mathbf{a}_\lambda) = 2\lambda\mathbf{a}_\lambda$. Therefore, RP solutions are exactly *radial points* of the GF.

Proposition 8 *A positive point $\mathbf{a} > 0$ is radial for the ODE in Eq. (8), that is satisfies $\dot{\mathbf{a}} = \mu\mathbf{a}$ for some $\mu \in (0, 2)$, iff it is the unique RP solution of (3) with $\lambda = \mu/2$.*

For (8), radially requires $\mathbf{b}_\pm = (\mu/2)\mathbf{1} = \lambda\mathbf{1}$ and since positive radial points coincide with RP solutions, the symmetry argument of Thm. 2 further implies that radial points satisfy $\mathbf{a}_\pm \propto \mathbf{1}$. Although this need not generically hold for GF (unlike for the RP), we can prove that the symmetric subspace is invariant (Fig. 9) and the flow contracts deviations from it measured in KL; App. F.2.

Theorem 9 (Maj. and min. tend to uniformity) *If the GF (8) is initialized so that $\mathbf{a}_\pm \propto \mathbf{1}$ then this remains true. More generally, for positive initialization, $\dot{D}_{KL}(\bar{\mathbf{1}}\|\bar{\mathbf{a}}_\pm) \leq 0$, with equality iff $\mathbf{a}_\pm \propto \mathbf{1}$.*

Proof of Thm. 9 is in App. F.2. The two preceding results show that GF must converge, after normalization, to a logit matrix with three distinct singular values, in agreement with the regularization path. Moreover, among potential limiting configurations, the regularization path itself is the only full-rank case once the logit norm is sufficiently large. The remaining low-rank critical points correspond to networks whose widths are too narrow for the SEL matrix to emerge. Since rank-constrained minima are known to remain critical points in the unconstrained regime [58], and are typically saddles, we expect these low-rank solutions to be unstable. We therefore anticipate that gradient flow converges to the regularization path, which would lead to the emergence dynamics

being inherited more broadly throughout the loss surface. We pose formalizing attractiveness of the regularization path as an important open question. Although it can likely be made rigorous using arguments similar to those of the balanced case [16], we note that the diagonal-plus-rank-one parameterization needed here makes the argument more delicate. As evidence for this added difficulty, App. F.4 shows that even when $R = 1$, the KL divergence to the ETF direction, which is a Lyapunov function under the exact diagonalization of Garrod et al. [16], need not be a Lyapunov function in our coordinates.

Appendix B. Related Work and Limitations

B.1. Detailed Related Work

Our work connects three lines of literature: neural collapse and the UFM, mode emergence in two-layer models, and regularization paths for cross-entropy objectives. While we point out specific connections throughout the main text, we collect and discuss them in more detail here.

Neural collapse, UFM, and data-dependent geometry. Our primary motivation comes from the neural-collapse geometry phenomenon, first empirically identified and formalized by Pappayan et al. [38]. This now rich literature supports the idea that sufficiently large neural networks, trained for sufficiently long, learn training-data representations that converge to specific geometries reflecting the information encoded in the label matrix. This picture is supported both theoretically [23, 25, 47, 49, 56] and empirically [1, 11–14, 20, 31, 38, 52, 61, 62]. A central theoretical route for identifying these geometries is the Unconstrained Features Model (UFM), a two-layer linear model with standard-basis inputs and one output per class. In the balanced one-hot setting, the UFM correctly predicts that learned representations form a simplex equiangular tight frame (ETF) [21, 35, 63]. Its predictions extend to imbalanced classes [8, 12, 13, 52], multilabel data [32], language data [57, 61], and other losses [5, 14, 29].

Deep networks exhibiting neural-collapse geometries are typically trained with cross-entropy (CE), so the faithful UFM reduction is also CE training. This introduces features that are absent from the MSE theory. The softmax map adds a nonlinear coupling across classes on top of the nonconvexity coming from the bilinear factorization $\mathbf{Z} = \mathbf{W}\mathbf{H}$, and unregularized CE drives the optima, and hence the training trajectory, to infinity. This is different from MSE, where the loss structure often enables explicit decompositions and closed-form or nearly closed-form dynamics [21]. In the symmetric balanced one-hot setting, the limiting neural-collapse geometry is the same for MSE-UFM and CE-UFM [21, 35]. Moreover, in that setting, regularization does not change the qualitative conclusion [16, 26, 63]; subtleties start to appear, for example, in deeper UFM variants [15, 48]. Once one departs from the balanced symmetric setting, however, the geometry can become sensitive to whether regularization is present [11–13, 38]. In this paper, we study both the entire ridge-regularization path and the unregularized GF trajectory.

Most of the neural-collapse/UFM literature focuses on convergence: either identifying the global minimizer of a regularized UFM objective [11, 12, 22, 63], or characterizing the limiting direction of GD/GF in the unregularized case [16, 26, 52]. Our focus is different: in addition to convergence, we study *mode emergence* along the training path. Our motivation on this is Zhao and Thrapoulidis [60], who identify the step-imbalanced one-hot setting studied here as the simplest UFM setup with nontrivial mode structure. In this setting, the representation geometry is organized around three semantic axes: a majority mode, a majority–minority mode, and a minority mode. This interpretation is visible directly from the columns of the basis \mathbf{U} : there are $M - 1$ columns supported on the

majority classes, one column separating majorities from minorities with opposite signs, and $M - 1$ columns supported on the minority classes. Thus, the step-imbalanced UFM provides a minimal setting in which one can ask not only where CE training converges, but also which semantic directions are learned first.

Mode emergence and the MSE framework. The mode interpretation above connects the neural-collapse/UFM literature to the long line of work on mode emergence in two-layer linear networks. This line was pioneered by Saxe et al. [44] in terms of technical tools and by Saxe et al. [46] in terms of semantic interpretation. In that framework, the singular modes of the target matrix represent semantic directions, and training under MSE learns these modes sequentially. The resulting picture is mathematically clean and has proved robust across many extensions [4, 6, 18, 36, 53]; see also the recent discussion in Bach [4]. However, this framework is largely restricted to MSE. One of the drivers of the present work is to push this mode-emergence program toward the CE loss.

The difficulty is that the MSE framework relies on two special properties. First, under spectral initialization, the MSE residual $\mathbf{Y} - \mathbf{W}\mathbf{H}$, the logits $\mathbf{W}\mathbf{H}$, and the label matrix \mathbf{Y} are simultaneously diagonalizable in any singular basis of \mathbf{Y} . This fixes the parameter subspaces along the trajectory and reduces training to the evolution of singular values. Second, the resulting singular-value dynamics often admit closed-form or nearly closed-form descriptions. For CE, even the first property fails in general because the softmax nonlinearity couples the logits across classes.

Towards a CE analogy. A recent step toward overcoming this obstacle was taken very recently by Garrod et al. [16]. Their work inspired ours in its technical core. In the balanced one-hot setting, the authors of [16] showed that if one chooses a particular Sylvester–Hadamard singular basis of the label matrix, then the CE residual $\mathbf{Y} - \mathbb{S}(\mathbf{W}\mathbf{H})$ remains diagonalizable together with $\mathbf{W}\mathbf{H}$ and \mathbf{Y} . This preserves the parameter subspaces under spectral initialization and reduces GF to singular-value dynamics. Although these dynamics are not closed form as in MSE, they are still tractable. However, their setting is balanced one-hot data. From the perspective of semantic mode emergence, this case is limited because there is only one nontrivial mode. More importantly, even the smallest departure from balance, namely the step-imbalanced setting studied here, breaks pure Hadamard diagonalization: the imbalanced target matrix is no longer diagonalized by the Sylvester–Hadamard basis in the same way.

Our contribution on this front is to show that the program remains viable beyond the balanced setting. We identify a decomposition of the step-imbalanced label matrix involving Sylvester–Hadamard components, which lets us leverage the same interaction between Hadamard structure, entrywise products, and softmax equivariance that underlies the proof of Garrod et al. [16]. However, exact diagonalization is no longer preserved. Instead, Theorem 1 identifies the correct replacement: a diagonal-plus-rank-one decomposition of the CE residual. We then show that this structure is still sufficient to reduce the full matrix ODE to a vector ODE. Thus, while pure diagonalization breaks under imbalance, the analysis can continue after enlarging the spectral structure in the right way.

Regularization paths and sparse mode activation. A further distinction of our work is that the same spectral structure also solves the regularization path. This is not an initialization assumption: for every regularization strength, we prove that the global minimizer of the logit-space nuclear-norm problem lies in the diagonal-plus-rank-one subspace identified by the softmax analysis. A priori, this is not obvious. The regularization-path problem is a convex logit-space problem, and it is not clear in advance why its solution should admit a tight reduction to an ℓ_1 -style CE minimization, let alone to a four-dimensional problem with analytic phase formulas.

There is related work on regularization paths for logistic regression and other exponentially tailed losses [41, 42]. In particular, Park and Hastie [39] develop an algorithm for numerically producing the entire trajectory of solutions to ℓ_1 -regularized logistic regression for general input data. There, the ℓ_1 -style penalty induces sparsity and sequential entry of coordinates. In our specific data setting, the coordinates are singular modes of the imbalanced label geometry, and we can analyze the path explicitly. Thus, the sparsity induced by the nuclear-norm/logit formulation becomes a mode-emergence phenomenon: majority first, then majority–minority, and finally minority.

Comparison with finite-regularization CE-UFM results. The closest finite-regularization CE-UFM results we are aware of are [12, 22]. Dang et al. [12] characterize global minimizers of a CE objective under an unconstrained ReLU-feature model, i.e., a UFM variant with entrywise nonnegative features. Their Theorem 4.1 gives a class-wise thresholding result: when the effective regularization is large relative to the size of a class, the corresponding class mean collapses to zero. Thus their path-like structure is a class-mean collapse/activation phenomenon in a nonnegative-feature model. Our result concerns a different object: we study the standard unconstrained UFM through the logit-space nuclear-norm formulation and characterize the regularization path spectrally, mode by mode. Modes are not classes. In our path, the majority singular mode activates first, then a majority–minority mode, and finally the minority mode. By contrast, the thresholding in [12] is class-wise: class means become nonzero sequentially, and when active they satisfy the angular structure specified by their theorem. The model and mechanism are also different: entrywise nonnegativity changes the feasible geometry and enables a lower-bound/equality-condition proof strategy, whereas our analysis works in the standard unconstrained UFM through a logit-space spectral reduction. The result closest to our RP analysis is [22]. They study the convexified CE-UFM logit/nuclear-norm problem (\mathbf{Z} -UFM $_\lambda$) and characterize the block-support regimes of the global minimizer under class imbalance. In the bias-free two-cluster specialization relevant to our setting, their Theorem 3.2 gives the same qualitative sequence of supports, after scaling the regularization parameter by the number of samples: zero, majority-only, majority plus cross-cluster block, and fully active. In this sense, we can map their intermediate threshold, after a change of coordinates, to our $\lambda_{K,R}$. However, their result is a block-support characterization: it determines which blocks of the mean prediction matrix are active in each regime. First, this does *not* provide a mode-emergence interpretation. Second, they do *not* analyze how this threshold behaves as function of K, R as we do in App. E.4. Third, and more important, it does *not* solve the RP. This distinction is essential for our purposes. Our Hadamard-coordinate reduction turns the path into a four-dimensional ℓ_1 -style problem whose coordinates are the majority, majority–minority, and minority singular modes. We then solve the active optimality equations phase by phase, obtaining scalar formulas for the mode amplitudes throughout the path. Our analytic mode-wise solution characterizes the evolution of the singular values, the $\lambda \rightarrow 0$ expansion of the normalized singular profile, the logarithmic directional convergence rate, and is also the structure that lets us connect the regularization path to gradient-flow mode emergence. Finally, [52] introduced SELI as the limiting imbalanced neural-collapse geometry and showed that the UFM converges to this geometry in the vanishing-regularization limit. This was the first work to explicitly characterize the imbalanced limiting geometry through explicit analysis of the UFM and connect it to minority collapse of Fang et al. [13] as a special case. Their analysis identifies the endpoint of the path. We study the path itself: the finite-regularization phases, mode-wise singular-value evolution, convergence rate to the SELI direction, and gradient-flow emergence.

B.2. Limitations

We analyze how CE training learns data modes in a deliberately minimal setting: Step-imbalanced and orthogonal data exhibit non-trivial mode structure, and the two-layer linear architecture allows this structure to be learned sequentially. Given the lack of CE trajectory-level results we deem this a reasonable first step. In this setting, we solved the ridge-regularization path and studied GF under a spectral initialization that preserves the subspaces identified by this path. The resulting dynamics reveal both quantitative and qualitative differences from the MSE dynamics most commonly studied in the literature. Methodologically, the main obstacle in extending the MSE theory to CE is the loss of closed-form decoupled dynamics. By extending the softmax diagonalization perspective of [16] to a non-symmetric setting, we further showcase that this obstacle is not necessarily a dead end. In the imbalanced UFM, closed-form expressions can be replaced by low-dimensional scalar equations for the regularization path and low-dimensional ODEs for gradient flow, both of which remain amenable to analysis. Our regularization-path analysis assumes that K is a power of two, although empirically the formulas appear more general (Fig. 6). Our GF analysis uses a spectral initialization that keeps the trajectory in the subspace identified by the regularization path. This is an exact CE trajectory, and empirically resembles the dynamics from small random initialization (Fig. 8), but proving this connection remains open. Future directions include extending the analysis to nonlinear and deeper models, richer mode structures, non-orthogonal inputs, and other softmax-based settings such as attention.

Appendix C. Notation and basic facts

For simplicity, we introduce additional necessary notation here. All matrix indexing starts from 0 instead of 1, and Φ_2 is commonly denoted as Φ , with columns $[\phi_0, \phi_1]$. In addition, Ψ_k is defined as $\Psi_K := \mathbf{1}_{K-1} \mathbf{1}_{K-1}^T - \Phi_{[1:,1:]}$. Note that $\Psi_K \in \{0, 2\}^{(K-1) \times (K-1)}$ is invertible and satisfies $\Psi \mathbf{1}_{K-1} = K \mathbf{1}_{K-1}$, $\Psi^{-1} \mathbf{1}_{K-1} = \frac{1}{K} \mathbf{1}_{K-1}$.

The following properties of Hadamard matrices shown in the following lemmas will be useful later in the proofs.

Lemma 10 *For integers m, i, j , consider their $l = \log_2(M)$ -bit binary representations $m_1, \dots, m_l, i_1, \dots, i_l$ and j_1, \dots, j_l , respectively. Accordingly define the permutation matrices*

$$\mathbf{\Pi}_m := \mathbf{J}^{m_1} \otimes \dots \otimes \mathbf{J}^{m_l}, \quad \mathbf{J} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

Then, any Sylvester Hadamard matrix of order M satisfies

$$\mathbf{\Pi}_m \bar{\Phi}_M \mathbf{e}_i = (-1)^{m \cdot (i+1)} \bar{\Phi}_M \mathbf{e}_i, \quad (9)$$

$$\mathbf{\Pi}_{i \oplus j} (\Phi_M \mathbf{e}_m \mathbf{e}_m^T \Phi_M^T \mathbf{e}_j) = \Phi_M \mathbf{e}_m \mathbf{e}_m^T \Phi_M^T \mathbf{e}_i, \quad (10)$$

where $i \oplus j$ is the xor operation.

In addition, the permutation matrices $\mathbf{\Pi}_m$ are symmetric and satisfy

$$\sum_{m=0}^{M-1} \mathbf{\Pi}_m = (\mathbf{J} + \mathbf{I}_2) \otimes \dots \otimes (\mathbf{J} + \mathbf{I}_2) = \mathbf{1}_M \mathbf{1}_M^T. \quad (11)$$

Proof The key to prove (9) is to expand the RHS as a Kronecker product by using the fact that $\bar{\Phi}_M e_i = \Phi_M e_{i+1}$:

$$\Pi_m \bar{\Phi}_M e_i = \bigotimes_{p=0}^l \mathbf{J}^{m_p} \phi_{(i+1)_p} = \bigotimes_{p=0}^l (-1)^{m_p(i+1)_p} \phi_{(i+1)_p} = (-1)^{m(i+1)} \bar{\Phi}_M e_i.$$

The second identity can be proved from the previous one as follows:

$$\begin{aligned} \Pi_{i \oplus j} (\Phi_M e_m e_m^T \Phi_M^T) e_j &= (-1)^{m \cdot j} \Pi_{i \oplus j} (\Phi_M e_m) = (-1)^{m \cdot j} (-1)^{(i+j) \cdot m} (\Phi_M e_m) \\ &= (-1)^{m \cdot i} (\Phi_M e_m) = (\Phi_M e_m e_m^T \Phi_M^T) e_i. \end{aligned}$$

■

Lemma 11 *The matrix $\mathbf{Z} \in \mathbb{R}^{K \times n}$ which is defined below, can be diagonalized as follows:*

$$\begin{aligned} \mathbf{Y} &= \mathbf{I}_M \otimes \begin{bmatrix} \mathbf{1}_R^T & 0 \\ \mathbf{0}_R^T & 1 \end{bmatrix}, \quad \hat{\mathbf{Z}} = (\mathbf{I}_k - \mathbf{1}_k \mathbf{1}_k^T) \mathbf{Y} \\ \mathbf{U} &:= \begin{bmatrix} \frac{1}{\sqrt{2M}} \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} & \frac{1}{\sqrt{M}} \bar{\Phi}_M \otimes \mathbf{I}_2 \end{bmatrix} \Sigma_{\mathbf{Z}} := \begin{bmatrix} \sqrt{\frac{R+1}{2}} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{M-1} \otimes \Sigma_2 \end{bmatrix} \\ \Sigma_2 &:= \begin{bmatrix} \sqrt{R} & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{V} := \begin{bmatrix} \frac{1}{\sqrt{M(R+1)}} \mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ -1 \end{pmatrix} & \frac{1}{\sqrt{M}} \bar{\Phi}_M \otimes \begin{pmatrix} \frac{\mathbf{1}_R}{\sqrt{R}} & \frac{\mathbf{0}_R}{\sqrt{R}} \\ 0 & 1 \end{pmatrix} \end{bmatrix}. \end{aligned} \quad (12)$$

Proof The diagonalization in (12) can be proved by direct matrix multiplication and the proof is omitted for brevity. Alternatively, $\mathbf{Z}\mathbf{Z}^T$ can be diagonalized by noting that its eigenvectors are of the form of Kronecker products, which follows from standard linear algebra. ■

Lemma 12 *For any two coordinates $r, s \in [M-1]$, there exist permutation matrices*

$$\mathbf{P} \in \mathbb{R}^{M \times M}, \quad \mathbf{C} \in \mathbb{R}^{(M-1) \times (M-1)}$$

such that

$$\mathbf{P} \bar{\Phi}_M = \bar{\Phi}_M \mathbf{C}, \quad \mathbf{C} e_r = e_s.$$

Proof Let $\ell = \log_2 M$. We use the standard representation of the Sylvester–Hadamard matrix $(\Phi_M)_{\mathbf{x}, \mathbf{y}} = (-1)^{\mathbf{x}^\top \mathbf{y}}$, where $\mathbf{x}, \mathbf{y} \in \{0, 1\}^\ell$ are binary encodings: \mathbf{x} indexes rows and \mathbf{y} indexes columns. The column $\mathbf{y} = \mathbf{0}$ is the all-ones column, so the columns of $\bar{\Phi}_M$ are indexed by the nonzero vectors $\mathbf{y} \neq \mathbf{0}$.

Let \mathbf{y}_r and \mathbf{y}_s be the nonzero binary vectors corresponding to coordinates r and s . We choose an invertible binary matrix \mathbf{A} such that

$$\mathbf{A}^\top \mathbf{y}_r = \mathbf{y}_s.$$

Such an \mathbf{A} exists because any nonzero binary vector can be mapped to any other nonzero binary vector by an invertible linear change of coordinates.

Since \mathbf{A} is invertible, the map $\mathbf{x} \mapsto \mathbf{A}\mathbf{x}$ is a bijection of the row index set $\{0, 1\}^\ell$. Therefore it corresponds to a row permutation matrix \mathbf{P} .

Now compute the entries of the row-permuted Hadamard matrix:

$$(\mathbf{P}\overline{\Phi}_M)_{\mathbf{x},\mathbf{y}} = (\overline{\Phi}_M)_{\mathbf{A}\mathbf{x},\mathbf{y}} = (-1)^{(\mathbf{A}\mathbf{x})^\top \mathbf{y}} = (-1)^{\mathbf{x}^\top \mathbf{A}^\top \mathbf{y}}.$$

But the last expression is exactly the Hadamard column indexed by $\mathbf{A}^\top \mathbf{y}$. Therefore, after the row permutation \mathbf{P} , the column indexed by \mathbf{y} becomes the column indexed by $\mathbf{A}^\top \mathbf{y}$.

Since \mathbf{A}^\top is invertible, it permutes all nonzero binary vectors and fixes $\mathbf{0}$. Hence the all-ones column stays fixed, and the remaining columns are merely permuted. Therefore there exists a permutation matrix \mathbf{C} such that

$$\mathbf{P}\overline{\Phi}_M = \overline{\Phi}_M \mathbf{C}.$$

Moreover, because we chose \mathbf{A} so that $\mathbf{A}^\top \mathbf{y}_r = \mathbf{y}_s$, the induced column permutation sends coordinate r to coordinate s , i.e. $\mathbf{C}\mathbf{e}_r = \mathbf{e}_s$. \blacksquare

Appendix D. Theorem 1 Proof and Discussions

In this section, we prove that the diagonal-plus-rank-one factorization is in fact a special case of the more general block-diagonal-plus-rank-one factorization which is made possible by Sylvester Hadamard matrices. In addition, to highlight the appearance of the rank-one term, we first prove the general permutation structure *without* having this term in the softmax argument, and later show that the results are retained, with minor changes to the expressions.

D.1. The Mapping of Theorem 1 and Proof Outline

The map $(\beta, b, \mathbf{b}_+, \mathbf{b}_-) := \mathcal{F}((\alpha, a, \mathbf{a}_+, \mathbf{a}_-))$ of Theorem 1 is given explicitly by

$$\beta = \sqrt{\frac{R}{2(R+1)}} (X_+ - X_-), \quad b = \sqrt{\frac{1}{2(R+1)}} (RX_+ + X_-), \quad (13a)$$

$$\mathbf{b}_+ = \sqrt{R}(\mathbf{1}_{M-1} - \frac{1}{2M} X_+ \overline{\Phi}_M^\top \mathbf{e}^{\mathbf{x}_+}), \quad \mathbf{b}_- = \mathbf{1}_{M-1} - \frac{1}{2M} X_- \overline{\Phi}_M^\top \mathbf{e}^{\mathbf{x}_-}, \quad (13b)$$

with $X_\pm := \frac{2M}{\mathbf{1}_M^\top \mathbf{e}^{\mathbf{x}_\pm + M}}$, $\mathbf{x}_+ := \frac{1}{M}(\gamma \mathbf{1}_M + \frac{1}{\sqrt{R}} \overline{\Phi}_M \mathbf{a}_+)$, $\mathbf{x}_- := \frac{1}{M}(\delta \mathbf{1}_M + \overline{\Phi}_M \mathbf{a}_-)$, and,

$$(\gamma, \delta) := \sqrt{\frac{2}{R+1}} \left(a + \frac{1}{\sqrt{R}} \alpha, a - \alpha \sqrt{R} \right). \quad (14)$$

As mentioned in Sec. 3 the main difficulty in the proof is identifying the right decomposition and obtaining the formulas above. The key in the proof in App. D was to realize that the residual matrix \mathbf{Q} in the LHS of (2) can decompose its $n = (R+1)M$ columns in M blocks \mathbf{Q}_i satisfying $\mathbf{Q}_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [\mathbf{p}_+ \mathbf{1}_R^\top \quad \mathbf{p}_-] \in \mathbb{R}^{K \times (R+1)}$, for zero-sum $\mathbf{p}_+ := \mathbf{e}_0 - \mathbb{S}(\mathbf{x}_+ \otimes \mathbf{f}_+)$ and $\mathbf{p}_- := \mathbf{e}_1 - \mathbb{S}(\mathbf{x}_- \otimes \mathbf{f}_-)$. Finally, it can be shown that any matrix with such structure can be decomposed into a block-diagonal-plus-rank-one structure, with the off-diagonal entries vanishing thanks to the specific structure of $\mathbf{p}_+, \mathbf{p}_-$.

Note that exact diagonalization is possible by using the singular vector $(a\mathbf{v}_0 + \alpha\tilde{\mathbf{v}})/\sqrt{a^2 + \alpha^2}$ with corresponding singular value $\sqrt{a^2 + \alpha^2}$. But what makes the theorem's formulation suitable for subsequent analysis is that it keeps the right subspace $\tilde{\mathbf{V}}$ fixed. Note also that Eq. (14) defines a one-to-one mapping between (α, a) and (γ, δ) pair; we often find it convenient to work with the latter.

D.2. Proof of Theorem 1

Lemma 13 For any $\mathbf{A} \in \mathbb{R}^{(K-1) \times (K-1)}$ having the block structure

$$\mathbf{A} = \text{diag}(a, \mathbf{A}_0, \dots, \mathbf{A}_{M-2}) = \text{diag}\left(a, \sum_{m=0}^{M-2} (\mathbf{e}_m \mathbf{e}_m^T) \otimes \mathbf{A}_m\right), \quad \mathbf{A}_m \in \mathbb{R}^{2 \times 2},$$

the i -th block column of $\mathbf{Y} - \mathbb{S}(\mathbf{UAV}^T)$, denoted as $(\mathbf{Y} - \mathbb{S}(\mathbf{UAV}^T))_i \in \mathbb{R}^{K \times (R+1)}$ satisfies

$$(\mathbf{Y} - \mathbb{S}(\mathbf{UAV}^T))_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [\mathbf{p} \mathbf{1}_R^T \quad \mathbf{q}], \quad \mathbf{p} = \mathbf{e}_0 - \mathbb{S}(\mathbf{x}), \quad \mathbf{q} = \mathbf{e}_1 - \mathbb{S}(\mathbf{y}),$$

where \mathbf{x}, \mathbf{y} are defined as

$$\begin{aligned} \mathbf{x} &:= \frac{1}{2M} \left(\sqrt{\frac{2}{R+1}} a \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} \bar{\Phi}_M \mathbf{e}_m \otimes \frac{\mathbf{a}_{m,0}}{\sqrt{R}} \right) \\ \mathbf{y} &:= \frac{1}{2M} \left(-\sqrt{\frac{2}{R+1}} a \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} \bar{\Phi}_M \mathbf{e}_m \otimes \mathbf{a}_{m,1} \right). \end{aligned} \quad (15)$$

Proof The first step is to calculate \mathbf{UAV}^T , to reach

$$\begin{aligned} \mathbf{UAV}^T &= \frac{a}{\sqrt{R+1}} \left(\frac{1}{M} \mathbf{1}_M \mathbf{1}_M^T \right) \otimes \left[\frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \mathbf{1}_R^T \quad -\frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right] \\ &\quad + \frac{1}{2} \sum_{m=0}^{M-2} (\bar{\Phi}_M \mathbf{e}_m \mathbf{e}_m^T \bar{\Phi}_M^T) \otimes \left(\mathbf{A}_m \begin{bmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} \frac{1}{\sqrt{R}} & \begin{pmatrix} 0 \\ 1 \end{pmatrix} \end{bmatrix} \right). \end{aligned} \quad (16)$$

Next, we analyze the structure of $\bar{\Phi}_M \mathbf{e}_m \mathbf{e}_m^T \bar{\Phi}_M^T$. Using $\bar{\Phi}_M \mathbf{e}_m = \Phi \mathbf{e}_{m+1}$ and defining $\mathbf{a}_{m,0} := \mathbf{A}_m \mathbf{e}_0$, $\mathbf{a}_{m,1} := \mathbf{A}_m \mathbf{e}_1$, the summation can be changed to

$$\sum_{m=0}^{M-2} (\Phi_M \mathbf{e}_{m+1} \mathbf{e}_{m+1}^T \Phi_M^T) \otimes \begin{bmatrix} \mathbf{a}_{m,0} \frac{1}{\sqrt{R}} & \mathbf{a}_{m,1} \end{bmatrix}.$$

The result of Lemma (10) shows that all columns of $\Phi_M \mathbf{e}_m \mathbf{e}_m^T \Phi_M^T$ are permutations of each other. Consequently, we fix the first column as reference, to expand the first term in the Kronecker product as

$$\Phi_M \mathbf{e}_{m+1} \mathbf{e}_{m+1}^T \Phi_M^T = [\Phi_M \mathbf{e}_{m+1} \quad \dots \quad \mathbf{\Pi}_{M-1}(\Phi_M \mathbf{e}_{m+1})] \quad (17)$$

where we have used the fact that $\mathbf{\Pi}_{i \oplus 0} = \mathbf{\Pi}_i$. Calculating the above summation and combining it with the term involving a would yield

$$\mathbf{UAV}^T = \begin{bmatrix} \mathbf{x} \mathbf{1}_R^T & \mathbf{y} & (\tilde{\mathbf{\Pi}}_1 \mathbf{x}) \mathbf{1}_R^T & \tilde{\mathbf{\Pi}}_1 \mathbf{y} & \dots & (\tilde{\mathbf{\Pi}}_{M-1} \mathbf{x}) \mathbf{1}_R^T & \tilde{\mathbf{\Pi}}_{M-1} \mathbf{y} \end{bmatrix}, \quad (18)$$

where $\tilde{\mathbf{\Pi}}_m := \mathbf{\Pi}_m \otimes \mathbf{I}_2$ and \mathbf{x}, \mathbf{y} are defined as above.

Since for any permutation matrix $\mathbf{\Pi}$, the softmax operator satisfies $\mathbb{S}(\mathbf{\Pi x}) = \mathbf{\Pi} \mathbb{S}(\mathbf{x})$, the expression $\mathbb{S}(\mathbf{UAV}^T)$ will also have the same block column structure. In addition, using $(\mathbf{\Pi}_i \otimes \mathbf{I}_2) \mathbf{e}_0 = \mathbf{e}_{2i}$ and $(\mathbf{\Pi}_i \otimes \mathbf{I}_2) \mathbf{e}_1 = \mathbf{e}_{2i+1}$, it can be shown that $(\mathbf{Y} - \mathbb{S}(\mathbf{UAV}^T))_i$ satisfies

$$(\mathbf{Y} - \mathbb{S}(\mathbf{UAV}^T))_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [(\mathbf{e}_0 - \mathbb{S}(\mathbf{x})) \mathbf{1}_R^T \quad \mathbf{e}_1 - \mathbb{S}(\mathbf{y})]$$

■

Corollary 14 *If the matrix \mathbf{A} is diagonal, with each block $\mathbf{A}_m = \text{diag}([a_{m,0}, a_{m,1}])$, the softmax expressions simplify to*

$$\begin{aligned}\tilde{\mathbf{x}} &:= \frac{1}{M} \left(\sqrt{\frac{2}{R+1}} a \mathbf{1}_M + \frac{1}{\sqrt{R}} \bar{\Phi}_M \mathbf{a}_+ \right), \quad \mathbb{S}(\mathbf{x}) = \mathbb{S} \left(\tilde{\mathbf{x}} \otimes \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right) \\ \tilde{\mathbf{y}} &:= \frac{1}{M} \left(\sqrt{\frac{2}{R+1}} a \mathbf{1}_M + \bar{\Phi}_M \mathbf{a}_- \right) \quad \mathbb{S}(\mathbf{y}) = \mathbb{S} \left(\tilde{\mathbf{y}} \otimes \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right),\end{aligned}\tag{19}$$

where $\mathbf{a}_+, \mathbf{a}_- \in \mathbb{R}^{M-2}$ and $(\mathbf{a}_+)_m = a_{m,0}$, $(\mathbf{a}_-)_m = a_{m,1}$.

Proof Direct substitution of diagonal \mathbf{A}_m in (15) leads to

$$\begin{aligned}\mathbf{x} &= \frac{1}{2M} \left(\sqrt{\frac{2}{R+1}} a \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} \frac{a_{m,0}}{\sqrt{R}} \bar{\Phi}_M \mathbf{e}_m \otimes \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right) \\ \mathbf{y} &= \frac{1}{2M} \left(-\sqrt{\frac{2}{R+1}} a \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} a_{m,1} \bar{\Phi}_M \mathbf{e}_m \otimes \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right).\end{aligned}\tag{20}$$

At this point, notice that adding or subtracting multiples of $\mathbf{1}_K$ from \mathbf{x}, \mathbf{y} does not change the softmax value. As a result, by adding appropriate constant vectors to \mathbf{x}, \mathbf{y} to change their first term, we reach the desired result. \blacksquare

The following theorem establishes that for any matrix \mathbf{Q} with the specified permutation structure, the projection $\mathbf{U}^T \mathbf{Q} \mathbf{V}$ is block diagonal with the same pattern as the argument \mathbf{A} .

Theorem 15 *Suppose a matrix \mathbf{Q} has the block structure $\mathbf{Q}_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [\mathbf{p} \mathbf{1}_R^T \quad \mathbf{q}]$. Then, for \mathbf{U}, \mathbf{V} as defined in (12)*

$$\mathbf{U}^T \mathbf{Q} \mathbf{V} = \mathbf{B} = \text{diag}(b, \mathbf{B}_0, \dots, \mathbf{B}_{M-2}), \quad \mathbf{B}_i \in \mathbb{R}^{2 \times 2},\tag{21}$$

$$b = \frac{1}{\sqrt{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (R\mathbf{p} - \mathbf{q}), \quad \mathbf{B}_i = ((\bar{\Phi}_M \mathbf{e}_i)^T \otimes \mathbf{I}_2) [\sqrt{R}\mathbf{p} \quad \mathbf{q}]\tag{22}$$

Proof The matrix \mathbf{Q} can be expressed as $\mathbf{Q} = \sum_{m=0}^{M-1} \mathbf{e}_m^T \otimes ((\mathbf{\Pi}_m \otimes \mathbf{I}_2) [\mathbf{p} \mathbf{1}_R^T \quad \mathbf{q}])$.

Using this expansion and (11), the first column of \mathbf{B} can be established as

$$\begin{aligned}\mathbf{Q} \mathbf{v}_0 &= \frac{1}{\sqrt{M(R+1)}} \sum_{m=0}^{M-1} (\mathbf{\Pi}_m \otimes \mathbf{I}_2) (R\mathbf{p} - \mathbf{q}) = \frac{1}{\sqrt{M(R+1)}} ((\mathbf{1}_M \mathbf{1}_M^T) \otimes \mathbf{I}_2) (R\mathbf{p} - \mathbf{q}) \\ \mathbf{B} \mathbf{e}_0 &= \mathbf{U}^T \mathbf{Q} \mathbf{v}_0 = \frac{1}{\sqrt{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (R\mathbf{p} - \mathbf{q}) \mathbf{e}_0 = b \mathbf{e}_0.\end{aligned}\tag{23}$$

Next, we will show that the first row is also the same multiple b of \mathbf{e}_0^T :

$$\begin{aligned}\mathbf{u}_0^T \mathbf{Q} &= \sum_{m=0}^{M-1} \mathbf{e}_m^T \otimes \underbrace{(\mathbf{u}_0^T (\mathbf{\Pi}_m \otimes \mathbf{I}_2) [\mathbf{p} \mathbf{1}_R^T \quad \mathbf{q}])}_{\mathbf{u}_0^T} = \mathbf{1}_M^T \otimes [(\mathbf{u}_0^T \mathbf{p}) \mathbf{1}_R^T \quad \mathbf{u}_0^T \mathbf{q}] \\ \mathbf{u}_0^T \mathbf{Q} \mathbf{V} &= \frac{\sqrt{2M} \mathbf{u}_0^T (R\mathbf{p} - \mathbf{q})}{\sqrt{2(R+1)}} \mathbf{e}_0^T = \frac{(\bar{\Phi}_K \mathbf{e}_0)^T (R\mathbf{p} - \mathbf{q})}{\sqrt{2(R+1)}} \mathbf{e}_0^T = b \mathbf{e}_0^T.\end{aligned}\tag{24}$$

The final phase would be to compute each block of the sub-matrix $\mathbf{B}_{[1:,1:]}$. Recall that the matrix \mathbf{Q} can be expressed as $\mathbf{Q} = \sum_{m=0}^{M-1} \mathbf{e}_m^T \otimes ((\mathbf{\Pi}_m \otimes \mathbf{I}_2) [\mathbf{p}\mathbf{1}_R^T \ \mathbf{q}])$. Expanding the expression for each block would result in

$$\begin{aligned} (\mathbf{B}_{[1:,1:]})_{ij} &= (\mathbf{U}_{[1:,1:]}^T \mathbf{Q} \mathbf{V}_{[1:,1:]})_{ij} = (\overline{\mathbf{\Phi}}_M \mathbf{e}_i \otimes \mathbf{I}_2)^T \mathbf{Q} \left((\overline{\mathbf{\Phi}}_M \mathbf{e}_j) \otimes \frac{1}{\sqrt{2}} \begin{bmatrix} \frac{1_R}{\sqrt{R}} & \frac{0_R}{\sqrt{R}} \\ 0 & 1 \end{bmatrix} \right) \\ &= \frac{1}{M} (\overline{\mathbf{\Phi}}_M \mathbf{e}_i \otimes \mathbf{\Phi})^T \sum_{m=0}^{M-1} (\mathbf{e}_m^T \overline{\mathbf{\Phi}}_M \mathbf{e}_j) (\mathbf{\Pi}_m \otimes \mathbf{I}_2) [\sqrt{R} \mathbf{p} \ \mathbf{q}] \\ &= \frac{1}{M} \sum_{m=0}^{M-1} (-1)^{m \cdot (j+1)} ((\mathbf{\Pi}_m \overline{\mathbf{\Phi}}_M \mathbf{e}_i)^T \otimes \mathbf{\Phi}) [\sqrt{R} \mathbf{p} \ \mathbf{q}], \end{aligned}$$

where we have used the symmetry of $\mathbf{\Pi}_m$, $\mathbf{\Phi}$, and the mixed-product Kronecker identity. Next, we use the identity (9) for $\mathbf{\Pi}_m$ to simplify the expression to

$$\begin{aligned} (\mathbf{B}_{[1:,1:]})_{ij} &= \frac{1}{M} \sum_{m=0}^{M-1} (-1)^{m \cdot (j+i)} ((\overline{\mathbf{\Phi}}_M \mathbf{e}_i)^T \otimes \mathbf{I}_2) [\sqrt{R} \mathbf{p} \ \mathbf{q}] \\ &= \delta_{ij} ((\overline{\mathbf{\Phi}}_M \mathbf{e}_i)^T \otimes \mathbf{\Phi}) [\sqrt{R} \mathbf{p} \ \mathbf{q}], \end{aligned}$$

and the desired result is proved. \blacksquare

The following lemma proves the rank-one difference for matrices \mathbf{Q} with zero-sum columns and the specific permutation structure of this section.

Lemma 16 *The matrices \mathbf{U}, \mathbf{V} as defined in (12), satisfy*

$$\begin{aligned} \mathbf{U} \mathbf{U}^T &= \mathbf{I}_k - \frac{1}{K} \mathbf{1}_K \mathbf{1}_K^T \\ \mathbf{V} \mathbf{V}^T - \mathbf{I}_n &= \mathbf{I}_M \otimes \begin{bmatrix} \frac{1_R \mathbf{1}_R^T}{R} - \mathbf{I}_R & \mathbf{0}_R \\ \mathbf{0}_R^T & 0 \end{bmatrix} - \frac{1}{M(R+1)} (\mathbf{1}_M \mathbf{1}_M^T) \otimes \begin{bmatrix} \frac{1_R \mathbf{1}_R^T}{R} & \mathbf{1}_R \\ \mathbf{1}_R^T & R \end{bmatrix}. \end{aligned} \quad (25)$$

Furthermore, it can be shown that for zero-sum vectors \mathbf{p} and \mathbf{q} , and the matrix $\mathbf{Q} \in \mathbb{R}^{K \times n}$ defined such that its i -th column block $\mathbf{Q}_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [\mathbf{p}\mathbf{1}_R^T \ \mathbf{q}]$, the following hold:

$$\mathbf{U} \mathbf{U}^T \mathbf{Q} = \mathbf{Q}, \quad \mathbf{Q} (\mathbf{V} \mathbf{V}^T - \mathbf{I}_k) = -\beta \mathbf{u}_0 \tilde{\mathbf{v}}^T, \quad (26)$$

$$\beta = \sqrt{\frac{R}{2(R+1)}} (\overline{\mathbf{\Phi}}_K \mathbf{e}_0)^T (\mathbf{p} + \mathbf{q}), \quad \mathbf{u}_0 = \mathbf{U} \mathbf{e}_0 = \overline{\mathbf{\Phi}}_K \mathbf{e}_0, \quad \tilde{\mathbf{v}} = \frac{\mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ R \end{pmatrix}}{\sqrt{MR(R+1)}} \quad (27)$$

Proof Equation (25) can be proved by direct calculation. The identity $\mathbf{U} \mathbf{U}^T \mathbf{Q} = \mathbf{Q}$ can be proved that every column of \mathbf{Q} consists of permutations of zero-sum vectors \mathbf{p}, \mathbf{q} which do not change under the projection $\mathbf{U} \mathbf{U}^T$.

Hence, we focus on proving the rank-1 difference between \mathbf{Q} and $\mathbf{Q}\mathbf{V}\mathbf{V}^T$:

$$\begin{aligned} \mathbf{Q}(\mathbf{V}\mathbf{V}^T - \mathbf{I}_n) &= \sum_{m=0}^{M-1} \left(\mathbf{e}_m^T \otimes \mathbf{0}_{K \times (R+1)} - \frac{1}{M(R+1)} \mathbf{1}_M^T \otimes ((\mathbf{\Pi}_m \otimes \mathbf{I}_2)(\mathbf{p} + \mathbf{q}) [\mathbf{1}_R^T \ R]) \right) \\ &= -\frac{1}{M(R+1)} \mathbf{1}_M^T \otimes (((\mathbf{1}_M \mathbf{1}_M^T) \otimes \mathbf{I}_2)(\mathbf{p} + \mathbf{q}) [\mathbf{1}_R^T \ R]) \\ &= -\frac{1}{M(R+1)} (\mathbf{1}_M \otimes ((\mathbf{1}_M^T \otimes \mathbf{I}_2)(\mathbf{p} + \mathbf{q})) (\mathbf{1}_M^T \otimes [\mathbf{1}_R^T \ R]). \end{aligned} \quad (28)$$

In addition, the fact that \mathbf{p}, \mathbf{q} are zero-sum can be used to reach the following simplification:

$$\begin{aligned} (\mathbf{1}_M^T \otimes \mathbf{I}_2)(\mathbf{p} + \mathbf{q}) &= (\mathbf{1}_M^T \otimes (\mathbf{I}_2 - \frac{1}{2} \mathbf{1}_2 \mathbf{1}_2^T))(\mathbf{p} + \mathbf{q}) = \frac{1}{2} \left(\mathbf{1}_M^T \otimes \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \right) (\mathbf{p} + \mathbf{q}) \\ &= \frac{1}{2} \begin{pmatrix} 1 \\ -1 \end{pmatrix} (\mathbf{1}_M^T \otimes (1 \ -1)) (\mathbf{p} + \mathbf{q}) = \frac{1}{2} \begin{pmatrix} 1 \\ -1 \end{pmatrix} (\overline{\Phi}_K \mathbf{e}_0)^T (\mathbf{p} + \mathbf{q}). \end{aligned} \quad (29)$$

Normalizing each row and column vector using the above simplification in (29) would lead to the desired result. \blacksquare

Remark 17 *Even in the case of $R = 1$, the rank-1 difference does not vanish if a block diagonal \mathbf{A} is used, and it only vanishes for a diagonal matrix inside the softmax structure.*

The following lemma will show how Theorem 15 and Theorem 16 can be used to give a factorization of the matrix \mathbf{Q} :

Lemma 18 *Suppose a matrix \mathbf{M} , for given $\mathbf{U}, \mathbf{V}, \tilde{\mathbf{v}}$ and $\mathbf{B} := \mathbf{U}^T \mathbf{Q} \mathbf{V}$ satisfies*

$$\mathbf{U}\mathbf{U}^T \mathbf{Q} = \mathbf{Q}, \quad \mathbf{Q}(\mathbf{V}\mathbf{V}^T - \mathbf{I}_k) = -\beta(\mathbf{U}\mathbf{e}_0)\tilde{\mathbf{v}}^T.$$

then, the matrix \mathbf{Q} can be decomposed as $\mathbf{Q} = \mathbf{U} [\beta \mathbf{e}_0 \ \mathbf{B}] [\tilde{\mathbf{v}} \ \mathbf{V}]^T$.

Proof This lemma can be proved by direct calculation based on the assumptions:

$$\mathbf{U}\mathbf{B}\mathbf{V}^T = \mathbf{U}\mathbf{U}^T \mathbf{Q} \mathbf{V}\mathbf{V}^T = \mathbf{Q} - \beta \mathbf{u}_0 \tilde{\mathbf{v}}^T = \mathbf{Q} - \mathbf{U}(\beta \mathbf{e}_0) \tilde{\mathbf{v}}^T,$$

after which rearranging and rewriting in block form will prove the result. \blacksquare

Since we want the softmax spectral argument and output to share their structure, we add a term involving α to the block matrix,

Theorem 19 *Suppose the matrix $\mathbf{A} \in \mathbb{R}^{(K-1) \times (K-1)}$ has the block structure*

$$\mathbf{A} = \text{diag}(a, \mathbf{A}_0, \dots, \mathbf{A}_{M-2}), \quad \mathbf{A}_m \in \mathbb{R}^{2 \times 2}$$

Then, for \mathbf{U}, \mathbf{V} defined as in (12), the block columns of the matrix $\mathbf{Q} = \mathbf{Y} - \mathbb{S} \left(\mathbf{U} [\alpha \mathbf{e}_0 \ \mathbf{A}] \tilde{\mathbf{V}}^T \right)$ satisfy

$$\mathbf{Q}_i = (\mathbf{\Pi}_i \otimes \mathbf{I}_2) [\mathbf{p} \mathbf{1}_R^T \ \mathbf{q}], \quad \mathbf{p} = \mathbf{e}_0 - \mathbb{S}(\mathbf{x}), \quad \mathbf{q} = \mathbf{e}_1 - \mathbb{S}(\mathbf{y}),$$

where \mathbf{x}, \mathbf{y} are defined as

$$\begin{aligned}\mathbf{x} &:= \frac{1}{2M} \left(\sqrt{\frac{2}{R+1}} \left(a + \frac{\alpha}{\sqrt{R}} \right) \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} \bar{\Phi}_M \mathbf{e}_m \otimes \frac{\mathbf{a}_{m,0}}{\sqrt{R}} \right) \\ \mathbf{y} &:= \frac{1}{2M} \left(\sqrt{\frac{2}{R+1}} \left(-a + \alpha\sqrt{R} \right) \mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} + 2 \sum_{m=0}^{M-2} \bar{\Phi}_M \mathbf{e}_m \otimes \mathbf{a}_{m,1} \right).\end{aligned}\quad (30)$$

In addition, it also holds that $\mathbf{Q} = \mathbf{U} [\beta \mathbf{e}_0 \quad \mathbf{B}] [\tilde{\mathbf{v}} \quad \mathbf{V}]^T$, with parameters

$$\mathbf{B} = \text{diag}(b, \mathbf{B}_0, \dots, \mathbf{B}_{M-2}), \quad \mathbf{B}_i \in \mathbb{R}^{2 \times 2}, \quad (31)$$

$$b = \frac{1}{\sqrt{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (R\mathbf{p} - \mathbf{q}), \quad \beta = \sqrt{\frac{R}{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (\mathbf{p} + \mathbf{q}), \quad (32)$$

$$\mathbf{B}_i = ((\bar{\Phi}_M \mathbf{e}_i)^T \otimes \mathbf{I}_2) [\sqrt{R}\mathbf{p} \quad \mathbf{q}] \quad \tilde{\mathbf{v}} = \frac{\mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ R \end{pmatrix}}{\sqrt{MR(R+1)}} \quad (33)$$

Proof All parts directly follow from Theorem 13, Theorem 15, Theorem 16, Theorem 18, and Theorem 14. Note that the additional $\alpha \mathbf{u}_0 \tilde{\mathbf{v}}_0^T$ introduces an additional

$$\frac{\alpha}{\sqrt{R+1}} \left(\frac{1}{M} \mathbf{1}_M \mathbf{1}_M^T \right) \otimes \left[\frac{1}{\sqrt{2R}} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \mathbf{1}_R^T \quad -\frac{\sqrt{R}}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right]$$

into the expansion (16), which has been accounted for in the the adjusted \mathbf{x}, \mathbf{y} expressions. \blacksquare

Theorem 20 (Restatement of Theorem 1) Fix any $\alpha, a \in \mathbb{R}$ and $\mathbf{a}_+, \mathbf{a}_- \in \mathbb{R}^{M-1}$. With these, define vectors $\mathbf{p}, \mathbf{q} \in \mathbb{R}^K$ as follows

$$\begin{aligned}\mathbf{p}_+ &:= \mathbf{e}_0 - \mathbb{S} \left(\mathbf{x}_+ \otimes \begin{pmatrix} 1 \\ 0 \end{pmatrix} \right), & \mathbf{x}_+ &:= \frac{1}{M} \left(\sqrt{\frac{2}{R+1}} \left(a + \frac{\alpha}{\sqrt{R}} \right) \mathbf{1}_M + \frac{1}{\sqrt{R}} \bar{\Phi}_M \mathbf{a}_+ \right), \\ \mathbf{p}_- &:= \mathbf{e}_1 - \mathbb{S} \left(\mathbf{x}_- \otimes \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right), & \mathbf{x}_- &:= \frac{1}{M} \left(\sqrt{\frac{2}{R+1}} \left(a - \alpha\sqrt{R} \right) \mathbf{1}_M + \bar{\Phi}_M \mathbf{a}_- \right),\end{aligned}\quad (34)$$

Finally, for diagonal matrix $\mathbf{A} = \text{diag}(a, \mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-) \in \mathbb{R}^{(K-1) \times (K-1)}$, define the matrix

$$\mathbf{Q} = \mathbf{Y} - \mathbb{S} \left(\mathbf{U} [\alpha \mathbf{e}_0 \quad \mathbf{A}] \tilde{\mathbf{V}}^T \right) \in \mathbb{R}^{K \times (M)(R+1)}.$$

Then, the following two statements are true about \mathbf{Q} :

1. The M block columns of the matrix \mathbf{Q} satisfy

$$\mathbf{Q}_i = (\Pi_i \otimes \mathbf{I}_2) [\mathbf{p}_+ \mathbf{1}_R^T \quad \mathbf{p}_-] \in \mathbb{R}^{K \times (R+1)}, \quad i \in [M], \quad \mathbf{p}_+ = \mathbf{e}_0 - \mathbb{S}(\mathbf{x}), \quad \mathbf{p}_- = \mathbf{e}_1 - \mathbb{S}(\mathbf{y})$$

2. \mathbf{Q} admits a decomposition as follows

$$\mathbf{Q} = \mathbf{U} [\beta \mathbf{e}_0 \quad \mathbf{B}] \tilde{\mathbf{V}}^\top$$

where $\beta \in \mathbb{R}$ and the diagonal matrix $\mathbf{B} := \text{diag}(b, \mathbf{b}_+ \otimes \mathbf{f}_+ + \mathbf{b}_- \otimes \mathbf{f}_-)$ with $\mathbf{b}_+ \in \mathbb{R}^{M-1}$, and $\mathbf{b}_- \in \mathbb{R}^{M-1}$ are given in terms of $\alpha, a, \mathbf{a}_+, \mathbf{a}_-$ as follows

$$\begin{aligned} \beta &= \sqrt{\frac{R}{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (\mathbf{p}_+ + \mathbf{p}_-), & b &= \frac{1}{\sqrt{2(R+1)}} (\bar{\Phi}_K \mathbf{e}_0)^T (R\mathbf{p}_+ - \mathbf{p}_-), \\ \mathbf{b}_+ &= \sqrt{R} (\bar{\Phi}_M \otimes \mathbf{f}_+)^T \mathbf{p}_+, & \mathbf{b}_- &= (\bar{\Phi}_M \otimes \mathbf{f}_-)^T \mathbf{p}_-. \end{aligned} \quad (35)$$

D.3. Simplified form of expressions

Lemma 21 For β, b, \mathbf{b}_+ , and \mathbf{b}_- defined as in (35), the following simplified expressions hold:

$$\begin{aligned} \beta &= \sqrt{\frac{R}{2(R+1)}} \left(\frac{2M}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} - \frac{M}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M} \right) \\ b &= \frac{1}{\sqrt{2(R+1)}} \left(\frac{RM}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} + \frac{M}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M} \right) \\ \mathbf{b}_+ &= \sqrt{R} \left(\mathbf{1}_{M-1} - \frac{\bar{\Phi}_M^T \exp(\mathbf{x}_+)}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} \right), & \mathbf{b}_- &= \mathbf{1}_{M-1} - \frac{\bar{\Phi}_M^T \exp(\mathbf{x}_-)}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M}. \end{aligned}$$

It holds in addition that $\mathbf{b}_\pm \geq \mathbf{0}$, $\mathbf{b}_+ = \sqrt{R}\mathbf{1}$ if $\mathbf{a}_+ = \mathbf{0}$ and $\mathbf{b}_- = \mathbf{1}$ if $\mathbf{a}_- = \mathbf{0}$, and $\mathbf{b}_+ \leq \sqrt{R}\mathbf{1}$ if $\mathbf{a}_+ \geq \mathbf{0}$ and $\mathbf{b}_- \leq \mathbf{1}$ if $\mathbf{a}_- \geq \mathbf{0}$.

Proof Notice the fact that the vectors $\mathbf{p}_+, \mathbf{p}_-$ can be written in terms of $\mathbf{x}_+, \mathbf{x}_-$ as

$$\mathbf{p} = \mathbf{e}_0 \otimes \mathbf{f}_+ - \frac{\exp(\mathbf{x}_+) \otimes \mathbf{f}_+ + \mathbf{1}_M \otimes \mathbf{f}_-}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M}, \quad \mathbf{q} = \mathbf{e}_0 \otimes \mathbf{f}_- - \frac{\exp(\mathbf{x}_-) \otimes \mathbf{f}_- + \mathbf{1}_M \otimes \mathbf{f}_+}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M}.$$

Using (35), for \mathbf{b}_+ and \mathbf{b}_- we have,

$$\mathbf{b}_+ = \sqrt{R} \left(\mathbf{1}_{M-1} - \frac{\bar{\Phi}_M^T \exp(\mathbf{x}_+)}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} \right), \quad \mathbf{b}_- = \mathbf{1}_{M-1} - \frac{\bar{\Phi}_M^T \exp(\mathbf{x}_-)}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M}.$$

In order to do the same for β and b , we first write the simplified form of $(\bar{\Phi} \mathbf{e}_0)^T \mathbf{p}_+$ and $(\bar{\Phi} \mathbf{e}_0)^T \mathbf{p}_-$:

$$\begin{aligned} (\bar{\Phi} \mathbf{e}_0)^T \mathbf{p} &= (\mathbf{1}_M^T \otimes (1 \quad -1)) \mathbf{p} = 1 - \frac{\mathbf{1}_M^T \exp(\mathbf{x}_+) - M}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} = \frac{2M}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M}, \\ (\bar{\Phi} \mathbf{e}_0)^T \mathbf{p}_- &= (\mathbf{1}_M^T \otimes (1 \quad -1)) \mathbf{q} = -1 - \frac{-\mathbf{1}_M^T \exp(\mathbf{x}_-) + M}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M} = \frac{-2M}{\mathbf{1}_M^T \exp(\mathbf{x}_-) + M}. \end{aligned}$$

By substituting the above in the equations for b, β , the desired relations are proved.

Our next step is to prove the auxiliary equations and inequalities. Note that $\mathbf{a}_\pm = \mathbf{0}$ would result in $\exp(\mathbf{x}_+) = \exp(\gamma/M)\mathbf{1}$, $\exp(\mathbf{x}_-) = \exp(\delta/M)\mathbf{1}$. In both cases, $\bar{\Phi}_M^T \mathbf{1} = \mathbf{0}$, and thus $\mathbf{b}_+ = \sqrt{R}\mathbf{1}_{M-1}$, $\mathbf{b}_- = \mathbf{1}_{M-1}$. Also,

$$\mathbf{b}_+ = \sqrt{R} \left(\frac{[\mathbf{0}_{M-1} \quad \Psi_M]^T \exp(\mathbf{x}_+)}{\mathbf{1}_M^T \exp(\mathbf{x}_+) + M} \right) \geq \mathbf{0},$$

since Ψ_M only has entries in $\{0, 2\}^{(M-1) \times (M-1)}$. The same reasoning applies to \mathbf{b}_- to show $\mathbf{b}_- \geq \mathbf{0}$. Finally, to prove that $\mathbf{b}_- \leq \mathbf{1}_{M-1}$ for $\mathbf{a}_- \geq \mathbf{0}$, we proceed to show that $\bar{\Phi}_M^T \exp(\mathbf{x}_-) \geq \mathbf{0}$. Expanding the exponential would result yield

$$\exp(\bar{\Phi}_M \mathbf{a}_-) = e^{\delta/M} \bigodot_{i=1}^{M-1} \left(\frac{1}{2} \cosh\left(\frac{a_{-,i-1}}{M}\right) + \frac{1}{2} \sinh\left(\frac{a_{-,i-1}}{M}\right) \Phi \mathbf{e}_i \right).$$

Since $\mathbf{a}_- \geq \mathbf{0}$, all hyperbolic functions are positive. Also, given that $(\Phi \mathbf{e}_i) \odot (\Phi \mathbf{e}_j) = \Phi \mathbf{e}_{i \oplus j}$, we conclude that $\exp(\mathbf{x}_-) = \bar{\Phi}_M \mathbf{r}$ for some $\mathbf{r}_- \geq \mathbf{0}$. As a result, $\bar{\Phi}_M^T \exp(\mathbf{x}_-) = M \mathbf{r}_- \geq \mathbf{0}$. The same argument applies for \mathbf{b}_+ to show that for $\mathbf{b}_+ \leq \sqrt{R}\mathbf{1}$ for $\mathbf{a}_+ \geq \mathbf{0}$. \blacksquare

Appendix E. Proofs for Regularization-path

E.1. Proof of Theorem 2 and basic properties of the solution

We begin with a proof sketch before giving the formal argument. The theorem splits into three claims. The theorem can be understood by splitting in three claims. The first—that if \mathbf{Z}_λ has the form given in the theorem statement, then $(\mathbf{Z}\text{-UFM}_\lambda)$ reduces to the vector problem (3)—is direct. The second, and substantive, claim is that this restriction is *tight*: the global minimizer of $(\mathbf{Z}\text{-UFM}_\lambda)$ lies on the $(\mathbf{U}, \tilde{\mathbf{V}})$ subspace $\forall \lambda > 0$. This is where Thm. 1 enters. The optimality condition for $(\mathbf{Z}\text{-UFM}_\lambda)$ is $-\nabla_{\mathbf{Z}} \mathcal{L}(\mathbf{Z}) = \mathbf{Y} - \mathbb{S}(\mathbf{Z}) \in \lambda \partial \|\mathbf{Z}\|_*$. For any candidate $\mathbf{Z} = \mathbf{U} \text{diag}_{+1}(\mathbf{a}) \tilde{\mathbf{V}}^T$, Thm. 1 rewrites the loss gradient as $\mathbf{U} \text{diag}_{+1}(\mathcal{F}(\mathbf{a})) \tilde{\mathbf{V}}^T$, so the condition becomes $\mathbf{U} \text{diag}_{+1}(\mathcal{F}(\mathbf{a})) \tilde{\mathbf{V}}^T \in \lambda \partial \|\mathbf{U} \text{diag}_{+1}(\mathbf{a}) \tilde{\mathbf{V}}^T\|_*$. By orthogonality of $\mathbf{U}, \tilde{\mathbf{V}}$, this is equivalent to $\mathcal{F}(\mathbf{a}) \in \lambda \partial \|\mathbf{a}\|_{+1}$, where $\|\cdot\|_{+1}$ denotes the regularizer in (3). Since this is precisely the optimality condition for (3), any of its solution certifies a solution of $(\mathbf{Z}\text{-UFM}_\lambda)$. Uniqueness on both sides (due to strict convexity of the CE loss on the subspace orthogonal to $\mathbf{1}$) closes the loop. Finally, the third claim, the reduction to four variables follows by symmetry: By leveraging the structure of the $\mathbf{U}, \tilde{\mathbf{V}}$ subspace with respect to the Sylvester-Hadamard matrix, we show that the objective in (3) is invariant to swapping the coordinates of \mathbf{a}_\pm . Because the minimizer is unique, this symmetry forces all elements within these vectors to be identical.

We now formalize each claim in turn. After that, at the end of this section, we further derive basic but important properties of the solution, which we use in proofs of later stages.

Our first step is to characterize the regularization path solution, while assuming it is constrained to the Hadamard factorization discussed in previous sections,

Theorem 22 (KKT conditions for regularized CE in the diagonal basis) *Consider the optimization problem*

$$\min_{\alpha \in \mathbb{R}, \mathbf{A} \in \mathbb{R}^{(k-1) \times (k-1)}} \text{CE}(\mathbf{L}) + \lambda \|\mathbf{L}\|_*, \quad \text{s. t. } \mathbf{L} = \mathbf{U} [\alpha \mathbf{e}_0 \quad \mathbf{A}] \tilde{\mathbf{V}}^T, \quad (36)$$

where $\lambda > 0$, the matrix \mathbf{A} is restricted to the diagonal form $\mathbf{A} = \text{diag}(a, \mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-)$. Then, the unique minimizer of (36) satisfies

$$\mathbf{Y} - \mathbb{S}(\mathbf{L}) = \mathbf{U} [\beta \mathbf{e}_0 \quad \mathbf{B}] \tilde{\mathbf{V}}^T, \quad \mathbf{B} = \text{diag}(b, \mathbf{b}_+ \otimes \mathbf{f}_+ + \mathbf{b}_- \otimes \mathbf{f}_-), \quad (37)$$

where

$$(\beta, b) \in \lambda \partial \sqrt{\alpha^2 + a^2}, \quad (38)$$

$$\mathbf{b}_\pm \in \lambda \partial \|\mathbf{a}_\pm\|_1, \quad (39)$$

Explicitly, this yields

$$(\beta, b) = \lambda \frac{(\alpha, a)}{\sqrt{\alpha^2 + a^2}} \text{ if } (\alpha, a) \neq (0, 0), \quad \beta^2 + b^2 \leq \lambda^2 \text{ otherwise}, \quad (40)$$

$$\mathbf{b}_\pm[j] = \lambda \text{sign}(\mathbf{a}_\pm[j]) \text{ if } \mathbf{a}_\pm[j] \neq 0, \quad |\mathbf{b}_\pm[j]| \leq \lambda \text{ otherwise}. \quad (41)$$

Proof Theorem 20 guarantees that (37) for any α, \mathbf{A} , without imposing any constraints on β, \mathbf{B} yet. Denote the linear mapping $\mathcal{T}(\mathbf{X}) = \mathbf{U} \mathbf{X} \tilde{\mathbf{V}}^T$ and its adjoint map as $\mathcal{T}^*(\mathbf{X}) = \mathbf{U}^T \mathbf{X} \tilde{\mathbf{V}}$. Since $\mathbf{U}, \tilde{\mathbf{V}}$ are fixed, the objective in (36) is convex in the free parameters (α, \mathbf{A}) , and the KKT condition reduces to first-order subgradient optimality:

$$\mathbf{0} \in \nabla_{\alpha, \mathbf{A}} \text{CE}(\mathcal{T}([\alpha \mathbf{e}_0 \quad \mathbf{A}])) + \lambda \partial_{\alpha, \mathbf{A}} \|\mathcal{T}([\alpha \mathbf{e}_0 \quad \mathbf{A}])\|_*. \quad (42)$$

The nuclear norm in the above expression is invariant with respect to the transformation \mathcal{T} , and using the subgradient chain rule for affine transformations [7, Theorem 3.43 (b)], the above expression can be simplified to

$$\mathbf{U}^T (\mathbf{Y} - \mathbb{S}(\mathcal{T}([\alpha \mathbf{e}_0 \quad \mathbf{A}])) \tilde{\mathbf{V}} = [\beta \mathbf{e}_0 \quad \mathbf{B}] \in \lambda \partial \|\alpha \mathbf{e}_0 \quad \mathbf{A}\|_*. \quad (43)$$

Given the diagonal structure of \mathbf{A} , the nuclear norm expression simplifies as

$$\|\alpha \mathbf{e}_0 \quad \mathbf{A}\|_* = \sqrt{\alpha^2 + a^2} + \|\mathbf{a}_+\|_1 + \|\mathbf{a}_-\|_1, \quad (44)$$

Applying the subgradient formula to (43) yields

$$(\beta, b) \in \lambda \partial \sqrt{\alpha^2 + a^2}, \quad \mathbf{b}_\pm \in \lambda \partial \|\mathbf{a}_\pm\|_1,$$

which are exactly (38) and (39). Uniqueness follows from the convexity of the nuclear norm and the strict convexity of the cross-entropy over the subspace of logit matrices with zero-sum columns, a constraint that our factorization automatically satisfies. \blacksquare

Our next step is to show that removing the factorization constraint $\mathbf{L} = \mathbf{U} [\alpha \mathbf{e}_0 \quad \mathbf{A}] \tilde{\mathbf{V}}^T$ does not change the solution via a KKT certificate proof, and thus the solution of $(\mathbf{Z}\text{-UFM}_\lambda)$ admits the Hadamard decomposition. Before proceeding with the proof, two utility lemmas are introduced:

Lemma 23 Let $\mathbf{U} \in \mathbb{R}^{m \times r}$ and $\mathbf{V} \in \mathbb{R}^{n \times s}$ be matrices with orthonormal columns, i.e., they satisfy

$$\mathbf{U}^\top \mathbf{U} = \mathbf{I}_r, \quad \mathbf{V}^\top \mathbf{V} = \mathbf{I}_s.$$

Then, for any $\mathbf{C} \in \mathbb{R}^{r \times s}$ and any $\mathbf{D} \in \lambda \partial \|\mathbf{C}\|_*$,

$$\|\mathbf{UCV}^\top\|_* = \|\mathbf{C}\|_*, \quad \mathbf{UDV}^\top \in \lambda \partial \|\mathbf{UCV}^\top\|_*.$$

Proof Let $D \in \lambda \partial \|C\|_*$. By the definition of the subdifferential, for every $Z \in \mathbb{R}^{r \times s}$,

$$\lambda \|Z\|_* \geq \lambda \|C\|_* + \langle D, Z - C \rangle.$$

Fix any $Y \in \mathbb{R}^{m \times n}$, and set

$$Z = U^\top Y V.$$

Using the fact that multiplication by matrices with orthonormal columns does not increase the nuclear norm (see [9, Proposition IV.2.4]),

$$\|Y\|_* \geq \|U^\top Y V\|_*,$$

we obtain

$$\lambda \|Y\|_* \geq \lambda \|U^\top Y V\|_* \geq \lambda \|C\|_* + \langle D, U^\top Y V - C \rangle$$

Since U and V have orthonormal columns, the matrices UCV^\top and C have the same nonzero singular values. Therefore,

$$\|UCV^\top\|_* = \|C\|_*,$$

and hence,

$$\lambda \|Y\|_* \geq \lambda \|UCV^\top\|_* + \langle UDV^\top, Y - UCV^\top \rangle.$$

Since this holds for every Y , it follows that

$$UDV^\top \in \lambda \partial \|UCV^\top\|_*.$$

■

Lemma 24 Let $L \in \mathbb{R}^{k \times n}$ be a matrix satisfying $L^T \mathbf{1}_k = \mathbf{0}$. Then, $c = \mathbf{0}$ is the unique solution to

$$\min_{c \in \mathbb{R}^n} \|L + \mathbf{1}_k c^T\|_*.$$

Proof The nuclear norm expression can be simplified using $L^T \mathbf{1} = \mathbf{0}$ to reach

$$\|L + \mathbf{1}_k c^T\|_* = \text{tr} \sqrt{L^T L + k c c^T} \geq \text{tr} \sqrt{L^T L} = \|L\|_*,$$

which establishes $c = \mathbf{0}$ as a minimizer. Next, Weyl's monotonicity theorem [9, Corollary 3.2.3] results in $\lambda_i(L^T L + k c c^T) \geq \lambda_i(L^T L)$, where λ_i denotes the i -th eigenvalue of each matrix. This, alongside the fact that $\text{tr}(L^T L + k c c^T) > \text{tr}(L^T L)$ for nonzero c , results in at least one eigenvalue to strictly increase. Consequently, for $c \neq \mathbf{0}$

$$\|L + \mathbf{1}_k c^T\|_* = \text{tr} \sqrt{L^T L + k c c^T} > \text{tr} \sqrt{L^T L} = \|L\|_*,$$

which proves the uniqueness. ■

Proposition 25 (KKT certificate for the nuclear-norm-regularized CE problem) *Let*

$$\mathbf{L}^* = \mathbf{U} \begin{bmatrix} \alpha \mathbf{e}_0 & \mathbf{A} \end{bmatrix} \tilde{\mathbf{V}}^T, \quad \mathbf{A} = \text{diag}(a, \mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-),$$

and

$$\mathbf{Y} - \mathbb{S}(\mathbf{L}^*) = \mathbf{U} \begin{bmatrix} \beta \mathbf{e}_0 & \mathbf{B} \end{bmatrix} \tilde{\mathbf{V}}^T, \quad \mathbf{B} = \text{diag}(b, \mathbf{b}_+ \otimes \mathbf{f}_+ + \mathbf{b}_- \otimes \mathbf{f}_-),$$

be the decomposition given by Theorem 20. If

$$(\beta, b) \in \lambda \partial \sqrt{\alpha^2 + a^2}, \quad (45)$$

$$\mathbf{b}_\pm[j] \in \lambda \partial |\mathbf{a}_\pm[j]|, \quad j \in [k/2] \quad (46)$$

then \mathbf{L}^* is the unique minimizer of the convex optimization problem

$$\min_{\mathbf{L} \in \mathbb{R}^{k \times n}} \text{CE}(\mathbf{L}) + \lambda \|\mathbf{L}\|_*.$$

Proof Consider the objective

$$f(\mathbf{L}) := \text{CE}(\mathbf{L}) + \lambda \|\mathbf{L}\|_*.$$

Notice that the cross-entropy term is invariant to centering each column; that is, $\text{CE}(\mathbf{L}) = \text{CE}(\widehat{\mathbf{L}})$, where $\widehat{\mathbf{L}}$ is the \mathbf{L} with all columns centered around zero. Moreover, the strict convexity of cross-entropy over this subspace alongside the overall objective convexity requires that $\widehat{\mathbf{L}}_1 = \widehat{\mathbf{L}}_2$ for any two solutions $\mathbf{L}_1, \mathbf{L}_2$.

Therefore, all solutions are of the form $\mathbf{L} = \mathbf{L}^* + \mathbf{1}_k \mathbf{r}^T$ for a unique $\widehat{\mathbf{L}}$ that satisfies $\widehat{\mathbf{L}}^T \mathbf{1}_k = \mathbf{0}$, and all possible \mathbf{r} are determined from $\min_{\mathbf{r} \in \mathbb{R}^n} \|\widehat{\mathbf{L}} + \mathbf{1}_k \mathbf{r}^T\|_*$. This minimization is shown in Theorem 24 to have $\mathbf{r} = \mathbf{0}$ as its only solution, which proves the solution is unique and satisfies $\mathbf{L}^T \mathbf{1} = \mathbf{0}$. The proposed solution automatically satisfies this requirement; therefore, it suffices to verify the first-order optimality condition

$$\mathbf{0} \in \nabla \text{CE}(\mathbf{L}^*) + \lambda \partial \|\mathbf{L}^*\|_* = \mathbb{S}(\mathbf{L}^*) - \mathbf{Y} + \lambda \partial \|\mathbf{L}^*\|_*$$

Define

$$\mathbf{C}^* := \begin{bmatrix} \alpha \mathbf{e}_0 & \mathbf{A} \end{bmatrix}, \quad \mathbf{D}^* := \begin{bmatrix} \beta \mathbf{e}_0 & \mathbf{B} \end{bmatrix}.$$

Then, by definition in the statement of the proposition,

$$\mathbf{L}^* = \mathbf{U} \mathbf{C}^* \tilde{\mathbf{V}}^T, \quad \mathbf{Y} - \mathbb{S}(\mathbf{L}^*) = \mathbf{U} \mathbf{D}^* \tilde{\mathbf{V}}^T, \quad \mathbf{D}^* \in \lambda \partial \|\mathbf{C}^*\|_*,$$

Since \mathbf{U} and $\tilde{\mathbf{V}}$ have orthonormal columns, the nuclear norm is invariant under multiplication by these matrices, and its subdifferential transforms equivariantly. Hence, using Theorem 23, we obtain

$$\mathbf{D}^* \in \lambda \partial \|\mathbf{C}^*\|_* \implies \mathbf{U} \mathbf{D}^* \tilde{\mathbf{V}}^T \in \lambda \partial \|\mathbf{U} \mathbf{C}^* \tilde{\mathbf{V}}^T\|_*.$$

Thus \mathbf{L}^* satisfies the KKT condition for the convex problem

$$\min_{\mathbf{L} \in \mathbb{R}^{k \times n}} \text{CE}(\mathbf{L}) + \lambda \|\mathbf{L}\|_*,$$

■

Now that the Hadamard factorization for the regularized problem (\mathbf{Z} -UFM $_\lambda$) has been proven, basic properties of the solution are discussed. These properties include sparsity patterns and useful nonnegativity results for variables in the optimal solution.

Lemma 26 (Basic properties of regularization path solution) *The optimal solution to (36) has the following properties:*

1. *The following quantities are nonnegative: $a_\lambda \geq 0$, $\mathbf{a}_+ \geq \mathbf{0}$, $\mathbf{a}_- \geq \mathbf{0}$, and $\gamma_\lambda, \delta_\lambda \geq 0$.*
2. *$\mathbf{a}_{\lambda,+} = \mathbf{0} \iff \lambda \geq \sqrt{R}$ and $\mathbf{a}_{\lambda,-} = \mathbf{0} \iff \lambda \geq 1$.*
3. *A threshold $\lambda_{K,R} \leq \sqrt{\frac{R+1}{2}}$ exists such that $(\alpha_\lambda, a_\lambda) = (0, 0)$ for $\lambda \geq \lambda_{K,R}$.*

Proof

1. Notice that $b \geq 0$, $\mathbf{b}_\pm \geq \mathbf{0}$ regardless of the problem variables, and $b + \frac{\beta}{\sqrt{R}} = \sqrt{\frac{R+1}{2}} X_+ > 0$ and $b - \beta\sqrt{R} = \sqrt{\frac{R+1}{2}} X_- > 0$. The KKT condition for (α, a) is either $(\alpha, a) = (0, 0)$ or $(\beta, b) = \frac{\lambda}{\sqrt{a^2 + \alpha^2}}(\alpha, a)$. It follows that in either case, $a \geq 0$, and the definition of (γ, δ) in (14) implies $\gamma, \delta \geq 0$.
Likewise, the conditions for \mathbf{a}_\pm state either $\mathbf{b}_\pm[j] = \lambda \text{sign}(\mathbf{a}_\pm[j])$ or $\mathbf{a}_\pm[j] = 0 \forall j \in [M]$. Since $\mathbf{b}_\pm \geq \mathbf{0}$, this implies $\mathbf{a}_\pm \geq \mathbf{0}$.
2. We shall prove the second statement by contradiction. Note that $\mathbf{a}_\pm \geq \mathbf{0}$ yields $\mathbf{b}_+ \leq \sqrt{R}\mathbf{1}$, $\mathbf{b}_- \leq \mathbf{1}$. Suppose for $\lambda > \sqrt{R}$, there exists $j \in [M]$ such that $\mathbf{a}_+[j] \neq 0$. The KKT conditions would give $\lambda = \mathbf{b}_+[j] \leq \sqrt{R}$, which leads to a contradiction. To prove the other direction, suppose that $\mathbf{a}_+ = \mathbf{0}$ for some $\lambda < \sqrt{R}$, which leads to $\mathbf{b}_+ = \sqrt{R}\mathbf{1}$. However, the KKT conditions would require $\mathbf{b}_+ \leq \lambda\mathbf{1}$, leading to a contradiction. An analogous argument proves that $\mathbf{a}_- = \mathbf{0} \iff \lambda \leq 1$.
3. Based on the KKT equations for (α, a) either $(\alpha, a) = 0, \beta^2 + b^2 \leq \lambda^2$ or $\beta^2 + b^2 = \lambda^2$, where $\beta^2 + b^2 = \frac{1}{2}(RX_+^2 + X_-^2)$. Jensen's inequality on the strictly convex exponential function gives $\mathbf{1}^T \exp(\mathbf{x}_+) > Me^\gamma$, and likewise $\mathbf{1}^T \exp(\mathbf{x}_-) > Me^\delta$. Since $\gamma, \delta \geq 0$, the definition of X_\pm implies $X_\pm < 1 \Rightarrow \beta^2 + b^2 < \frac{R+1}{2}$. This in turn establishes $\beta^2 + b^2 = \lambda^2$ is not possible for $\lambda \geq \sqrt{\frac{R+1}{2}}$. ■

E.1.1. REDUCTION TO FOUR VARIABLES

We prove that the vectorized regularization path reduces to an optimization over four variables.

Proposition 27 (Reduction to four variables) *There exist scalars $a_{+,\lambda}$ and $a_{-,\lambda}$ such that the minimizer of Eq. (3) satisfies $\mathbf{a}_{\pm,\lambda} = a_{\pm,\lambda}\mathbf{1}_{M-1}$.*

Proof Suppose $(\alpha_\lambda, a_\lambda, \mathbf{a}_{+,\lambda}, \mathbf{a}_{-,\lambda})$ is a minimizer. Let $r, s \in [M-1]$. By Proposition 28 (see also 12), there exist \mathbf{P}, \mathbf{C} such that $\mathbf{P}\bar{\Phi}_M = \bar{\Phi}_M\mathbf{C}$ and $\mathbf{C}\mathbf{e}_r = \mathbf{e}_s$, and, $(\alpha_\lambda, a_\lambda, \mathbf{C}\mathbf{a}_{+,\lambda}, \mathbf{C}\mathbf{a}_{-,\lambda})$ is also a minimizer. By uniqueness of the minimizer of Eq. (3) it must be that

$$\mathbf{C}\mathbf{a}_{\pm,\lambda} = \mathbf{a}_{\pm,\lambda} \implies (\mathbf{a}_{\pm,\lambda})_s = (\mathbf{a}_{\pm,\lambda})_r.$$

Since r, s were arbitrary, all coordinates of $\mathbf{a}_{\pm,\lambda}$ are equal completing the proof ■

The key observation behind the proof is showing invariance of the optimization objective under admissible Hadamard permutations.

Proposition 28 (Reduced-objective invariance) *Let $\mathbf{C} \in \mathbb{R}^{(M-1) \times (M-1)}$ be a permutation matrix for which there exists a permutation matrix $\mathbf{P} \in \mathbb{R}^{M \times M}$ satisfying*

$$\mathbf{P}\bar{\Phi}_M = \bar{\Phi}_M \mathbf{C}. \quad (47)$$

For example, this is guaranteed to exist by Lemma 12. Then the vectorized objective in Eq. (3)

$$f_\lambda(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) := \mathcal{L}\left(\mathbf{U} \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \tilde{\mathbf{V}}^\top\right) + \lambda \left(\sqrt{\alpha^2 + a^2} + \|\mathbf{a}_+\|_1 + \|\mathbf{a}_-\|_1\right)$$

satisfies

$$f_\lambda(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) = f_\lambda(\alpha, a, \mathbf{C}\mathbf{a}_+, \mathbf{C}\mathbf{a}_-).$$

Proof Define the induced class and sample permutations

$$\mathbf{\Pi} := \mathbf{P} \otimes \mathbf{I}_2, \quad \mathbf{\Omega} := \mathbf{P} \otimes \mathbf{I}_{R+1}.$$

Since every majority/minority class-pair block has the same structure, the label matrix \mathbf{Y} (Eq. (1)) satisfies

$$\mathbf{\Pi}\mathbf{Y}\mathbf{\Omega}^\top = \mathbf{Y}.$$

Using the fact that $\mathbf{\Pi}, \mathbf{\Omega}$ are permutation matrices, together with the permutation equivariance of the softmax and the entrywise logarithm, we obtain for every logit matrix \mathbf{Z} :

$$\mathbf{Y} \odot \log \mathbb{S}(\mathbf{\Pi}\mathbf{Z}\mathbf{\Omega}^\top) = \mathbf{\Pi}\mathbf{Y}\mathbf{\Omega}^\top \odot \mathbf{\Pi} \log \mathbb{S}(\mathbf{Z}) \mathbf{\Omega}^\top = \mathbf{\Pi}(\mathbf{Y} \odot \log \mathbb{S}(\mathbf{Z}))\mathbf{\Omega}^\top.$$

Recalling that $\mathcal{L}(\mathbf{Z}) = -\mathbf{1}^\top(\mathbf{Y} \odot \log \mathbb{S}(\mathbf{Z}))\mathbf{1}$ yields

$$\mathcal{L}(\mathbf{\Pi}\mathbf{Z}\mathbf{\Omega}^\top) = \mathcal{L}(\mathbf{Z}). \quad (48)$$

We next check how this symmetry acts on the structured parametrization. Recall that

$$\mathbf{U} = \begin{bmatrix} \frac{1}{\sqrt{K}}\mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} & \frac{1}{\sqrt{M}}\bar{\Phi}_M \otimes \mathbf{I}_2 \end{bmatrix}.$$

For the first block, since $\mathbf{P}\mathbf{1}_M = \mathbf{1}_M$:

$$\mathbf{\Pi} \left(\frac{1}{\sqrt{K}}\mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right) = \frac{1}{\sqrt{K}}(\mathbf{P}\mathbf{1}_M) \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \frac{1}{\sqrt{K}}\mathbf{1}_M \otimes \begin{pmatrix} 1 \\ -1 \end{pmatrix},$$

For the Hadamard block, using the assumption (47)

$$\mathbf{\Pi}(\bar{\Phi}_M \otimes \mathbf{I}_2) = (\mathbf{P}\bar{\Phi}_M) \otimes \mathbf{I}_2 = (\bar{\Phi}_M \mathbf{C}) \otimes \mathbf{I}_2 = (\bar{\Phi}_M \otimes \mathbf{I}_2)(\mathbf{C} \otimes \mathbf{I}_2).$$

Hence

$$\mathbf{\Pi}\mathbf{U} = \mathbf{U}\mathbf{D}, \quad \mathbf{D} := \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{C} \otimes \mathbf{I}_2 \end{bmatrix}. \quad (49)$$

Now write $\tilde{\mathbf{V}} = [\tilde{\mathbf{v}} \quad \mathbf{v}_0 \quad \mathbf{V}_H]$, where

$$\tilde{\mathbf{v}} = \frac{1}{\sqrt{MR(R+1)}} \mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ R \end{pmatrix}, \quad \mathbf{v}_0 = \frac{1}{\sqrt{M(R+1)}} \mathbf{1}_M \otimes \begin{pmatrix} \mathbf{1}_R \\ -1 \end{pmatrix},$$

and

$$\mathbf{V}_H = \frac{1}{\sqrt{M}} \bar{\Phi}_M \otimes M, \quad M := \begin{pmatrix} \mathbf{1}_R/\sqrt{R} & \mathbf{0}_R/\sqrt{R} \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{(R+1) \times 2}.$$

The first two columns are invariant because they are constant over the M class-pair blocks, i.e., $\Omega \tilde{\mathbf{v}} = \tilde{\mathbf{v}}$, and $\Omega \mathbf{v}_0 = \mathbf{v}_0$. For the Hadamard block, again using (47)

$$\Omega \mathbf{V}_H (P \bar{\Phi}_M) \otimes M = \frac{1}{\sqrt{M}} (\bar{\Phi}_M C) \otimes M = \left(\frac{1}{\sqrt{M}} \bar{\Phi}_M \otimes M \right) (C \otimes I_2).$$

Therefore

$$\Omega \tilde{\mathbf{V}} = \tilde{\mathbf{V}} \mathbf{E}, \quad \mathbf{E} := \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & C \otimes I_2 \end{bmatrix}. \quad (50)$$

Combining Eq. (49) and (50) yields

$$\Pi U \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \tilde{\mathbf{V}}^\top \Omega^\top = U D \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \mathbf{E}^\top \tilde{\mathbf{V}}^\top.$$

Recall now that

$$\text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) = [\alpha e_0 \quad \text{diag}(a, \mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-)],$$

Thus, under pre/post multiplying by D/E the first column αe_0 and the scalar diagonal entry a are unchanged, while

$$(C \otimes I_2) \text{diag}(\mathbf{a}_+ \otimes \mathbf{f}_+ + \mathbf{a}_- \otimes \mathbf{f}_-) (C \otimes I_2)^\top = \text{diag}((C \mathbf{a}_+) \otimes \mathbf{f}_+ + (C \mathbf{a}_-) \otimes \mathbf{f}_-),$$

implying $D \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \mathbf{E}^\top = \text{diag}_{+1}(\alpha, a, C \mathbf{a}_+, C \mathbf{a}_-)$.

Tracing back the identities, we have shown that

$$\begin{aligned} \mathcal{L}(U \text{diag}_{+1}(\alpha, a, C \mathbf{a}_+, C \mathbf{a}_-) \mathbf{V}^\top) &= \mathcal{L}(U D \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \mathbf{E}^\top \mathbf{V}^\top) \\ &= \mathcal{L}(\Pi U \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \tilde{\mathbf{V}}^\top \Omega^\top) \\ &= \mathcal{L}(U \text{diag}_{+1}(\alpha, a, \mathbf{a}_+, \mathbf{a}_-) \tilde{\mathbf{V}}^\top), \end{aligned}$$

where the last line uses Eq. (48). Thus the CE term in f_λ is invariant under $(\mathbf{a}_+, \mathbf{a}_-) \mapsto (C \mathbf{a}_+, C \mathbf{a}_-)$. The regularizer is also clearly invariant because C is a permutation matrix. This proves the desired. \blacksquare

E.2. Proof of Theorem 3: Analytic solution

In this subsection we solve the four-dimensional convex problem from Eq. (3) explicitly. The argument has three steps. First, we isolate the scalar relations that encode the active KKT system. Second, we show that these scalar equations admit a unique solution in the relevant λ -ranges. Third, we substitute these scalar solutions back into the KKT conditions and verify the phase-wise formulas stated in the main theorem. This organization separates the existence of the scalar roots from the algebra needed to recover the optimization variables.

First, we prove a utility lemma, which shows that whenever $\mathbf{b}_\pm \propto \mathbf{1}$, $\mathbf{a}_\pm \propto \mathbf{1}_{M-1}$, and determine the multiplication coefficient as well:

Lemma 29 *Let $0 \leq t \leq 1$, $\gamma \in \mathbb{R}$, and suppose $\mathbf{r} \in \mathbb{R}^M$ satisfies*

$$\begin{aligned} \mathbf{r} &= \exp\left(\frac{1}{M}(\gamma\mathbf{1}_M + \bar{\Phi}_M\boldsymbol{\eta})\right), \\ \bar{\Phi}_M^\top \mathbf{r} &= t(M + \mathbf{1}_M^\top \mathbf{r})\mathbf{1}_{M-1}, \end{aligned} \quad (51)$$

Then, the solution to this equation satisfies

$$\boldsymbol{\eta} = (\gamma - M \log(s))\mathbf{1}_{M-1}, \quad s^{M-1}((1+\theta)s + \theta) = e^\gamma, \quad \theta := \frac{Mt}{1-t}, \quad (52)$$

where $0 < s \leq \exp(\gamma/M)$ is the unique solution of the scalar equation in (52)

Proof The first step is to multiply (51) by $\bar{\Phi}_m$ given that $\bar{\Phi}_M \bar{\Phi}_M^\top = M\mathbf{I}_M - \mathbf{1}_M \mathbf{1}_M^\top$ to reach

$$\begin{aligned} (M\mathbf{I}_M - \mathbf{1}_M \mathbf{1}_M^\top) \mathbf{r} &= t(M + \mathbf{1}_M^\top \mathbf{r})(\bar{\Phi}_M \mathbf{1}_{M-1}) \\ \Rightarrow \mathbf{r} &= \frac{\mathbf{1}_M^\top \mathbf{r}}{M} \mathbf{1}_M + t\left(1 + \frac{\mathbf{1}_M^\top \mathbf{r}}{M}\right)(M\mathbf{e}_0 - \mathbf{1}_M) \end{aligned}$$

After defining $\kappa := \frac{\mathbf{1}_M^\top \mathbf{r}}{M}$ and using that $\mathbf{1}_M^\top \log(\mathbf{r}) = \gamma$, $\bar{\Phi}_M^\top \log(\mathbf{r}) = \boldsymbol{\eta}$, standard algebra leads to

$$\log(\mathbf{r}) = \log\left(\frac{\kappa + t(1+\kappa)(M-1)}{\kappa - t(1+\kappa)}\right)\mathbf{e}_0 + \log(\kappa - t(1+\kappa))\mathbf{1}_M, \quad (53)$$

$$\boldsymbol{\eta} = \bar{\Phi}_M^\top \log(\mathbf{r}) = \log\left(\frac{\kappa + t(1+\kappa)(M-1)}{\kappa - t(1+\kappa)}\right)\mathbf{1}_M, \quad (54)$$

$$\gamma = \mathbf{1}_M^\top \log(\mathbf{r}) = \log\left((\kappa + t(1+\kappa)(M-1))(\kappa - t(1+\kappa))^{M-1}\right), \quad (55)$$

after which the desired result can be reached by defining $s := \kappa - t(1+\kappa)$ and rewriting the above equations in terms of s . Note that the uniqueness of s for a given γ is the result the monotonicity of $s^{M-1}((1+\theta)s + \theta)$ for $s > 0$, and $s \leq e^{\gamma/M}$ follows by $e^\gamma = s^{M-1}((1+\theta)s + \theta) \geq s^M$. ■

Auxiliary scalar functions. For $\lambda > 0$, define the following variables and functions:

$$\Psi_\lambda(x) := \left(\frac{M}{\lambda(\theta_{-, \lambda} + M)}\right) \sqrt{\frac{2(1+x)^2}{(1+x)^2 - 2}} - 1, \quad x > \sqrt{2} - 1 \quad (56a)$$

$$H_\lambda(x) := Q(\Psi_\lambda(x), \theta_{-, \lambda}) - \left(\frac{M\sqrt{R}}{\lambda(\theta_{-, \lambda} + M)}\right) \frac{1+x}{1+\Psi_\lambda(x)} Q(x, \theta_{+, \lambda}) \quad (56b)$$

The role of Ψ_λ is to express the minority-side scalar in terms of the majority-side scalar once the first KKT relation is imposed, while H_λ packages the remaining compatibility condition. In the mixed and fully active phases, solving the regularization path therefore reduces to finding a root of H_λ and then setting $r_\lambda = \Psi_\lambda(s_\lambda)$.

The following lemma determines that the equation governing $\lambda_{K,R}$, as will be derived in the main theorem, has a unique solution.

Lemma 30 *For any $M > 1$ and $R > 1$, there exists a unique $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ such that*

$$Q\left(\frac{2\lambda_{K,R}}{\sqrt{2\lambda_{K,R}^2 - 1}} - 1, \theta_{+, \lambda_{K,R}}\right) = 0.$$

Proof Let $\xi_\lambda := \frac{2\lambda}{\sqrt{2\lambda^2 - 1}} - 1$ and $\theta_\lambda := \theta_{+, \lambda}$. Then, we have

$$\begin{aligned} \frac{d}{d\lambda} Q(\xi_\lambda, \theta_\lambda) &= \frac{\partial Q}{\partial \xi_\lambda} \frac{d\xi_\lambda}{d\lambda} + \frac{\partial Q}{\partial \theta_\lambda} \frac{d\theta_\lambda}{d\lambda} \\ &= \frac{\partial Q}{\partial \xi_\lambda} \frac{-2}{(2x^2 - 1)^{3/2}} + \frac{\partial Q}{\partial \theta_\lambda} \frac{-m\sqrt{R}}{\lambda^2} < 0, \end{aligned} \quad (57)$$

where we have used the fact that $\frac{\partial Q}{\partial \theta_\lambda} > 0$ and $\frac{\partial Q}{\partial \xi_\lambda} > 0$. Therefore, the function $Q(\xi_\lambda, \theta_\lambda)$ is strictly decreasing in λ . Moreover,

$$Q(\xi_1, \theta_1) = Q(1, \theta_1) = \log(1 + 2\theta_1) > 0,$$

and at $\lambda^* = \sqrt{\frac{R+1}{2}}$ the relation $\theta_{\lambda^*} = M \frac{1 - \xi_{\lambda^*}}{1 + \xi_{\lambda^*}}$ would result in

$$Q(\xi_{\lambda^*}, \theta_{\lambda^*}) = \log(\xi_{\lambda^*}^{M-1} M - (M-1)\xi_{\lambda^*}^M) < \log(1) = 0,$$

which implies the existence of a unique root in $[1, \sqrt{\frac{R+1}{2}}]$. ■

The previous lemma identifies the transition point between phases (ii) and (iii). Below this threshold, the mixed mode becomes active, and one must solve a genuinely coupled scalar system. The next two propositions show that this system is well posed in the two remaining nontrivial regimes and record the positivity properties needed later when reconstructing the active coordinates.

Proposition 31 *Let $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ be the point defined in Theorem 30. Then, for any $\lambda \in [1, \lambda_{K,R})$, the function $H_\lambda(\cdot)$ has a root $s_\lambda \in (\sqrt{2} - 1, \frac{2\lambda}{\sqrt{2\lambda^2 - 1}} - 1)$. Furthermore, if $r_\lambda = \Psi_\lambda(s_\lambda)$, then*

$$0 < Q(s_\lambda, \theta_{+, \lambda}), \quad 0 < Q(r_\lambda, \theta_{-, \lambda}), \quad 0 < s_\lambda \exp\left(\frac{-Q(s_\lambda, \theta_{+, \lambda})}{M}\right) < 1.$$

Proof Let $\xi_\lambda := \frac{2\lambda}{\sqrt{2\lambda^2 - 1}} - 1$. We apply the intermediate value theorem by establishing that $H_\lambda(x) \rightarrow +\infty$ as $x \downarrow (\sqrt{2} - 1)$ and $H_\lambda(\xi_\lambda) < 0$.

For the first claim, note that $\Psi_\lambda(x) \rightarrow +\infty$ as $x \downarrow (\sqrt{2} - 1)$, which implies $Q(\Psi_\lambda(x), \theta_{-, \lambda}) \rightarrow +\infty$. Since $\frac{1+x}{1+\Psi_\lambda(x)} \rightarrow 0$ in the same limit, the second term of H_λ vanishes, and hence $H_\lambda(x) \rightarrow +\infty$.

For the second claim, since $\lambda \geq 1$ we have $\theta_{-, \lambda} = 0$, and substituting $\Psi_\lambda(\xi_\lambda) = 1$ into the definition of H_λ yields

$$H_\lambda(\xi_\lambda) = -\frac{\sqrt{R}(1 + \xi_\lambda)}{2\lambda} Q(\xi_\lambda, \theta_{+, \lambda}).$$

Since $\lambda \leq \lambda_{K,R}$, (57) implies $Q(\xi_\lambda, \theta_{+, \lambda}) \geq 0$, with equality only at $\lambda = \lambda_{K,R}$, so $H_\lambda(\xi_\lambda) < 0$ for all $\lambda < \lambda_{K,R}$. The intermediate value theorem then guarantees a root $s_\lambda \in (\sqrt{2} - 1, \xi_\lambda)$ of H_λ .

It remains to verify the sign conditions. Since $\theta_{-, \lambda} = 0$, and Ψ_λ is strictly decreasing with $\Psi_\lambda(\xi_\lambda) = 1$. Because $s_\lambda < \xi_\lambda$, it follows that $r_\lambda = \Psi_\lambda(s_\lambda) > 1$, and therefore

$$Q(r_\lambda, \theta_{-, \lambda}) = Q(r_\lambda, 0) = M \log r_\lambda > 0.$$

From $H_\lambda(s_\lambda) = 0$ with $\theta_{-, \lambda} = 0$ and $\frac{M}{\lambda(\theta_{-, \lambda} + M)} = \frac{1}{\lambda}$, we obtain

$$Q(s_\lambda, \theta_{+, \lambda}) = \frac{\lambda(1 + r_\lambda)}{\sqrt{R}(1 + s_\lambda)} Q(r_\lambda, 0) > 0.$$

Finally, writing $\gamma := Q(s_\lambda, \theta_{+, \lambda})$ and using the definition of Q ,

$$\exp(\gamma) = s_\lambda^{M-1} ((1 + \theta_{+, \lambda})s_\lambda + \theta_{+, \lambda}) > s_\lambda^M,$$

since $\theta_{+, \lambda} > 0$ (as $\lambda < \lambda_{K,R} \leq \sqrt{R}$). Hence $0 < s_\lambda \exp(-\gamma/M) < 1$. ■

Proposition 32 *For any $\lambda \in (0, 1)$, the function $H_\lambda(\cdot)$ has a root $s_\lambda \in (\sqrt{2} - 1, 1)$. Furthermore, if $r_\lambda = \Psi_\lambda(s_\lambda)$, then,*

$$\begin{aligned} 0 < Q(s_\lambda, \theta_{+, \lambda}), \quad 0 < Q(r_\lambda, \theta_{-, \lambda}), \\ 0 < s_\lambda \exp\left(\frac{-Q(s_\lambda, \theta_{+, \lambda})}{M}\right) < 1, \quad 0 < r_\lambda \exp\left(\frac{-Q(r_\lambda, \theta_{-, \lambda})}{M}\right) < 1. \end{aligned}$$

Proof Fix $\lambda \in (0, 1)$. We drop the subscript λ for better readability. We apply the intermediate value theorem by showing that $H(x) \rightarrow +\infty$ as $x \downarrow \sqrt{2} - 1$ and $H(1) < 0$.

Since $\lambda < 1$, we have $\theta_- = M(1/\lambda - 1) > 0$ and $\frac{M}{\lambda(\theta_- + M)} = 1$, so the definitions of Ψ and H reduce to

$$\Psi(x) = \sqrt{\frac{2(1+x)^2}{(1+x)^2 - 2}} - 1, \quad H(x) = Q(\Psi(x), \theta_-) - \sqrt{R} \frac{1+x}{1+\Psi(x)} Q(x, \theta_+).$$

For the first claim, as $x \downarrow \sqrt{2} - 1$ we have $\Psi(x) \rightarrow +\infty$, which implies $Q(\Psi(x), \theta_-) \rightarrow +\infty$. Since $Q(x, \theta_+)$ remains bounded while $\frac{1+x}{1+\Psi(x)} \rightarrow 0$, the second term vanishes and hence $H(x) \rightarrow +\infty$.

For the second claim, since $\Psi(1) = 1$ we have

$$H(1) = Q(1, \theta_-) - \sqrt{R} Q(1, \theta_+) = \log(1 + 2\theta_-) - \sqrt{R} \log(1 + 2\theta_+).$$

Since $R > 1$, we have $\theta_+ = M(\sqrt{R}/\lambda - 1) > M(1/\lambda - 1) = \theta_- > 0$, and therefore

$$\sqrt{R} \log(1 + 2\theta_+) > \log(1 + 2\theta_+) > \log(1 + 2\theta_-),$$

which gives $H(1) < 0$. The intermediate value theorem then guarantees a root $s_\lambda \in (\sqrt{2} - 1, 1)$ of H .

It remains to verify the sign conditions. Since Ψ is strictly decreasing and $s_\lambda < 1$, we have $r_\lambda = \Psi(s_\lambda) > \Psi(1) = 1$. Because $Q(\cdot, \theta_-)$ is strictly increasing and $Q(1, \theta_-) = \log(1 + 2\theta_-) > 0$, it follows that

$$Q(r_\lambda, \theta_-) > Q(1, \theta_-) > 0.$$

From $H(s_\lambda) = 0$ we then obtain

$$Q(s_\lambda, \theta_+) = \frac{(1 + r_\lambda)}{\sqrt{R}(1 + s_\lambda)} Q(r_\lambda, \theta_-) > 0.$$

Finally, writing $\gamma := Q(s_\lambda, \theta_+)$ and using the definition of Q ,

$$\exp(\gamma) = s_\lambda^{M-1}((1 + \theta_+)s_\lambda + \theta_+) = s_\lambda^M + \theta_+ s_\lambda^{M-1}(s_\lambda + 1) > s_\lambda^M,$$

since $\theta_+ > 0$. Combining with $s_\lambda > 0$, we imply that $0 < s_\lambda \exp(-\gamma/M) < 1$. An identical argument with $\delta := Q(r_\lambda, \theta_-)$ and r_λ in place of θ_+ and s_λ gives $0 < r_\lambda \exp(-\delta/M) < 1$. ■

Lemma 33 *For each λ , let z_λ denote the unique zero of $Q(x, \theta_+)$, defined in (58). Then, for any $0 < \lambda < \lambda' \leq \sqrt{R}$, we have $z_\lambda < z_{\lambda'}$.*

Proof Let $\xi_\lambda := \theta_+$. Notice that z_λ is the solution to $Q(z_\lambda, \xi_\lambda) = 0$. Differentiating it with respect to λ gives

$$\frac{\partial Q}{\partial z_\lambda} \frac{dz_\lambda}{d\lambda} + \frac{\partial Q}{\partial \xi_\lambda} \frac{d\xi_\lambda}{d\lambda} = 0$$

which implies that

$$\left(\frac{m-1}{z_\lambda} + \frac{1 + \xi_\lambda}{(1 + \xi_\lambda)z_\lambda + \xi_\lambda} \right) \frac{dz_\lambda}{d\lambda} = - \frac{1}{(1 + \xi_\lambda)z_\lambda + \xi_\lambda} \frac{d\xi_\lambda}{d\lambda}.$$

Since $\xi'_\lambda < 0$ it directly follows that $z'_\lambda > 0$, thereby proving the desired result. ■

With the scalar existence theory in place, we can now return to the vector optimization problem. The theorem below is obtained by imposing the sparsity pattern appropriate to each phase, simplifying the KKT system under that pattern, and then invoking the previous scalar results to certify that the resulting formulas are well defined.

E.3. Analytic Solution of Theorem 3

We first give the deferred definition of the pair (r_λ, s_λ) . Recall from the main body the definitions

$$\theta_+ := M[\sqrt{R}\lambda^{-1} - 1]_+, \quad \theta_- := M[\lambda^{-1} - 1]_+, \quad Q(x, \theta) := \log(x^{M-1}((1+\theta)x + \theta)). \quad (58)$$

For $\lambda \leq \lambda_{K,R}$, let (s_λ, r_λ) denote the unique pair satisfying

$$\begin{cases} 1 + r_\lambda = \frac{M}{\lambda(\theta_- + M)} \sqrt{\frac{2(1+s_\lambda)^2}{(1+s_\lambda)^2 - 2}}, \\ (1 + r_\lambda) \cdot Q(r_\lambda, \theta_-) = \sqrt{R} \left(\frac{M}{\lambda(\theta_- + M)} \right) (1 + s_\lambda) \cdot Q(s_\lambda, \theta_+), \end{cases} \quad (59)$$

and set $\gamma_\lambda := Q(s_\lambda, \theta_+)$ and $\delta_\lambda := Q(r_\lambda, \theta_-)$ for all λ . Notice that the above equations are equivalent to $H_\lambda(s_\lambda) = 0, r_\lambda = \Psi(s_\lambda)$ with the definitions above.

Theorem 34 (Phase-wise analytic solution) *Let $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ be the unique root of*

$$Q\left(\frac{2\lambda}{\sqrt{2\lambda^2 - 1}} - 1, \theta_+\right) \quad (60)$$

Then the unique solution \mathbf{a}_λ of the minimization in Eq. (3) is given as follows:

- (i) For $\lambda \in [\sqrt{R}, \infty)$, $\mathbf{a}_\lambda = [0 \ 0 \ 0 \ 0]^\top$
- (ii) For $\lambda \in [\lambda_{K,R}, \sqrt{R})$, $\mathbf{a}_\lambda = [0 \ 0 \ -M\sqrt{R} \log s_\lambda \ 0]^\top$, where s_λ is the unique root of $Q(x, \theta_{+, \lambda})$.
- (iii) For $\lambda \in [1, \lambda_{K,R})$,

$$\mathbf{a}_\lambda = \left[\frac{\sqrt{R}(\gamma_\lambda - \delta_\lambda)}{\sqrt{2(R+1)}} \quad \frac{R\gamma_\lambda + \delta_\lambda}{\sqrt{2(R+1)}} \quad \sqrt{R}(\gamma_\lambda - M \log s_\lambda) \quad 0 \right]^\top,$$

where for $s_\lambda \in (\sqrt{2} - 1, 1)$ the root of $H_\lambda(x)$ and $r_\lambda := \Psi_\lambda(s_\lambda) \geq 1$, we have set

$$\gamma_\lambda = Q(s_\lambda, \theta_{+, \lambda}), \quad \delta_\lambda = Q(r_\lambda, \theta_{-, \lambda}). \quad (61)$$

- (iv) For $\lambda \in (0, 1)$,

$$\mathbf{a}_\lambda = \left[\frac{\sqrt{R}(\gamma_\lambda - \delta_\lambda)}{\sqrt{2(R+1)}} \quad \frac{R\gamma_\lambda + \delta_\lambda}{\sqrt{2(R+1)}} \quad \sqrt{R}(\gamma_\lambda - M \log s_\lambda) \quad \delta_\lambda - M \log r_\lambda \right]^\top,$$

where for $s_\lambda \in (\sqrt{2} - 1, 1)$ the root of $H_\lambda(x)$, and $r_\lambda := \Psi(s_\lambda)$, we have set

$$\gamma_\lambda = Q(s_\lambda, \theta_{+, \lambda}), \quad \delta_\lambda = Q(r_\lambda, \theta_{-, \lambda}). \quad (62)$$

Proof For clarity, we suppress the subscript λ throughout the proof. The existence of a unique $\lambda_{K,R} \in [1, \sqrt{R}]$ is guaranteed by Theorem 30.

We treat each phase separately, deriving the equations that arise from the KKT conditions under the appropriate sparsity pattern. In each case the resulting system admits a unique solution; the uniqueness of the overall minimizer from Theorem 22 then guarantees that the proposed formulas hold in every phase.

- (i) This case would follow trivially from the sparsity features proved in Theorem 26.
- (ii) Based on Theorem 26, any solution would have to satisfy $\mathbf{a}_+ \neq \mathbf{0}$, $\mathbf{a}_- = \mathbf{0}$, $\gamma = \delta = 0$. Assuming all entries in \mathbf{a}_+ are nonzero, we have $\mathbf{b}_+ = \lambda \mathbf{1}$, after which Theorem 29 (with $t = 1 - \frac{\lambda}{\sqrt{R}}$) would result in $\mathbf{a}_+ = -M\sqrt{R} \log s_\lambda$, where $Q(s_\lambda, \theta_+) = 0$ and $s_\lambda \in (0, 1)$

At this point, we only need to check the complementary KKT condition for $(\alpha, a) = (0, 0)$, which is $\beta^2 + b^2 \leq \lambda^2$. Substituting for this would lead to

$$X_+ = \frac{2M}{M + s^{1-M} + (M-1)s} = \frac{2\lambda}{\sqrt{R}(1+s)}, \quad X_- = 1$$

$$\frac{\beta^2 + b^2}{\lambda^2} = \frac{1}{2\lambda^2}(RX_+^2 + X_-^2) = \frac{2}{(1+s)^2} + \frac{1}{2\lambda^2}.$$

Moreover, using Theorem 33, we have $z \geq \frac{2\lambda_{K,R}}{\sqrt{2\lambda_{K,R}^2 - 1}} - 1$, which implies that $\beta^2 + b^2 \leq \lambda^2$.

Thus, the unique solution in this phase has been found, validating the assumption of no zero entries in \mathbf{a}_+ .

- (iii) This is the first genuinely coupled regime: The majority coordinate and the mixed (α, a) -block are active, while the minority coordinate is still zero.

First, note that given $1 \leq \lambda \leq \sqrt{R}$, Theorem 26 would require $\mathbf{a}_+ \neq \mathbf{0}$ and $\mathbf{a}_- = \mathbf{0}$ for any solution. Assuming all entries in \mathbf{a}_+ are non-zero would again give $\mathbf{b}_+ = \lambda \mathbf{1}$, after which Theorem 29 would solve yield $\mathbf{a}_+ = \sqrt{R}(\gamma - M \log s)\mathbf{1} \geq \mathbf{0}$, where $\gamma = Q(s, \theta_+)$ and s is unique for a given γ .

The KKT conditions for (α, a) in this phase are $(\beta, b) = \frac{\lambda}{\sqrt{a^2 + \alpha^2}}(\alpha, a)$. Given that Theorem 26 guarantees $\gamma, \delta \geq 0$, they can be reformulated equivalently as

$$\lambda^2 = \beta^2 + b^2 = \frac{1}{2}(RX_+^2 + X_-^2), \quad \frac{\gamma}{X_+} = \frac{\delta}{X_-}. \quad (63)$$

Defining $r := \exp(\delta/M)$, the above conditions would simplify to

$$\frac{1}{2} = \frac{1}{(1+s)^2} + \frac{1}{\lambda^2(1+r)^2}, \quad \lambda\delta(1+r) = \sqrt{R}\gamma(1+s)$$

The existence of $s \in (\sqrt{2} - 1, \xi_\lambda)$ and the fact that all terms in (61) are well-defined follow from Theorem 31. Consequently, the above system has a unique solution, which via uniqueness of the solution to the original problem, verifies the assumption of all-nonzero entries in \mathbf{a}_+ .

Finally, we should note that by utilizing $\theta_- = 0$ in this phase, the above conditions can be written in the form of (59).

- (iv) The final phase follows the same template as phase (iii), with the minority coordinate becoming active its only new feature.

Similar to the previous phase, Theorem 26 requires $\mathbf{a}_\pm \neq \mathbf{0}$. By assuming \mathbf{a}_+ and \mathbf{a}_- have strictly positive entries, the KKT conditions give $\mathbf{b}_\pm = \lambda \mathbf{1}$, from which Theorem 29 would give

$$\mathbf{a}_+ = \sqrt{R}(\gamma - M \log s)\mathbf{1}, \quad \mathbf{a}_- = (\delta - M \log r)\mathbf{1},$$

$$\gamma = Q(s, \theta_+), \quad \delta = Q(r, \theta_-). \quad (64)$$

In addition, the reformulated KKT conditions for (α, a) of (63) would simplify to the slightly different form

$$\frac{1}{2} = \frac{1}{(1+s)^2} + \frac{1}{(1+r)^2}, \quad \delta(1+r) = \sqrt{R}\gamma(1+s) \quad (65)$$

The above system of four equations can be reformulated as a scalar equation in s , where Theorem 32 guarantees the existence of a unique $s \in (\sqrt{2} - 1, 1)$ and well-defined other quantities.

Finally, by using the fact that $\theta_- = 0$ for $\lambda \geq 1$, one can reach a unified form of equations for the third and fourth phase as in (59). \blacksquare

Next, we will discuss more high-level characteristics of the regularization path solution given the analytic parametrization for different phases.

Proposition 35 (More characteristics of the regularization path solution) *For all values of λ , the regularization path solution parameters satisfy the following properties:*

$$\forall \lambda \geq 0: \quad a_{+,\lambda} \geq \sqrt{R}a_{-,\lambda}, \quad \alpha_\lambda \leq 0 \quad (66)$$

Proof We first show that $\sqrt{R}a_- \leq a_+$. This is immediate for $\lambda \geq 1$, since $a_{\lambda,-} = 0$ and $a_{+,\lambda} \geq 0$. It remains to verify the claim for $\lambda \in (0, 1)$. In this regime, we have

$$\begin{aligned} a_+ &= \sqrt{R}(\gamma - M \log s), \quad s^{M-1}((1+\theta_+)s + \theta_+) = e^\gamma \Rightarrow a_+ = \sqrt{R} \log(1 + \theta_+ + \frac{\theta_+}{s}) \\ a_- &= (\delta - M \log r), \quad r^{M-1}((1+\theta_+)r + \theta_-) = e^\delta \Rightarrow a_- = \log(1 + \theta_- + \frac{\theta_-}{r}) \end{aligned} \quad (67)$$

We know that $s \leq r$ and $\theta_+ = \sqrt{R}\theta_- + M(\sqrt{R} - 1) \geq \theta_-$. For $\lambda \leq 1$, we therefore conclude that

$$\log(1 + \theta_- + \frac{\theta_-}{r}) \leq \log(1 + \theta_+ + \frac{\theta_+}{s}) \Rightarrow a_- \leq \frac{a_+}{\sqrt{R}}.$$

This establishes the first claim for all $\lambda > 0$. For all values of λ , either $(\alpha, a) = (0, 0)$ or at least one of them is nonzero, which gives

$$\begin{aligned} \beta &= \frac{\lambda\alpha}{\sqrt{\alpha^2 + a^2}} \Rightarrow \alpha \propto [\mathbf{1}_M^T \exp(\mathbf{x}_-) - \mathbf{1}_M^T \exp(\mathbf{x}_+)] \\ &\propto e^{-\frac{\alpha\sqrt{R}}{M}} \left(e^{-\frac{a_-}{M\sqrt{R}}} (e^{\frac{a_-}{\sqrt{R}}} + M - 1) \right) - e^{\frac{\alpha}{M\sqrt{R}}} \left(e^{-\frac{a_+}{M\sqrt{R}}} (e^{\frac{a_+}{\sqrt{R}}} + M - 1) \right), \end{aligned}$$

where $\mathbf{a}_\pm = a_\pm \mathbf{1}$ was used to simplify the above equation. Now the claim can be proved by contradiction. Suppose that $\alpha > 0$. Since the expression $a^{M-1} + (M-1)a$ is increasing in a and $a_- \leq \frac{a_+}{\sqrt{R}}$, the above expression would become negative, contradicting the prior assumption. Hence, $\alpha \leq 0$, with strict inequality in its active phase. \blacksquare

E.4. Behavior of the threshold $\lambda_{K,R}$

Proposition 36 Consider the threshold $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ between phases (iii) and (iv) in Theorem 34, with its value defined by (60). Then, it has the following properties:

1. $\lambda_{K,R}$ is increasing as a function of M for fixed R , and $\lim_{M \rightarrow \infty} \lambda_{K,R} = 1$
2. For fixed M , $\lambda_{K,R}$ is increasing in R , but $\frac{\lambda_{K,R}}{\sqrt{R}}$ decreases and $\lim_{R \rightarrow \infty} \frac{\lambda_{K,R}}{\sqrt{R}} = c_K$, where c_K solely depends on K .

Proof The threshold is the value of λ where phases (ii) and (iii) intersect; as a result, $\gamma = \delta = 0$ and at the same time $\beta^2 + b^2 = \lambda^2$. Using these to modify the phase (iii) equations leads us to

$$1 = \frac{2}{(1+s)^2} + \frac{1}{2\lambda_{K,R}^2}, \quad Q(s, \theta_+) = \log(s^{M-1}((1+\theta_+)s + \theta_+)) = 0, \quad \theta_+ = M[\sqrt{R}\lambda_{K,R}^{-1} - 1] \quad (68)$$

From this point, we drop the subscript from $\lambda_{M,R}$ to reduce clutter. Solving for λ in terms of s yields

$$\lambda(s) = \frac{1+s}{\sqrt{2(s^2 + 2s - 1)}}, \quad \lambda'(s) = -\frac{2}{\sqrt{2}(s^2 + 2s - 1)^{3/2}} < 0, \quad \frac{\partial \theta_+}{\partial s} = -\frac{M\sqrt{R}}{\lambda^2} \lambda'(s) > 0$$

Assuming M is continuous, implicitly differentiating $Q(s, \theta_+) = 0$ with respect to it and simplifying would result in

$$Q'(M) = \partial_M Q + (\partial_s Q) s'(M) = 0, \\ \partial_M Q = \frac{1 - s^M + \log(s^M)}{M} \leq 0, \quad \partial_s Q = \frac{M-1}{s} + \frac{(1+\theta_+) + (1+s) \partial_s \theta_+}{(1+\theta_+)s + \theta_+} \geq 0,$$

from which we can conclude $s'(M) \geq 0$ and thus $\lambda'(M) \leq 0$.

The same differentiation method can be applied to \sqrt{R} to reach

$$Q'(\sqrt{R}) = \partial_{\sqrt{R}} Q + (\partial_s Q) s'(\sqrt{R}) = 0, \quad \partial_s Q = \frac{(M + \theta_+)(1+s)}{\sqrt{R}((1+\theta_+)s + \theta_+)} \geq 0,$$

from which follows that s is decreasing in \sqrt{R} (and thus R) and λ is increasing in R . However, to determine the behavior of the ratio to λ_+ , we reformulate $Q(s, \theta_+) = 0$ as

$$\frac{\sqrt{R}}{\lambda} = \frac{s^{-(M-1)} - s}{M(1+s)} + 1,$$

where the RHS is decreasing in s . This will result in $\frac{\lambda_{K,R}}{\sqrt{R}}$ to be decreasing in R .

To derive the asymptotics, first note that by using $s \geq \sqrt{2} - 1$, the following can be established:

$$s^M(1+\theta_+) \leq s^{M-1}((1+\theta_+)s + \theta_+) \leq s^M(1+\theta_+)(1 + \frac{1}{s}) \leq \frac{\sqrt{2}}{\sqrt{2}-1} s^M(1+\theta_+), \\ (\frac{\sqrt{2}}{\sqrt{2}-1})^{1/M} \frac{1}{(1+\theta_+)^{1/M}} \leq s \leq \frac{1}{(1+\theta_+)^{1/M}}.$$

■

Since $\lambda_{K,R} \in [1, \sqrt{\frac{R+1}{2}}]$ stays bounded for all M , $\theta_+ = \mathcal{O}(M)$, and the limit of upper and lower bounds are 1. Consequently, $\lim_{M \rightarrow \infty} \lambda_{K,R} = 1$. To derive the limit for $R \rightarrow \infty$, we substitute for θ_+ in terms of s to reach

$$(s+1)^2 - 2 = \frac{1}{2R} \left(\frac{s^{-(M-1)} - s}{M} + 1 + s \right)^{-2} \Rightarrow \lim_{R \rightarrow \infty} (s+1)^2 - 2 = 0 \Rightarrow \lim_{R \rightarrow \infty} s = \sqrt{2} - 1,$$

where the boundedness of s for all R was used to calculate the LHS limit. Substituting this into $Q(s, \theta_+) = 0$ would yield

$$\lim_{R \rightarrow \infty} \frac{\lambda}{\sqrt{R}} = \left(1 + \frac{(\sqrt{2}-1)^{-M+1} - (\sqrt{2}-1)}{M\sqrt{2}} \right)^{-1} = c_K$$

E.5. Asymptotic behavior of RP solution and Proof of Theorem 4

Before analyzing the asymptotic behavior of the regularization path solution, it is useful to restate and refine the upper and lower bounds for s, r for $\lambda \in (0, 1]$ the following lemma:

Lemma 37 *For the regularization path characterized (64) and (65), the following bounds can be established:*

$$\sqrt{\frac{2(R+1)}{R}} - 1 \leq s \leq 1, \quad 1 \leq r \leq \sqrt{2(R+1)} - 1 \quad (69)$$

Proof Recall that $s \in (\sqrt{2}-1, 1)$ and $r \in (1, \infty)$ has already been established in Theorem 32. To reach the refined bounds in the lemma, notice that $\alpha \leq 0$ from Theorem 35 results in $\delta \geq \gamma$. Substituting this bound in (65) and simplifying proves the desired bounds. ■

Proposition 38 (Asymptotic direction as $\lambda \rightarrow 0^+$) *For each $\lambda \in (0, 1)$, consider the solution of the KKT system in (64) and (65) of the form*

$$\mathbf{a}_+ = \sqrt{R}(\gamma - M \log s) \mathbf{1}_{M-1}, \quad \mathbf{a}_- = (\delta - \log r) \mathbf{1}_{M-1}, \quad a = \frac{R\gamma + \delta}{\sqrt{2(R+1)}}, \quad \alpha = \frac{\sqrt{R}(\gamma - \delta)}{\sqrt{2(R+1)}},$$

where all parameters are nonzero. Define $\mathbf{a}_\infty := (0, \sqrt{\frac{R+1}{2}}, \sqrt{R}\mathbf{1}_{M-1}, \mathbf{1}_{M-1})$. Also define the norm operator $\|\mathbf{a}\|_{+1} := \sqrt{a^2 + \alpha^2} + \|\mathbf{a}_+\|_1 + \|\mathbf{a}_-\|_1$. The the following limit holds:

$$\lim_{\lambda \rightarrow 0^+} \frac{\mathbf{a}_\lambda}{\|\mathbf{a}_\lambda\|_{+1}} = \frac{\mathbf{a}_\infty}{\|\mathbf{a}_\infty\|_{+1}}$$

Proof We bound the quantities using (65), along with the bounds established in Theorem 37.

Using the fact that $a_0 = \sqrt{R}(\gamma - M \log(s))$ and $a_1 = \delta - M \log(r)$, the following bounds hold:

$$\begin{aligned} s^M(1 + \theta_+) &\leq s^{M-1}((1 + \theta_+)s + \theta_+) = e^\gamma \leq s^M(1 + \theta_+)\left(1 + \frac{1}{s}\right) \\ \Rightarrow \log(1 + \theta_+) &\leq \gamma - M \log(s) \leq \log(1 + \theta_+) - \log\left(1 - \sqrt{\frac{R}{2(R+1)}}\right), \\ \log(1 + \theta_-) &\leq \delta - M \log(r) \leq \log(1 + \theta_-) + \log(2), \end{aligned}$$

Since s, r are bounded from above and below by Theorem 37,

$$\begin{aligned}
 \gamma &= \log(1 + \theta_+) + c_\gamma, & M \log\left(\sqrt{\frac{2(R+1)}{R}} - 1\right) &\leq c_\gamma \leq -\log\left(1 - \sqrt{\frac{R}{2(R+1)}}\right) \\
 \delta &= \log(1 + \theta_-) + c_\delta, & 0 &\leq c_\delta \leq \log(2) + \log(\sqrt{2(R+1)} - 1) \\
 \frac{a_+}{\sqrt{R}} &= \log(1 + \theta_+) + c_{a_+}, & 0 &\leq c_{a_+} \leq -\log\left(1 - \sqrt{\frac{R}{2(R+1)}}\right) \\
 \frac{a_-}{\sqrt{R}} &= \log(1 + \theta_-) + c_{a_-}, & 0 &\leq c_{a_-} \leq \log(2)
 \end{aligned} \tag{70}$$

Since $\log(1 + \theta_\pm)$ scale as $\mathcal{O}(\log(\lambda^{-1}))$, it is evident from the above bounds in (70) that all quantities grow at this rate, which establishes the desired limit. Furthermore, it follows that

$$\lim_{\lambda \rightarrow 0^+} \frac{\mathbf{a}_\lambda}{\|\mathbf{a}_\lambda\|_{+1}} = \frac{\mathbf{a}_\infty}{\|\mathbf{a}_\infty\|_{+1}} + \mathcal{O}\left(\frac{1}{\log \lambda^{-1}}\right)$$

■

Beyond the above formal characterization of the limit of the regularization path, we further conjecture that the singular value associated with the majority class converges to its limit from above; in other words,

$$\frac{a_{+,\lambda}}{\|\mathbf{a}_\lambda\|_{+1}} - \frac{a_{+,\infty}}{\|\mathbf{a}_\infty\|_{+1}} \geq 0, \quad \lim_{\lambda \rightarrow 0^+} \frac{a_{+,\lambda}}{\|\mathbf{a}_\lambda\|_{+1}} - \frac{a_{+,\infty}}{\|\mathbf{a}_\infty\|_{+1}} = 0.$$

Proof [Supporting Argument] The limit statement follows directly from Theorem 38. For the inequality claim, we use the following inequality, observed consistently in empirical simulations of the regularization path solution:

$$\frac{a_+}{\sqrt{R}} \geq \sqrt{\frac{2}{R+1}} \sqrt{a^2 + \alpha^2}$$

Next, we simplify the expression as

$$\frac{a_{+,\lambda}}{\|\mathbf{a}_\lambda\|_{+1}} - \frac{a_{+,\infty}}{\|\mathbf{a}_\infty\|_{+1}} = \frac{(M-1)(a_{+,\lambda} - \sqrt{R}a_{-,\lambda}) + (\sqrt{\frac{R+1}{2}}a_{+,\lambda} - \sqrt{a_\lambda^2 + \alpha_\lambda^2})}{\|\mathbf{a}_\lambda\|_{+1}\|\mathbf{a}_\infty\|_{+1}}.$$

The denominator is positive. The first term in the numerator is nonnegative as a result of Theorem 35, and the second term is nonnegative as a result of the empirical inequality. ■

E.6. Alternative proof of regularization-path limit

For completeness, we verify here that the max-margin problem that is associated with the regularization path limit of Eq. (3) has SEL as its unique solution. Note that Theorem 4 proves that directly by taking an explicit limit of the regularization-path solution. Alternatively, we can show through proof by contradiction, which is a now classical argument due to [40, 42] and used extensively in the

recent literature [28, 30, 51, e.g.], that the solution of (3) directionally converges in the limit $\lambda \rightarrow 0$ to the (normalized) max-margin solution below:

$$\begin{aligned} \min_{\mathbf{a}_+, \mathbf{a}_-, a, \alpha} \quad & \|\mathbf{a}_+\|_1 + \|\mathbf{a}_-\|_1 + \sqrt{a^2 + \alpha^2} \\ \text{s.t.} \quad & \forall i \in [M-1]: \frac{1}{\sqrt{R}}(\mathbf{1} - \bar{\phi}_{i+1})^T \mathbf{a}_+ \geq 1, \quad \frac{1}{\sqrt{R}}\mathbf{1}^T \mathbf{a}_+ + \sqrt{\frac{2}{R+1}} \left(a + \frac{\alpha}{\sqrt{R}} \right) \geq 1 \quad (71) \\ & \forall i \in [M-1]: (\mathbf{1} - \bar{\phi}_{i+1})^T \mathbf{a}_- \geq 1, \quad \mathbf{1}^T \mathbf{a}_- + \sqrt{\frac{2}{R+1}} \left(a - \sqrt{R}\alpha \right) \geq 1 \end{aligned}$$

We omit the details of this standard argument because they bring no new insights. Instead, we verify here that the unique solution of this max-margin problem matches the SEL geometry of Theorem 4. This also yields an alternative proof of [52, Thm. 1, Prop. 2]. Note that compared to our approach in Theorem 4 the overall proof here is not as explicit (cf. proof by contradiction) and does not yield rates of convergence.

Theorem 39 *The optimization problem of (71) admits the unique solution of*

$$\mathbf{a}_+ = \frac{1}{M}\sqrt{R}\mathbf{1}, \quad \mathbf{a}_- = \frac{1}{M}\mathbf{1}, \quad a = \frac{1}{M}\sqrt{\frac{R+1}{2}}, \quad \alpha = 0 \quad (72)$$

Proof Define the rescaled variables $\hat{\mathbf{a}}_+ = \frac{1}{\sqrt{R}}\mathbf{a}_+$, $\hat{\mathbf{a}}_- = \mathbf{a}_-$, $(\hat{a}, \hat{\alpha}) = \sqrt{\frac{2}{R+1}}(a, \alpha)$. The problem can be written as

$$\begin{aligned} \min \quad & \sqrt{R}\|\hat{\mathbf{a}}_+\|_1 + \|\hat{\mathbf{a}}_-\|_1 + \sqrt{\frac{R+1}{2}}\sqrt{\hat{a}^2 + \hat{\alpha}^2} \\ \text{s.t.} \quad & \Psi\hat{\mathbf{a}}_+ \geq \mathbf{1}, \quad \Psi\hat{\mathbf{a}}_- \geq \mathbf{1}, \quad \mathbf{1}^T \hat{\mathbf{a}}_+ + \left(\hat{a} + \frac{\hat{\alpha}}{\sqrt{R}}\right) \geq 1, \quad \mathbf{1}^T \hat{\mathbf{a}}_- + (\hat{a} - \sqrt{R}\hat{\alpha}) \geq 1. \end{aligned}$$

The KKT conditions in terms of the dual variables are

$$\sqrt{R}\text{sign}(\hat{\mathbf{a}}_+) = \Psi\boldsymbol{\nu}_+ + \lambda_+\mathbf{1}, \quad \text{sign}(\hat{\mathbf{a}}_-) = \Psi\boldsymbol{\nu}_- + \lambda_-\mathbf{1}, \quad (73a)$$

$$\sqrt{\frac{2}{R+1}}(\lambda_+ + \lambda_-, \frac{\lambda_+}{\sqrt{R}} - \lambda_-\sqrt{R}) \in \partial_{(\hat{a}, \hat{\alpha})}(\sqrt{\hat{a}^2 + \hat{\alpha}^2}) \quad (73b)$$

$$\boldsymbol{\nu}_\pm \geq \mathbf{0}, \quad \lambda_\pm \geq 0, \quad \boldsymbol{\nu}_\pm \odot (\Psi\hat{\mathbf{a}}_\pm - \mathbf{1}) = \mathbf{0}, \quad (73c)$$

$$\lambda_+(\mathbf{1}^T \hat{\mathbf{a}}_+ + (\hat{a} + \frac{\hat{\alpha}}{\sqrt{R}}) - 1) = 0, \quad \lambda_-(\mathbf{1}^T \hat{\mathbf{a}}_- + (\hat{a} - \hat{\alpha}\sqrt{R}) - 1) = 0, \quad (73d)$$

where the sign function is treated in an extended value form. Notice that since the entries of Ψ are nonnegative, (73a) alongside the dual constraints in (73c) result in nonnegative signs for \mathbf{a}_+ , \mathbf{a}_- and thus $\mathbf{a}_+ \geq \mathbf{0}$, $\mathbf{a}_- \geq \mathbf{0}$. From this point on, our first step is to show the only solution satisfying $(\hat{a}, \hat{\alpha}) \neq (0, 0)$ and $\mathbf{a}_+ > \mathbf{0}$, $\mathbf{a}_- > \mathbf{0}$ is the solution in the theorem statement. Assuming that these conditions would result in the simplified KKT equations

$$\boldsymbol{\nu}_+ = \Psi^{-1}(\sqrt{R} - \lambda_+)\mathbf{1} = \frac{1}{M}(\sqrt{R} - \lambda_+)\mathbf{1}, \quad \boldsymbol{\nu}_- = \frac{1}{M}(1 - \lambda_-)\mathbf{1} \quad (74a)$$

$$\sqrt{\frac{2}{R+1}}(\lambda_+ + \lambda_-, \frac{\lambda_+}{\sqrt{R}} - \lambda_-\sqrt{R}) = \frac{1}{\sqrt{\hat{a}^2 + \hat{\alpha}^2}}(\hat{a}, \hat{\alpha}). \quad (74b)$$

Notice that (74b) implies that $\frac{\lambda_+^2}{R} + \lambda_-^2 = \frac{1}{2}$. This constraint eliminates the possibility for $\lambda_+ = \sqrt{R}$ or $\lambda_- = 1$, from which follows that $\nu_+ > \mathbf{0}, \nu_- > \mathbf{0}$. As a result of complementary slackness, the constraints $\Psi \hat{\mathbf{a}}_+ \geq \mathbf{1}$ and $\Psi \hat{\mathbf{a}}_- \geq \mathbf{1}$ are active, from which $\hat{\mathbf{a}}_+, \hat{\mathbf{a}}_-$ can be recovered as

$$\Psi \hat{\mathbf{a}}_+ = \mathbf{1}, \quad \Psi \hat{\mathbf{a}}_- = \mathbf{1} \Rightarrow \hat{\mathbf{a}}_+ = \frac{1}{M} \mathbf{1}, \quad \hat{\mathbf{a}}_- = \frac{1}{M} \mathbf{1}.$$

As a result, the remaining KKT conditions on $(\hat{a}, \hat{\alpha})$ can be simplified to

$$\hat{a} + \frac{\hat{\alpha}}{\sqrt{R}} \geq \frac{1}{M}, \quad \hat{a} - \hat{\alpha} \sqrt{R} \geq \frac{1}{M}, \quad \sqrt{\frac{2}{R+1}} (\lambda_+ + \lambda_-, \frac{\lambda_+}{\sqrt{R}} - \lambda_- \sqrt{R}) = \frac{1}{\sqrt{\hat{a}^2 + \hat{\alpha}^2}} (\hat{a}, \hat{\alpha}).$$

Standard case-by case checking will show that setting either of $\lambda_+ = 0$ or $\lambda_- = 0$ will lead to violation of one of the above simplified constraints; consequently, $\lambda_+, \lambda_- \neq 0$, which leads to the scalar constraints being active. Solving for $\hat{a}, \hat{\alpha}$ will give the solution in the theorem statement.

So far, (72) has been shown to be a valid KKT point for the optimization problem. Our final step is proving that it is the unique optimal solution to (71). Note that the strict convexity of the 2-norm over $(\hat{a}, \hat{\alpha})$ guarantees their uniqueness. Using this fact, the primal constraints can be simplified to

$$\Psi(\hat{\mathbf{a}}_+ - \frac{1}{M} \mathbf{1}) \geq \mathbf{0}, \quad \Psi(\hat{\mathbf{a}}_- - \frac{1}{M} \mathbf{1}) \geq \mathbf{0} \quad \mathbf{1}^T(\hat{\mathbf{a}}_+ - \frac{1}{M} \mathbf{1}) \geq 0, \quad \mathbf{1}^T(\hat{\mathbf{a}}_- - \frac{1}{M} \mathbf{1}) \geq 0.$$

The scalar constraints above hold for (72) and the terms $\mathbf{1}^T \hat{\mathbf{a}}_+, \mathbf{1}^T \hat{\mathbf{a}}_-$ in the objective function as well. As a result, the scalar constraints are active for every other solution. Use the fact that $\mathbf{1}^T \Psi = M \mathbf{1}^T$ to conclude that while every element of the vector $\Psi(\hat{\mathbf{a}}_+ - \frac{1}{M} \mathbf{1})$ is nonnegative, but the sum of its elements is zero. Hence, $\Psi(\hat{\mathbf{a}}_+ - \frac{1}{M} \mathbf{1}) = \mathbf{0}$ and likewise for $\hat{\mathbf{a}}_-$, which concludes the uniqueness argument. ■

Appendix F. Proofs for Gradient-flow path

F.1. Proof of Theorem 5

To reduce clutter in our proof, we drop the subscript t in our expressions. For the moment, we assume a general decomposition for \mathbf{W}, \mathbf{H} . Suppose $\mathbf{W} = \mathbf{U} \widetilde{\mathbf{W}} \mathbf{R}^T$ and $\mathbf{H} = \mathbf{R} \begin{bmatrix} \alpha_H \mathbf{e}_+ & \widetilde{\mathbf{H}} \end{bmatrix} \widetilde{\mathbf{V}}^T$. Then we have $\mathbf{W} \mathbf{H} = \mathbf{U} \begin{bmatrix} \widetilde{\mathbf{W}} \alpha_H \mathbf{e}_+ & \widetilde{\mathbf{W}} \widetilde{\mathbf{H}} \end{bmatrix} \widetilde{\mathbf{V}}^T$ and

$$\begin{aligned} \dot{\mathbf{W}} &= \mathbf{U} \begin{bmatrix} \beta \mathbf{e}_0 & \mathbf{B} \end{bmatrix} \underbrace{\widetilde{\mathbf{V}}^T \widetilde{\mathbf{V}}}_{\mathbf{I}} \begin{bmatrix} \alpha_H \mathbf{e}_0 & \widetilde{\mathbf{H}} \end{bmatrix}^T \mathbf{R}^T = \mathbf{U} (\beta \alpha_H \mathbf{e}_0 \mathbf{e}_0^T + \mathbf{B} \widetilde{\mathbf{H}}) \mathbf{R}^T \\ \dot{\mathbf{H}} &= \mathbf{R} \widetilde{\mathbf{W}} \begin{bmatrix} \beta \mathbf{e}_0 & \mathbf{B} \end{bmatrix} \widetilde{\mathbf{V}}^T = \mathbf{R} \begin{bmatrix} \beta (\widetilde{\mathbf{W}})_0 \mathbf{e}_0 & \widetilde{\mathbf{W}} \mathbf{B} \end{bmatrix} \widetilde{\mathbf{V}}^T. \end{aligned}$$

So, the subspaces and structure do not change over time, and $\widetilde{\mathbf{W}}, \widetilde{\mathbf{H}}$ can be tracked directly:

$$\frac{d\widetilde{\mathbf{W}}}{dt} = \mathbf{B} \widetilde{\mathbf{W}} + \beta \alpha_H \mathbf{e}_0 \mathbf{e}_0^T, \quad \frac{d\widetilde{\mathbf{H}}}{dt} = \widetilde{\mathbf{W}} \mathbf{B}, \quad \alpha_H = \beta \widetilde{\mathbf{W}} [0, 0]$$

In addition, the gradient flow equations result in

$$\begin{aligned}\frac{d\mathbf{W}^T\mathbf{W}}{dt} &= \mathbf{W}^T(\mathbf{Y} - \mathbb{S}(\mathbf{W}\mathbf{H}))\mathbf{H} + \mathbf{H}(\mathbf{Y} - \mathbb{S}(\mathbf{W}\mathbf{H}))^T\mathbf{W} \\ \frac{d\mathbf{H}\mathbf{H}^T}{dt} &= \mathbf{W}^T(\mathbf{Y} - \mathbb{S}(\mathbf{W}\mathbf{H}))\mathbf{H} + \mathbf{H}(\mathbf{Y} - \mathbb{S}(\mathbf{W}\mathbf{H}))^T\mathbf{W},\end{aligned}$$

which directly leads to the conserved quantity equation $\frac{d}{dt}(\mathbf{W}^T\mathbf{W} - \mathbf{H}_t\mathbf{H}_t^T) = \mathbf{0}$. Combining this with the factorizations of \mathbf{W} , \mathbf{H} , we have

$$\frac{d}{dt}(\mathbf{R}(\widetilde{\mathbf{W}}^T\widetilde{\mathbf{W}} - \alpha_H^2\mathbf{e}_0\mathbf{e}_0^T - \widetilde{\mathbf{H}}\widetilde{\mathbf{H}}^T)\mathbf{R}^T) = \mathbf{0} \Rightarrow \frac{d}{dt}(\widetilde{\mathbf{W}}^T\widetilde{\mathbf{W}} - \alpha_H^2\mathbf{e}_0\mathbf{e}_0^T - \widetilde{\mathbf{H}}\widetilde{\mathbf{H}}^T) = \mathbf{0}.$$

In the diagonal case, we denote the diagonal matrices $\overline{\mathbf{W}}$, $\overline{\mathbf{H}}$ as

$$\widetilde{\mathbf{W}} = \text{diag}(w, \mathbf{w}_+ \otimes \mathbf{f}_+ + \mathbf{w}_- \otimes \mathbf{f}_-), \quad \widetilde{\mathbf{H}} = \text{diag}(\tilde{h}, \tilde{\mathbf{h}}_+ \otimes \mathbf{f}_+ + \tilde{\mathbf{h}}_- \otimes \mathbf{f}_-),$$

As a result, the dynamics in the previous section can be simplified to

$$\dot{\mathbf{w}}_{\pm} = \mathbf{b}_{\pm} \odot \mathbf{h}_{\pm}, \quad \dot{\mathbf{h}}_{\pm} = \mathbf{b}_{\pm} \odot \mathbf{w}_{\pm}, \quad \dot{w} = \beta\alpha_H + bh, \quad \dot{h} = bw, \quad \dot{\alpha}_H = \beta\tilde{w} \quad (75)$$

In addition, the logits and the equations governing them would be

$$\mathbf{a}_{\pm} := \mathbf{w}_{\pm} \odot \mathbf{h}_{\pm}, \quad \dot{\mathbf{a}}_{\pm} = \mathbf{b}_{\pm} \odot (\mathbf{w}_{\pm}^2 + \mathbf{h}_{\pm}^2), \quad (76)$$

$$a := wh, \quad \alpha := \alpha_H w, \quad \dot{a} = b(w^2 + h^2) + \beta\alpha_H h, \quad \dot{\alpha} = \beta(w^2 + \alpha_H^2) + bh\alpha_H. \quad (77)$$

The matrix conserved quantity equations would also simplify to

$$\frac{d}{dt}(w^2 - \alpha_H^2 - h^2) = 0, \quad \frac{d}{dt}(\mathbf{w}_{\pm}^2 - \mathbf{h}_{\pm}^2) = \mathbf{0}.$$

Given the conserved quantities, assuming a balanced initialization would simplify the above equations to

$$\mathbf{a}_{\pm} := \mathbf{w}_{\pm}^2 \geq \mathbf{0}, \quad \dot{\mathbf{a}}_{\pm} = 2\mathbf{a}_{\pm} \odot \mathbf{b}_{\pm}, \quad w^2 = \sqrt{a^2 + \alpha^2}, \quad (78)$$

$$\dot{a} = b(w^2 + \frac{a^2}{w^2}) + \beta\frac{\alpha a}{w^2} = \frac{b(2a^2 + \alpha^2) + \beta a \alpha}{\sqrt{a^2 + \alpha^2}}, \quad (79)$$

$$\dot{\alpha} = \beta(w^2 + \frac{\alpha^2}{w^2}) + b\frac{a\alpha}{w^2} = \frac{\beta(2\alpha^2 + a^2) + ba\alpha}{\sqrt{a^2 + \alpha^2}}. \quad (80)$$

F.2. Proof of Theorem 9: Majority and Minority tend to Uniformity

Theorem 40 (Restatement of Theorem 9) *Let $\mathbf{a}_{\pm,t} \in \mathbb{R}^{M-1}$ be variables initialized with $\mathbf{a}_{\pm,0} > \mathbf{0}$, and evolve under (8), for arbitrary a_0, α_0 initialization. For every vector \mathbf{a} , define its normalization as $\bar{\mathbf{a}} := \mathbf{a}/\|\mathbf{a}\|_1$. Then, the Kullback-Leibler (KL) divergence $D_{KL}(\bar{\mathbf{1}}_{M-1} \|\bar{\mathbf{a}}_{\pm})$ is decreasing over time, and $\dot{D}_{KL}(\bar{\mathbf{1}}_{M-1} \|\bar{\mathbf{a}}_{\pm}) = 0$ iff $\bar{\mathbf{a}}_{\pm} = (M-1)^{-1}\mathbf{1}_{M-1}$.*

The proof closely follows the method of [16] for the balanced case. We drop the dimension subscripts from $\mathbf{1}_{M-1}$ and $\bar{\Phi}_M$ to reduce clutter in our proof. Our first step is to expand the KL divergence and its derivative as

$$D_{\pm} = D_{\text{KL}} \left(\frac{1}{M-1} \mathbf{1}_{M-1} \parallel \bar{\mathbf{a}}_{\pm} \right) = -\log(M-1) + \log(\mathbf{1}^T \mathbf{a}_{\pm}) - \frac{1}{M-1} \mathbf{1}^T \log(\mathbf{a}_{\pm}),$$

where we have used $\mathbf{a}_{\pm} > 0$ to simplify the norms with summations. As a result, using the dynamics equations $\dot{\mathbf{a}}_{\pm} = 2\mathbf{a}_{\pm} \odot \mathbf{b}_{\pm}$, the derivative expression is

$$\dot{D}_{\pm} = 2 \frac{\mathbf{b}_{\pm}^T \mathbf{a}_{\pm}}{\mathbf{1}^T \mathbf{a}_{\pm}} - 2 \frac{\mathbf{1}^T \mathbf{b}_{\pm}}{M-1} = \frac{2}{\mathbf{1}^T \mathbf{a}_{\pm}} \left(\mathbf{b}_{\pm}^T \mathbf{a}_{\pm} - \frac{\mathbf{1}^T \mathbf{a}_{\pm}}{M-1} \mathbf{1}^T \mathbf{b}_{\pm} \right).$$

The expressions simplify to

$$\begin{aligned} \dot{D}_+ &= \frac{2MR e^{\gamma/M}}{\mathbf{1}^T \mathbf{a}_+} \left(\frac{-\mathbf{q}_+^T \bar{\Phi}^T \exp(\bar{\Phi} \mathbf{q}_+) + \frac{1}{M-1} (\mathbf{1}^T \mathbf{q}_+) \mathbf{1}^T \bar{\Phi}^T \exp(\bar{\Phi} \mathbf{q}_+)}{\mathbf{1}^T \exp(\mathbf{x}_+) + M} \right) \\ \dot{D}_- &= \frac{2M e^{\delta/M}}{\mathbf{1}^T \mathbf{a}_-} \left(\frac{-\mathbf{q}_-^T \bar{\Phi}^T \exp(\bar{\Phi} \mathbf{q}_-) + \frac{1}{M-1} (\mathbf{1}^T \mathbf{q}_-) \mathbf{1}^T \bar{\Phi}^T \exp(\bar{\Phi} \mathbf{q}_-)}{\mathbf{1}^T \exp(\mathbf{x}_-) + M} \right), \end{aligned}$$

where $\mathbf{q}_+ := (\sqrt{RM})^{-1} \mathbf{a}_+$, $\mathbf{q}_- = \mathbf{a}_-$. The above expressions suggest a similar form in the numerator of both derivatives, with strictly positive denominators and coefficients. In the next step, define $\mathbf{z} := \bar{\Phi} \mathbf{q}$. Recall that the first row of $\bar{\Phi}$ is $\mathbf{1}^T$, from which $z_0 = \mathbf{1}^T \mathbf{q}$ follows, and $\bar{\Phi} \mathbf{1} = M \mathbf{e}_0 - \mathbf{1}$. Using these identities, the numerators can be simplified to

$$\begin{aligned} -\mathbf{z}^T \exp(\mathbf{z}) + \frac{z_0}{M-1} (M \mathbf{e}_0 - \mathbf{1})^T \exp(\mathbf{z}) &= -\sum_{i=1}^{M-1} z_i e^{z_i} + \frac{z_0}{M-1} \sum_{i=1}^{M-1} e^{z_i} \\ &= -\left(\sum_{i=1}^{M-1} z_i e^{z_i} - \frac{1}{M-1} \left(\sum_{i=1}^{M-1} z_i \right) \left(\sum_{j=1}^{M-1} e^{z_j} \right) \right) \\ &= -\frac{1}{2(M-1)} \sum_{i,j=1}^{M-1} (z_i - z_j) (e^{z_i} - e^{z_j}) \end{aligned}$$

Notice that each term $(z_i - z_j)(e^{z_i} - e^{z_j})$ is nonnegative regardless of the values of z_i, z_j , and vanishes if and only if $z_i = z_j$. As a result, the summation is nonnegative, with strict inequality unless $\forall i, j \geq 1 : z_i = z_j$. This in turn leads to $\mathbf{z} = (M \mathbf{e}_0 - \mathbf{1}) \mathbf{q}$ for some scalar q . Using the definition $\mathbf{z} = \bar{\Phi} \mathbf{q}$ and solving for \mathbf{q} gives the unique answer $\mathbf{q} = q \mathbf{1}$ for the derivative to vanish.

Hence the KL-divergence is monotonic decreasing and serves as a Lyapunov function pushing \mathbf{a}_{\pm} towards uniformity.

F.3. Proof of Theorem 7

We prove the theorem by separately showing that majority emerges before minority (Prop. 41), and then, we show it emerges before maj-min (Prop. 42).

Throughout, let $R > 1$ and

$$\mathbf{a}_t^\varepsilon = (\gamma_t^\varepsilon, \delta_t^\varepsilon, \mathbf{a}_{+,t}^\varepsilon, \mathbf{a}_{-,t}^\varepsilon)$$

be the ODE trajectory in Eq. (8) under positive initialization satisfying, for some constants $C, c > 0$,

$$\|\mathbf{a}^\varepsilon(0)\|_\infty \leq C\varepsilon, \quad \text{and} \quad c\varepsilon \leq \|\mathbf{a}_+^\varepsilon(0)\|_\infty. \quad (81)$$

For any $\eta > 0$, define

$$T_+^\varepsilon(\eta) := \inf\{t \geq 0 : \|\mathbf{a}_{+,t}^\varepsilon\|_\infty = \eta\},$$

and

$$T_0^\varepsilon(\eta) := \inf\{t \geq 0 : a_{+,t}^\varepsilon = \eta\},$$

where recall that $a_{+,t}^\varepsilon = \frac{1}{\sqrt{2}} \sqrt{R\gamma_t^\varepsilon(t)^2 + \delta_t^\varepsilon(t)^2}$. We will show that for small enough η , there exists small enough initialization scale ε such that $T_+^\varepsilon < T_-^\varepsilon$ and $T_+^\varepsilon < T_0^\varepsilon$.

Proposition 41 (Majority emerges first) *Assume (81). Then there exists $\eta_0 > 0$ such that for every $\eta \in (0, \eta_0)$, there exists $\varepsilon_0(\eta) > 0$ for which, for all $\varepsilon \in (0, \varepsilon_0(\eta))$,*

$$T_+^\varepsilon(\eta) < T_-^\varepsilon(\eta).$$

Proposition 42 (Majority mode emerges before the maj-min mode) *Assume (81). Then there exists $\eta_0 > 0$ such that for every $\eta \in (0, \eta_0)$, there exists $\varepsilon_0(\eta) > 0$ such that, for all $\varepsilon \in (0, \varepsilon_0(\eta))$,*

$$T_+^\varepsilon(\eta) < T_0^\varepsilon(\eta).$$

F.3.1. AUXILIARY LEMMAS

Before proving the two main propositions, we need some auxiliary results to locally bound the coefficients \mathbf{b}_\pm and X_\pm that enter the ODE. These will allow quantifying the escape times T_\pm, T_0 in the next section.

Lemma 43 (Local control of \mathbf{b}_\pm) *Fix any $\rho > 0$ and denote $E(u) := e^u(e^u - 1)$. Then,*

$$\forall t \in [0, T_+(\rho)], \quad \mathbf{0} \leq \sqrt{R}\mathbf{1} - \mathbf{b}_{+,t}(t) \leq \sqrt{R}E\left(\frac{\rho}{\sqrt{R}}\right) \cdot \mathbf{1}, \quad (82a)$$

$$\forall t \in [0, T_-(\rho)], \quad \mathbf{0} \leq \mathbf{1} - \mathbf{b}_{-,t}(t) \leq E(\rho) \cdot \mathbf{1}. \quad (82b)$$

Proof For convenience, recall that

$$\mathbf{b}_+ = \sqrt{R}\left(\mathbf{1} - \frac{X_+}{K}\bar{\Phi}^\top e^{\mathbf{x}_+}\right), \quad \mathbf{b}_- = \mathbf{1} - \frac{X_-}{K}\bar{\Phi}^\top e^{\mathbf{x}_-},$$

with

$$\mathbf{x}_+ = \frac{1}{M}\left(\gamma\mathbf{1} + \frac{1}{\sqrt{R}}\bar{\Phi}\mathbf{a}_+\right), \quad \mathbf{x}_- = \frac{1}{M}\left(\delta\mathbf{1} + \bar{\Phi}\mathbf{a}_-\right).$$

The non-negativity bounds in Eq. (82) are proved in Lemma 21. Here, we prove the upper bounds.

The key observation is that the estimates for \mathbf{b}_\pm on this interval depend only on \mathbf{a}_\pm , and *not* on γ, δ . To see this, define for convenience

$$\mathbf{y}_+ := \frac{1}{M\sqrt{R}}\bar{\Phi}\mathbf{a}_+$$

to write

$$\frac{X_+ \bar{\Phi}^\top e^{x_+}}{K} = \frac{e^{\gamma/M}}{e^{\gamma/M} \mathbf{1}^\top e^{y_+} + M} \bar{\Phi}^\top e^{y_+}.$$

Using that the multiplicative scalar is positive,

$$\left\| \frac{X_+ \bar{\Phi}^\top e^{x_+}}{K} \right\|_\infty \leq \frac{1}{\mathbf{1}^\top e^{y_+}} \|\bar{\Phi}^\top e^{y_+}\|_\infty.$$

Since $\bar{\Phi}^\top \mathbf{1} = 0 \implies \bar{\Phi}^\top e^{y_+} = \bar{\Phi}^\top (e^{y_+} - \mathbf{1})$, and $\bar{\Phi} \in \{\pm 1\}^{M \times (M-1)}$ we have

$$\|\bar{\Phi}^\top e^{y_+}\|_\infty \leq M \|e^{y_+} - \mathbf{1}\|_\infty.$$

Also,

$$\mathbf{1}^\top e^{y_+} \geq M e^{-\|y_+\|_\infty}.$$

Putting the above displays together,

$$\left\| \frac{X_+ \bar{\Phi}^\top e^{x_+}}{K} \right\|_\infty \leq e^{\|y_+\|_\infty} \|e^{y_+} - \mathbf{1}\|_\infty.$$

Now, on $[0, T_+)$, using again that $\bar{\Phi} \in \{\pm 1\}^{M \times (M-1)}$:

$$\|y_+\|_\infty = \left\| \frac{1}{M\sqrt{R}} \bar{\Phi} a_+ \right\|_\infty \leq \frac{1}{\sqrt{R}} \|a_+\|_\infty \leq \frac{\rho}{\sqrt{R}}.$$

Hence also, by the coordinatewise bound $|e^x - 1| \leq e^{|x|} - 1$,

$$\|e^{y_+} - \mathbf{1}\|_\infty \leq e^{\|y_+\|_\infty} - 1 \leq e^{\rho/\sqrt{R}} - 1.$$

Put together, we have shown that

$$\left\| \frac{X_+ \bar{\Phi}^\top e^{x_+}}{K} \right\|_\infty \leq e^{\rho/\sqrt{R}} (e^{\rho/\sqrt{R}} - 1)$$

which immediately implies Eq. (82a). The bound Eq. (82b) follows in the exact same way and is omitted for brevity. \blacksquare

Lemma 44 (Local control of X_\pm) Fix $\rho > 0$, and define the local stopping time

$$\tau_\rho := \inf \left\{ t \geq 0 : \max \left\{ \gamma_t(t), \delta_t(t), \|\mathbf{a}_{+,t}(t)\|_\infty, \|\mathbf{a}_{-,t}(t)\|_\infty \right\} = \rho \right\}.$$

Then, for every $t \in [0, \tau_\rho)$,

$$\frac{2}{1 + e^{2\rho}} =: \underline{X}(\rho) \leq X_{+,t}, X_{-,t} \leq \bar{X}(\rho) := \frac{2}{1 + e^{-2\rho}},$$

Proof Recall that

$$X_+ = \frac{2M}{\mathbf{1}^\top e^{x_+} + M}, \quad X_- = \frac{2M}{\mathbf{1}^\top e^{x_-} + M}.$$

Fix $t \in [0, \tau_\rho)$ such that $\gamma_t, \delta_t, \|\mathbf{a}_{+,t}\|_\infty, \|\mathbf{a}_{-,t}\|_\infty < \rho$. Using $\bar{\Phi} \in \{\pm 1\}^{M \times (M-1)}$, we have

$$\|\bar{\Phi} \mathbf{a}_{\pm,t}\|_\infty \leq M \|\mathbf{a}_{\pm,t}\|_\infty \implies \|\mathbf{x}_{+,t}\|_\infty \leq \frac{\gamma_t}{M} + \frac{1}{M\sqrt{R}} \|\bar{\Phi} \mathbf{a}_{+,t}\|_\infty \leq \frac{\rho}{M} + \frac{\rho}{\sqrt{R}} \leq 2\rho.$$

Similarly, $\|\mathbf{x}_{-,t}\|_\infty \leq 2\rho$. Therefore

$$M e^{-2\rho} \leq \mathbf{1}^\top e^{x_{\pm,t}} \leq M e^{2\rho}.$$

Substituting this into the definitions of X_\pm gives the desired. \blacksquare

F.3.2. PROOF OF PROPOSITION 41

Recall that

$$\dot{\mathbf{a}}_{\pm} = 2\mathbf{a}_{\pm} \odot \mathbf{b}_{\pm}.$$

Fix $\rho > 0$ and assume first that $\eta \in (0, \rho)$. By Lemma 43, for every $t \in [0, T_+^{\varepsilon}(\rho)]$,

$$\|\mathbf{b}_{+,t} - \sqrt{R}\mathbf{1}\|_{\infty} \leq \sqrt{R}E\left(\frac{\rho}{\sqrt{R}}\right) \leq \sqrt{R}E(\rho).$$

Hence, for every coordinate i and every $t \in [0, T_+^{\varepsilon}(\rho)]$,

$$\sqrt{R}(1 - E(\rho)) \leq b_{+,i}(t) \leq \sqrt{R},$$

and therefore

$$2\sqrt{R}(1 - E(\rho))a_{+,i}(t) \leq \dot{a}_{+,i}(t) \leq 2\sqrt{R}a_{+,i}(t).$$

Applying Gronwall's inequality coordinatewise and taking maximum over i gives, for every $t \in [0, T_+^{\varepsilon}(\rho)]$,

$$\|\mathbf{a}_{+,0}^{\varepsilon}\|_{\infty} e^{2\sqrt{R}(1-E(\rho))t} \leq \|\mathbf{a}_{+,t}^{\varepsilon}\|_{\infty} \leq \|\mathbf{a}_{+,0}^{\varepsilon}\|_{\infty} e^{2\sqrt{R}t}. \quad (83)$$

In the exact same manner, for every $t \in [0, T_-^{\varepsilon}(\rho)]$,

$$\|\mathbf{a}_{-,0}^{\varepsilon}\|_{\infty} e^{2(1-E(\rho))t} \leq \|\mathbf{a}_{-,t}^{\varepsilon}\|_{\infty} \leq \|\mathbf{a}_{-,0}^{\varepsilon}\|_{\infty} e^{2t}. \quad (84)$$

Because $\eta < \rho$, we have

$$T_+^{\varepsilon}(\eta) \leq T_+^{\varepsilon}(\rho), \quad T_-^{\varepsilon}(\eta) \leq T_-^{\varepsilon}(\rho),$$

and therefore (83) may be evaluated at $t = T_+^{\varepsilon}(\eta)$, while (84) may be evaluated at $t = T_-^{\varepsilon}(\eta)$. Doing this, using

$$\|\mathbf{a}_{+,T_+^{\varepsilon}(\eta)}^{\varepsilon}\|_{\infty} = \eta, \quad \|\mathbf{a}_{-,T_-^{\varepsilon}(\eta)}^{\varepsilon}\|_{\infty} = \eta,$$

and rearranging gives

$$\frac{1}{2\sqrt{R}} \log \frac{\eta}{\|\mathbf{a}_{+,0}^{\varepsilon}\|_{\infty}} \leq T_+^{\varepsilon}(\eta) \leq \frac{1}{2\sqrt{R}(1-E(\rho))} \log \frac{\eta}{\|\mathbf{a}_{+,0}^{\varepsilon}\|_{\infty}}, \quad (85)$$

$$\frac{1}{2} \log \frac{\eta}{\|\mathbf{a}_{-,0}^{\varepsilon}\|_{\infty}} \leq T_-^{\varepsilon}(\eta) \leq \frac{1}{2(1-E(\rho))} \log \frac{\eta}{\|\mathbf{a}_{-,0}^{\varepsilon}\|_{\infty}}, \quad (86)$$

Because $R > 1$ and $E(\rho) \rightarrow 0$ as $\rho \rightarrow 0$, we may choose $\eta_0 > 0$ such that for every $\rho \in (0, \eta_0)$,

$$\sqrt{R}(1 - E(\rho)) > 1. \quad (87)$$

Fix any $\eta \in (0, \eta_0)$ and choose any $\rho \in (\eta, \eta_0)$. Then (85), (86), together with

$$c\varepsilon \leq \|\mathbf{a}_{+,0}^{\varepsilon}\|_{\infty}, \quad \|\mathbf{a}_{-,0}^{\varepsilon}\|_{\infty} \leq \|\mathbf{a}^{\varepsilon}(0)\|_{\infty} \leq C\varepsilon,$$

imply

$$T_+^{\varepsilon}(\eta) \leq \frac{1}{2\sqrt{R}(1-E(\rho))} \log \frac{\eta}{c\varepsilon},$$

and

$$T_-^{\varepsilon}(\eta) \geq \frac{1}{2} \log \frac{\eta}{C\varepsilon}.$$

By (87) the upper bound for $T_+^{\varepsilon}(\eta)$ has a strictly smaller coefficient in front of $\log(1/\varepsilon)$ than the lower bound for $T_-^{\varepsilon}(\eta)$. Therefore, for all sufficiently small ε , $T_+^{\varepsilon}(\eta) < T_-^{\varepsilon}(\eta)$, as desired.

F.3.3. PROOF OF PROPOSITION 42

We suppress the superscript ε when clear from context. Recall that

$$\dot{\gamma} = \frac{(2R\gamma^2 + \delta^2)X_+ + \gamma\delta X_-}{\sqrt{2(R\gamma^2 + \delta^2)}}, \quad \dot{\delta} = \frac{(R\gamma^2 + 2\delta^2)X_- + R\gamma\delta X_+}{\sqrt{2(R\gamma^2 + \delta^2)}}.$$

By direct differentiation,

$$\dot{a}_0(t) = \frac{R\gamma(t)\dot{\gamma}(t) + \delta(t)\dot{\delta}(t)}{2a_0(t)}.$$

Substituting the formulas for $\dot{\gamma}$, $\dot{\delta}$, and simplifying using $2a_0(t)^2 = R\gamma(t)^2 + \delta(t)^2$, yields

$$\dot{a}_0(t) = R\gamma(t)X_{+,t} + \delta(t)X_{-,t}. \quad (88)$$

By Lemma 44, for every $t \in [0, \tau_\rho)$,

$$X_{+,t}, X_{-,t} \leq \bar{X}(\rho), \quad \bar{X}(\rho) = \frac{2}{1 + e^{-2\rho}}.$$

Hence, on $[0, \tau_\rho)$,

$$\dot{a}_0(t) \leq \bar{X}(\rho)(R\gamma(t) + \delta(t)) \leq \sqrt{2(R+1)} \bar{X}(\rho) a_0(t), \quad (89)$$

where the second inequality follows from $R\gamma + \delta \leq \sqrt{R+1}\sqrt{R\gamma^2 + \delta^2} = \sqrt{2(R+1)} a_0$.

Next, by Lemma 43, exactly as in the proof of Proposition 41,

$$\|\mathbf{a}_{+,t}(t)\|_\infty \geq \|\mathbf{a}_+(0)\|_\infty e^{2\sqrt{R}(1-E(\rho))t}, \quad t \in [0, T_+(\rho)]. \quad (90)$$

Let

$$\mu(\rho) := 2\sqrt{R}(1 - E(\rho)), \quad \nu(\rho) := \sqrt{2(R+1)} \bar{X}(\rho).$$

Since $E(\rho) \rightarrow 0$ and $\bar{X}(\rho) \rightarrow 1$ as $\rho \rightarrow 0$, we have

$$\mu(\rho) \rightarrow 2\sqrt{R}, \quad \nu(\rho) \rightarrow \sqrt{2(R+1)}.$$

Because $R > 1$,

$$2\sqrt{R} > \sqrt{2(R+1)}.$$

Hence there exists $\rho_{\text{rate}} > 0$ such that

$$\mu(\rho) > \nu(\rho) \quad \text{for all } \rho \in (0, \rho_{\text{rate}}). \quad (91)$$

By Proposition 41, there exists $\rho_{\text{maj}} > 0$ such that its conclusion holds for every threshold $\rho \in (0, \rho_{\text{maj}})$. Define

$$\rho_0 := \min\{\rho_{\text{rate}}, \rho_{\text{maj}}\}, \quad \eta_0 := \rho_0/2.$$

Fix any $\eta \in (0, \eta_0)$, and set

$$\rho := 2\eta.$$

Then $\rho \in (0, \rho_0)$, so both (91) and Proposition 41 apply at threshold ρ .

We claim that $T_+^\varepsilon(\eta) < T_0^\varepsilon(\eta)$ for all sufficiently small ε . Suppose, toward a contradiction, that

$$T_0^\varepsilon(\eta) \leq T_+^\varepsilon(\eta),$$

and set

$$T := T_0^\varepsilon(\eta).$$

For sufficiently small ε , the initialization satisfies

$$\|\mathbf{a}_\pm^\varepsilon(0)\|_\infty < \eta, \quad a_0^\varepsilon(0) < \eta.$$

Since $0 < \eta < \rho$, continuity gives

$$T_+^\varepsilon(\eta) < T_+^\varepsilon(\rho).$$

Hence $T < T_+^\varepsilon(\rho)$. Moreover, by Proposition 41, after possibly decreasing $\varepsilon_0(\eta)$,

$$T_+^\varepsilon(\rho) < T_-^\varepsilon(\rho).$$

Therefore

$$T < T_+^\varepsilon(\rho), \quad T < T_-^\varepsilon(\rho).$$

Also, for every $t \leq T$, we have $a_0(t) \leq \eta$, and hence, since $\rho = 2\eta$:

$$\gamma(t) \leq \sqrt{\frac{2}{R}}\eta < \rho, \quad \delta(t) \leq \sqrt{2}\eta < \rho,$$

Thus

$$T < \tau_\rho, \tag{92}$$

and we can plug in T to (89) and (90).

Specifically, By (89), (92), and Gronwall's inequality,

$$a_0(T) \leq a_0(0)e^{\nu(\rho)T}.$$

Since $a_0(T) = \eta$, this implies

$$e^T \geq \left(\frac{\eta}{a_0(0)}\right)^{1/\nu(\rho)}. \tag{93}$$

Using (90), (93), and $\|\mathbf{a}_+(0)\|_\infty \geq c\varepsilon$, we get $\|\mathbf{a}_{+,t}(T)\|_\infty \geq c\varepsilon \left(\frac{\eta}{a_0(0)}\right)^{\mu(\rho)/\nu(\rho)}$. Since

$$a_0(0) \leq C\sqrt{\frac{R+1}{2}}\varepsilon,$$

we obtain

$$\|\mathbf{a}_{+,T}\|_\infty \geq c \left(C\sqrt{\frac{R+1}{2}}\right)^{-\mu(\rho)/\nu(\rho)} \eta^{\mu(\rho)/\nu(\rho)} \varepsilon^{1-\mu(\rho)/\nu(\rho)}. \tag{94}$$

By (91), $\mu(\rho) > \nu(\rho)$. Hence the exponent $1 - \mu(\rho)/\nu(\rho)$ is negative, so the right-hand side of (94) diverges as $\varepsilon \rightarrow 0$. In particular, for all sufficiently small ε ,

$$\|\mathbf{a}_{+,T}\|_\infty > \eta.$$

But $T = T_0^\varepsilon(\eta) \leq T_+^\varepsilon(\eta)$, so by definition of $T_+^\varepsilon(\eta)$, $\|\mathbf{a}_{+,T}\|_\infty \leq \eta$, a contradiction. Therefore $T_+^\varepsilon(\eta) < T_0^\varepsilon(\eta)$ completing the proof.

F.4. Lyapunov function counterexample

Here, we provide a counterexample showing that the KL distance cannot be a Lyapunov function under our diagonal-plus-rank-one parameterization even for $R = 1$. This is in contrast to the diagonal parameterization for $R = 1$ by [16]. Of course, as made clear by our work, the diagonal-plus-rank-one parameterization is necessary to capture the RP and the GF evolution in the general imbalanced case.

Proposition 45 (Counterexample to KL monotonicity) *Let $R = 1$ and $K = 4$. Define*

$$\mathbf{z} = \left[\sqrt{\frac{\gamma^2 + \delta^2}{2}}, a_+, a_- \right], \quad \bar{\mathbf{z}} := \frac{\mathbf{z}}{\|\mathbf{z}\|_1}, \quad \bar{\mathbf{1}} := \frac{1}{3} \mathbf{1}_3.$$

Then $D_{\text{KL}}(\bar{\mathbf{1}} \|\bar{\mathbf{z}})$ is not necessarily decreasing along the ODE. More precisely, take

$$\gamma = \kappa \varepsilon, \quad \delta = \varepsilon, \quad a_+ = a_- = \varepsilon,$$

with $0 < \kappa < 1$ and $\varepsilon > 0$. Then, for all sufficiently small $\varepsilon > 0$, $\dot{D}_{\text{KL}}(\bar{\mathbf{1}} \|\bar{\mathbf{z}}) > 0$.

Proof For $K = 4$, $R = 1$, we have

$$\mathbf{x}_+ = \frac{1}{2} \begin{bmatrix} \gamma + a_+ \\ \gamma - a_+ \end{bmatrix}, \quad \mathbf{x}_- = \frac{1}{2} \begin{bmatrix} \delta + a_- \\ \delta - a_- \end{bmatrix}, \quad X_+ = \frac{4}{\mathbf{1}^\top e^{\mathbf{x}_+} + 2}, \quad X_- = \frac{4}{\mathbf{1}^\top e^{\mathbf{x}_-} + 2},$$

and

$$b_+ = 1 - \frac{X_+}{4} [1 \quad -1] e^{\mathbf{x}_+}, \quad b_- = 1 - \frac{X_-}{4} [1 \quad -1] e^{\mathbf{x}_-}.$$

As in the proposition's statement, let $\gamma = \kappa \varepsilon$, $\delta = \varepsilon$, $a_+ = a_- = \varepsilon$, and

$$c := \sqrt{\frac{\kappa^2 + 1}{2}}. \tag{95}$$

Then $\mathbf{z} = \varepsilon [c \quad 1 \quad 1]^\top$. Hence

$$\bar{\mathbf{z}} = \frac{\mathbf{z}}{\|\mathbf{z}\|_1} = \frac{1}{c + 2} \begin{bmatrix} c \\ 1 \\ 1 \end{bmatrix}.$$

Let $V(\mathbf{z}) := D_{\text{KL}}(\bar{\mathbf{1}} \|\bar{\mathbf{z}})$. With $r_i := \dot{z}_i / z_i$, we have

$$\dot{V} = \sum_{i=0}^2 \bar{z}_i r_i - \frac{1}{3} \sum_{i=0}^2 r_i \tag{96}$$

As $\varepsilon \rightarrow 0^+$, we have

$$\mathbf{x}_+ \rightarrow 0, \quad \mathbf{x}_- \rightarrow 0,$$

and therefore

$$X_+ \rightarrow 1, \quad X_- \rightarrow 1, \quad b_+ \rightarrow 1, \quad b_- \rightarrow 1.$$

Consequently,

$$r_1 := \frac{\dot{a}_+}{a_+} = 2b_+ \rightarrow 2, \quad r_2 := \frac{\dot{a}_-}{a_-} = 2b_- \rightarrow 2.$$

Moreover, calling $z_0 = \sqrt{\frac{\gamma^2 + \delta^2}{2}}$, we have

$$\frac{\dot{z}_0}{z_0} = \frac{\gamma\dot{\gamma} + \delta\dot{\delta}}{\gamma^2 + \delta^2},$$

and substituting $\gamma = \kappa\varepsilon$, $\delta = \varepsilon$, and using $X_+, X_- \rightarrow 1$, we obtain

$$r_0 := \frac{\dot{z}_0}{z_0} \rightarrow \frac{\kappa + 1}{c}.$$

Thus, from Eq. (96)

$$\lim_{\varepsilon \rightarrow 0^+} \dot{V} = \frac{c(\kappa + 1)/c + 2 + 2}{c + 2} - \frac{1}{3}((\kappa + 1)/c + 2 + 2) = \frac{\kappa + 5}{c + 2} - \frac{1}{3} \left(\frac{\kappa + 1}{c} + 4 \right).$$

To see that this limit is strictly positive for $\kappa \in (0, 1)$, we can rearrange the inequality $\lim_{\varepsilon \rightarrow 0^+} \dot{V} > 0$ to $c(\kappa + 3) > \kappa^2 + \kappa + 2$. Squaring both sides and substituting $c^2 = (\kappa^2 + 1)/2$ reduces the condition to:

$$\kappa^4 - 2\kappa^3 + 2\kappa - 1 < 0 \implies (\kappa + 1)(\kappa - 1)^3 < 0.$$

Since $\kappa \in (0, 1)$, this inequality strictly holds, completing the proof. \blacksquare

Appendix G. Numerical Results

G.1. Regularization Path Experiments

In this section, we provide additional numerical results supporting our analysis. We write σ_+ , σ_- , and $\sigma_{+/-}$ for the singular values corresponding, respectively, to the majority, minority, and majority–minority modes. Since these singular values diverge along the CE regularization path as $\lambda \rightarrow 0$, we compare their relative evolution after factoring out both the overall scale and the limiting SEL profile.

Specifically, denote in this section the singular values of the SEL matrix as

$$s_+^* = \sqrt{R}, \quad s_-^* = 1, \quad s_{+/-}^* = \sqrt{\frac{R + 1}{2}},$$

and define the path-dependent scale

$$C_\lambda := \frac{|\sigma_{+, \lambda}| + |\sigma_{-, \lambda}| + |\sigma_{+/-, \lambda}|}{s_+^* + s_-^* + s_{+/-}^*}.$$

We then plot the rescaled coordinates

$$\hat{\sigma}_{+, \lambda} := \frac{\sigma_{+, \lambda}}{C_\lambda \sqrt{R}}, \quad \hat{\sigma}_-(\lambda) := \frac{\sigma_{-, \lambda}}{C_\lambda}, \quad \hat{\sigma}_{+/-, \lambda} := \frac{\sigma_{+/-, \lambda}}{C_\lambda \sqrt{(R + 1)/2}}.$$

With this normalization, all three coordinates converge to 1 whenever the singular-value vector converges in direction to the SEL profile. Thus, the plots emphasize deviations from the limiting geometry rather than the overall divergence of the singular values.

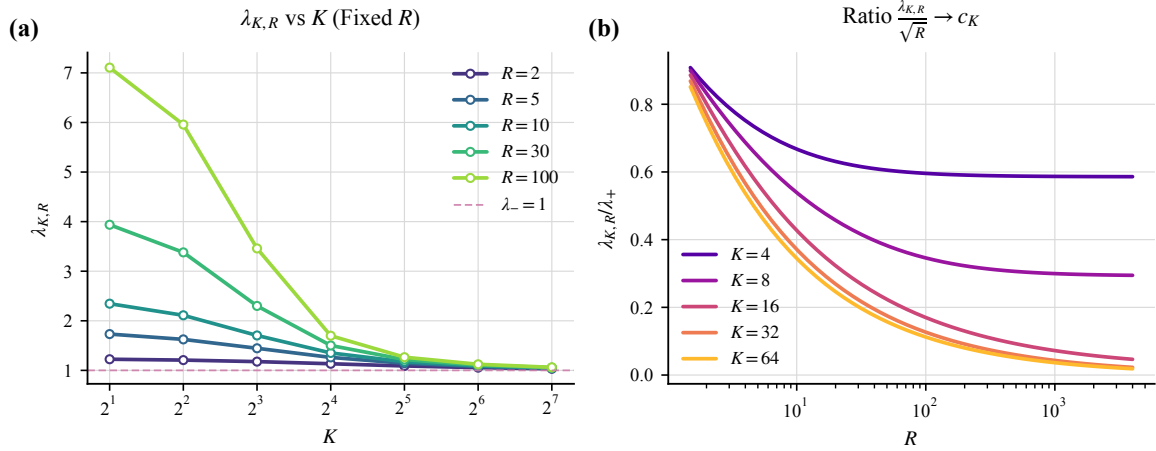


Figure 2: Behavior of regularization-path emergence thresholds for varying number of classes (K) and imbalance ratios (R).

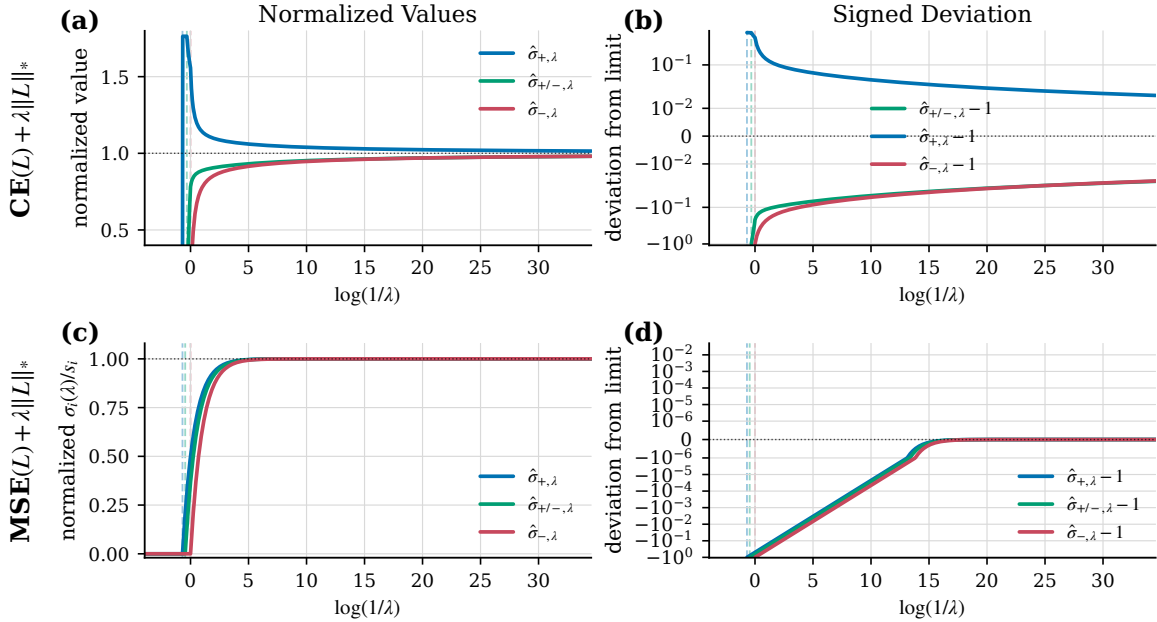


Figure 3: Comparison of the regularization-path trajectories for the three normalized singular coordinates under CE and MSE, shown for $R = 4$ and $K = 8$. The panels track the majority, majority–minority, and minority coordinates as functions of λ ; horizontal dashed lines indicate their limiting normalized values

G.1.1. REGULARIZATION PATH COMPARISON FOR DIFFERENT K AND R

We complement the theoretical analysis with numerical comparisons of three ways of tracing the regularization path. For each value of λ , we overlay the singular values predicted by our analytic solution, the singular values obtained by solving the convex logit-space problem ($Z\text{-UFM}_\lambda$) in CVXPY, and the singular values obtained by projected gradient descent on the sphere of radius B

HOW CROSS-ENTROPY LEARNS DATA MODES: EMERGENCE AND IMPLICIT BIAS IN THE UFM

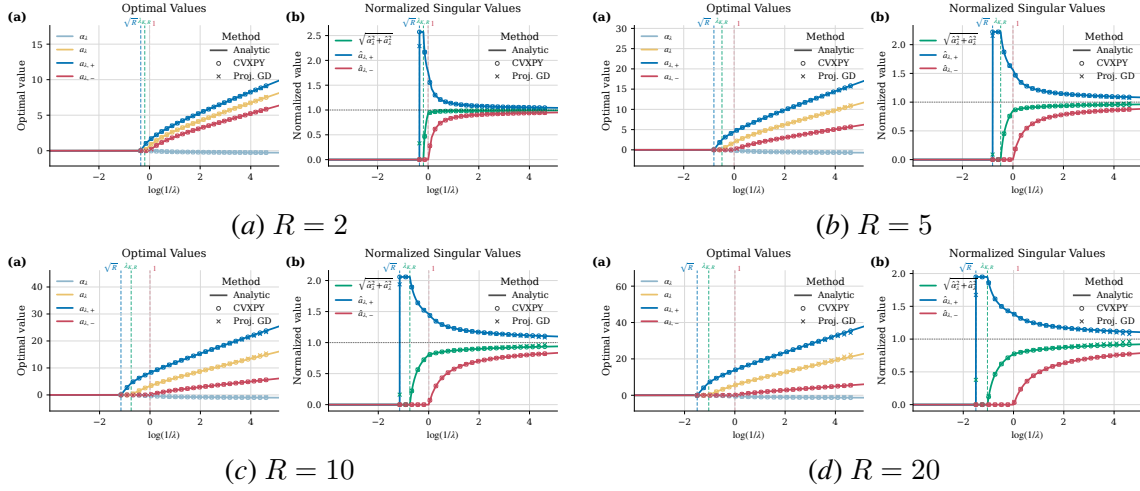


Figure 4: Numerical comparison of the regularization-path singular values for $K = 4$ across different imbalance ratios R . In each panel, the curves obtained from the analytic solution are overlaid with those produced by a direct CVXPY solve of $(\mathbf{Z}\text{-UFM}_\lambda)$ and by projected gradient descent on the sphere with radius B matched to the same value of λ .

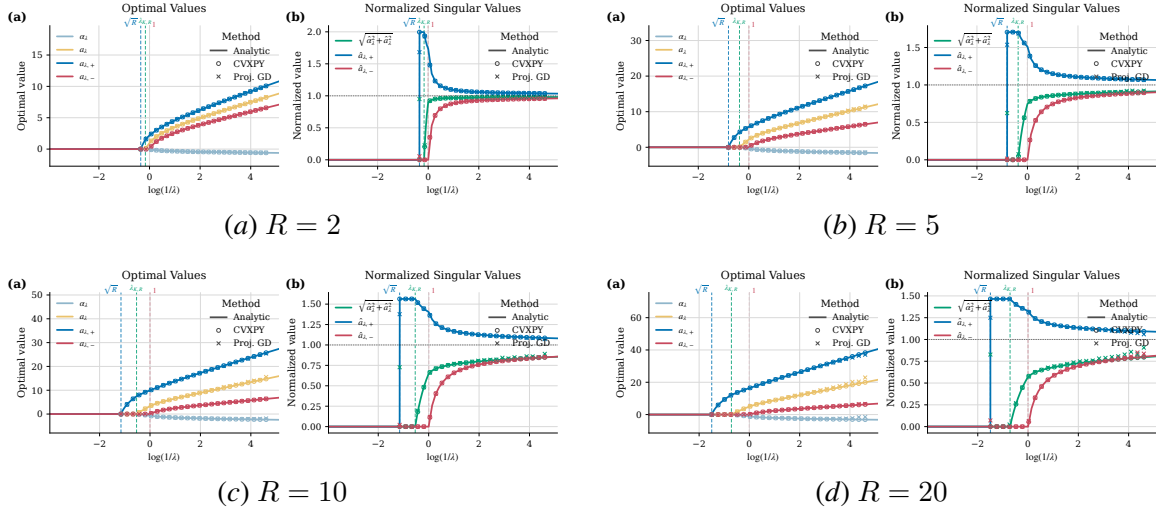


Figure 5: Numerical comparison of the regularization-path singular values for $K = 8$ across different imbalance ratios R . In each panel, the curves obtained from the analytic solution are overlaid with those produced by a direct CVXPY solve of $(\mathbf{Z}\text{-UFM}_\lambda)$ and by projected gradient descent on the sphere with radius B matched to the same value of λ .

corresponding to that value of λ . This comparison serves two purposes: it verifies the correctness of the analytic formulas, and it shows that the projected optimization procedure tracks the same path across a range of imbalance ratios and class counts. The resulting comparisons are shown in Figs. 4 and 5.

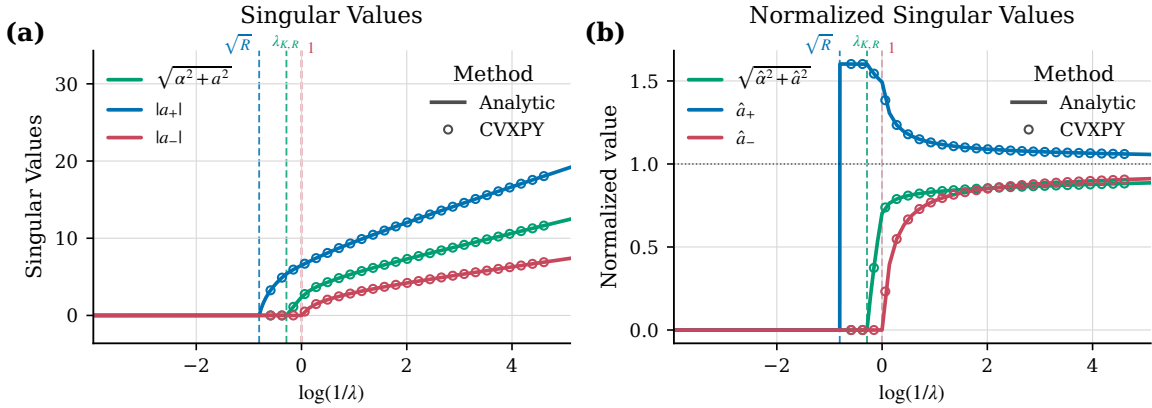


Figure 6: Regularization-path validation for $K = 12$ and $R = 5$. The figure compares the singular-coordinate trajectories predicted by the analytic formulas with those obtained from a direct CVXPY solve of the convex nuclear-norm problem ($\mathcal{Z}\text{-UFM}_\lambda$). The numerical agreement supports the claim that the regularization-path description is not an artifact of the Sylvester–Hadamard coordinate construction.

Beyond powers of two. Our presentation in the main text assumes $K = 2^m$, which allows us to use Sylvester–Hadamard matrices and keeps the spectral coordinates especially explicit. This assumption is convenient for the proof, but the resulting regularization-path formulas depend only on the three-mode structure induced by the step imbalance: majority, majority–minority, and minority. This raises the possibility that the same path description remains valid more generally, even when K is not a power of two and the particular Hadamard construction is unavailable. As a numerical check, we compare in Fig. 6 the analytic singular-coordinate trajectories with a direct CVXPY solution of the convex nuclear-norm problem for $K = 12$ and $R = 5$. The close agreement provides evidence that the regularization-path characterization is not merely an artifact of the Sylvester–Hadamard coordinates.

Role of the nuclear norm. The phase transitions identified in the regularization path are not a generic consequence of imbalance alone. To illustrate this, we repeat the same path computation after replacing the nuclear-norm penalty in ($\mathcal{Z}\text{-UFM}_\lambda$) by a Frobenius-norm penalty in logit space. The resulting trajectories are shown in Fig. 7 for both CE and MSE losses.

The contrast with the nuclear-norm path is clear. Under Frobenius regularization, the penalty controls the overall logit scale but does not separately threshold individual singular modes. As a result, once the solution leaves zero, the singular coordinates grow together in their limiting proportions. After normalization by the limiting singular-value profile, all coordinates remain essentially constant at one along the path. Thus, the staged activation observed for the nuclear-norm regularization path disappears.

This comparison highlights the mechanism behind mode emergence in our analysis. The nuclear norm acts as an ℓ_1 -type penalty on the singular values, allowing different modes to become active at different regularization strengths. By contrast, the Frobenius norm behaves as an ℓ_2 -type penalty on the singular-value vector, coupling the modes through a common scale and thereby suppressing sequential emergence.

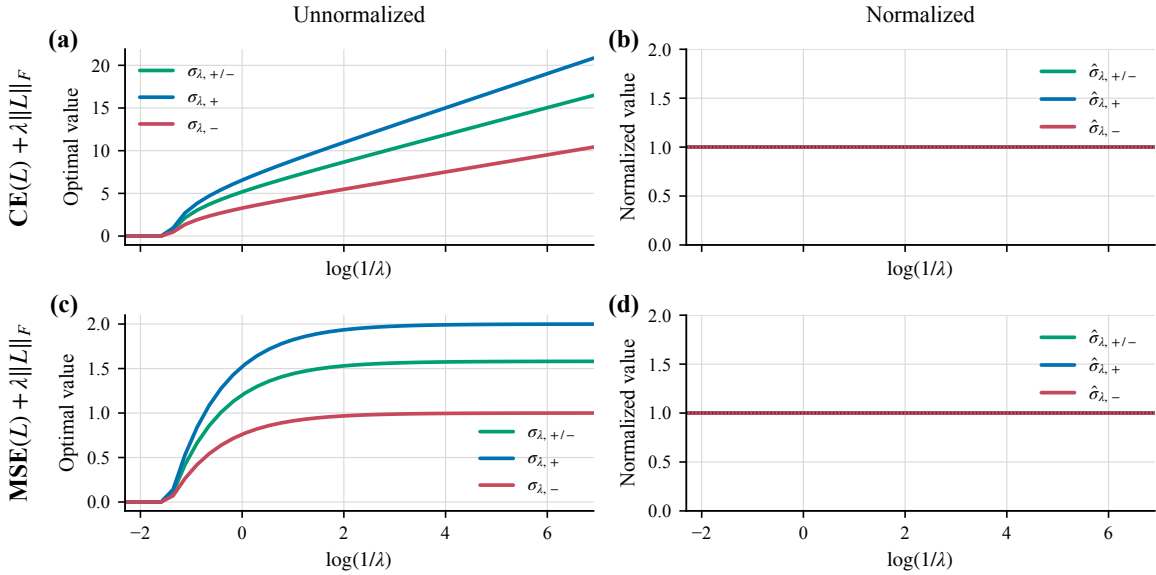


Figure 7: Regularization paths with a Frobenius-norm penalty in logit space. The top row shows the CE objective with Frobenius regularization, and the bottom row shows the corresponding MSE objective. Left panels plot the unnormalized singular coordinates, while right panels plot the same coordinates after normalization by their limiting profile. In contrast to the nuclear-norm path, the normalized coordinates remain aligned throughout the trajectory, indicating that the modes do not emerge sequentially.

G.2. Gradient-flow experiments

G.2.1. SETUP

We numerically simulate the gradient flow in (4) using gradient descent with stepsize 10^{-2} . Throughout this experiment we set $K = 4$. We compare two initialization schemes.

Hadamard initialization. We initialize the model according to the spectrally aligned construction in Theorem 5, using randomly chosen initial singular coordinates. For this initialization, the matrix-valued gradient flow is predicted to remain in the invariant subspace induced by the $(\mathbf{U}, \tilde{\mathbf{V}})$ coordinates, and its evolution should therefore be exactly captured by the reduced vector ODE derived in Theorem 5. We simulate both the original gradient descent dynamics and the reduced ODE from the same initial coordinates.

Random spectral initialization. As a comparison, we also initialize \mathbf{W} and \mathbf{H} with the same initial singular values as in the Hadamard initialization, but with random left and right singular directions. This keeps the initial scale and spectrum fixed while removing the spectral alignment required by the theory.

All results are averaged over 10 independent runs. Shaded regions indicate the standard error of the mean. Metrics are recorded at 92 geometrically spaced iterations, which allows the full trajectory to be visualized on a logarithmic time axis.

Metrics. We track six quantities. First, we record the cross-entropy loss. Second, we measure neural-collapse-style alignment through the Frobenius distances between the normalized Gram

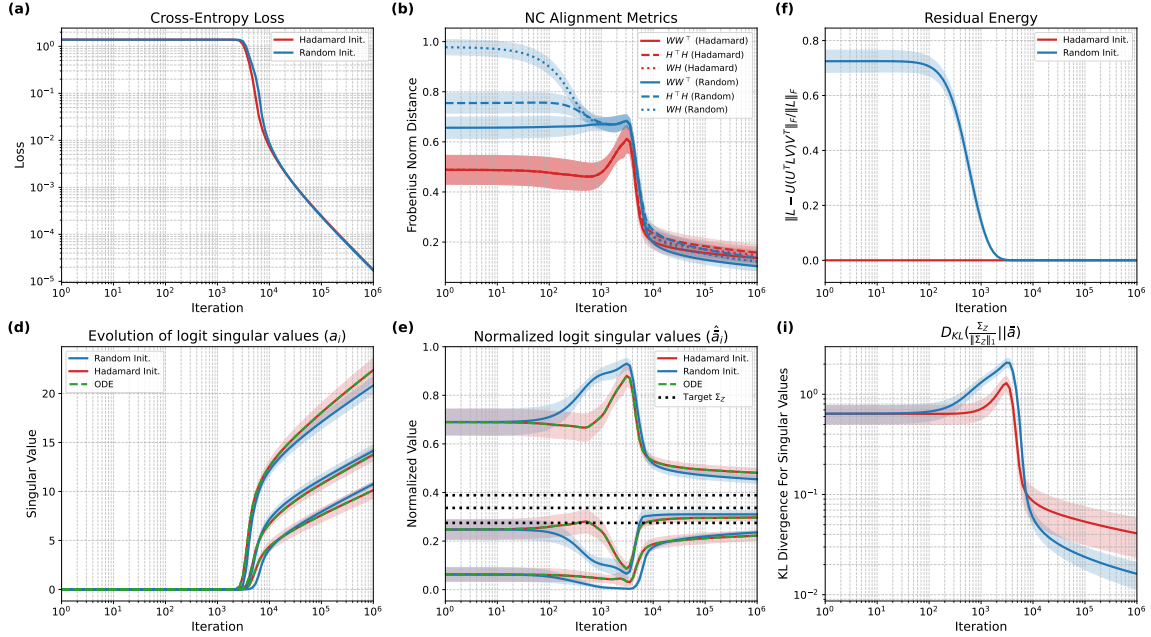


Figure 8: Gradient-flow experiments for $K = 4$ and $R = 2$. We compare Hadamard initialization, random spectral initialization with the same initial singular values, and the reduced vector ODE initialized from the Hadamard coordinates. The panels show the cross-entropy loss, neural-collapse-style alignment metrics, residual energy outside the Hadamard signal subspace, raw singular-coordinate trajectories, ℓ_1 -normalized singular-coordinate trajectories, and the KL discrepancy between the normalized singular-coordinate profile and the limiting SELI profile. Shaded regions denote standard errors over 10 runs.

matrices of \mathbf{W} , \mathbf{H} , and \mathbf{WH} and their corresponding limiting SELI geometries. Third, we measure the relative residual energy outside the signal subspace associated with the Hadamard coordinates,

$$\frac{\|\mathbf{L} - \mathbf{U}\mathbf{U}^\top \mathbf{L} \tilde{\mathbf{V}} \tilde{\mathbf{V}}^\top\|_F}{\|\mathbf{L}\|_F},$$

where $\mathbf{L} = \mathbf{WH}$ denotes the logit matrix. This quantity is zero exactly when the logits lie in the reduced subspace. Fourth, we track the singular coordinates of the logits. Fifth, we plot their ℓ_1 -normalized versions, which removes the diverging scale and isolates convergence in direction. Finally, we report the KL divergence between the normalized singular-coordinate profile and the limiting SELI singular-value profile. In the present imbalanced setting, this KL quantity is used only as a directional discrepancy measure; unlike in the balanced setting, it is not a Lyapunov function for the dynamics.

G.2.2. RESULTS

The results are summarized in Fig. 8. They show three main trends.

Reduction to the vector ODE. Under Hadamard initialization, the residual energy outside the $(\mathbf{U}, \tilde{\mathbf{V}})$ signal subspace remains at its initial value of zero throughout training. This verifies the

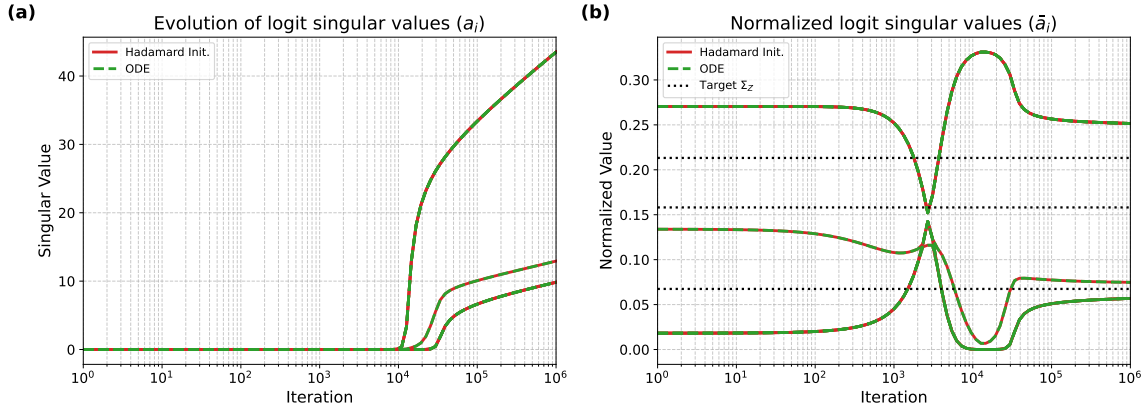


Figure 9: Preservation of block-constant spectral coordinates under gradient flow. The model uses $K = 8$ and $R = 10$, and is initialized with $\mathbf{a}_{+,0} = a_{+,0}\mathbf{1}_{M-1}$ and $\mathbf{a}_{-,0} = a_{-,0}\mathbf{1}_{M-1}$. The trajectories show that the coordinates within each majority and minority block remain identical throughout training, confirming preservation of the block-constant subspace. The same run also exhibits the qualitative phase structure observed along the regularization path: majority modes emerge first, followed by the majority–minority mode, and then the minority modes.

invariance predicted by the softmax diagonalization result in Theorem 1: the Hadamard coordinates close under the gradient-flow dynamics, so the matrix-valued trajectory does not leave the reduced subspace.

Validation of the reduced dynamics. The singular-coordinate trajectories obtained from the full gradient descent simulation under Hadamard initialization agree with the solution of the reduced vector ODE. This agreement is visible both for the raw singular coordinates and for their normalized versions. Thus, the reduced ODE in Theorem 5 accurately captures the evolution of the original matrix-valued system when the initialization is spectrally aligned.

Behavior beyond the invariant initialization. Random spectral initialization does not begin in the Hadamard signal subspace, and therefore initially has substantial residual energy outside that subspace. Nevertheless, this residual energy decreases during training, and the normalized singular-coordinate profile moves toward the same limiting SELI profile. This suggests that the Hadamard-reduced dynamics capture not only an exactly invariant initialization, but also a useful organizing description of the late-time behavior of more generic initializations. The KL diagnostic exhibits a transient increase before decreasing, which is consistent with the fact that it is not a Lyapunov function in this imbalanced setting.

G.2.3. PRESERVATION OF BLOCK-CONSTANT SPECTRAL COORDINATES

The reduced dynamics in Theorem 5 still contain vector-valued coordinates

$$\mathbf{a}_{+,t}, \mathbf{a}_{-,t} \in \mathbb{R}^{M-1},$$

corresponding to the majority and minority singular blocks. A particularly important special case is the block-constant initialization

$$\mathbf{a}_{+,0} = a_{+,0}\mathbf{1}_{M-1}, \quad \mathbf{a}_{-,0} = a_{-,0}\mathbf{1}_{M-1}.$$

This condition is preserved by the reduced gradient-flow dynamics, using the fact that under this initialization, we get

$$\mathbf{b}_{+,0} = b_{+,0}\mathbf{1}_{M-1}, \quad \mathbf{b}_{-,0} = b_{-,0}\mathbf{1}_{M-1},$$

which implies

$$\mathbf{a}_{+,t} = a_{+,t}\mathbf{1}_{M-1}, \quad \mathbf{a}_{-,t} = a_{-,t}\mathbf{1}_{M-1} \quad \text{for all } t \geq 0.$$

Consequently, the dynamics further reduce from vector-valued coordinates to four scalar variables. This is the gradient-flow analogue of the symmetry reduction used in the regularization-path analysis, where the optimizer is forced by the same permutation symmetry to have block-constant majority and minority coordinates. See Thm. 9.

This block-constant initialization also makes the phase structure visible at the level of gradient-flow time. Although gradient flow is not parametrized by the regularization strength λ , the same mode-wise ordering observed along the regularization path appears empirically in the trajectory: the majority coordinates grow first, followed by the majority–minority coordinate, and finally the minority coordinates. Thus, under the all-ones spectral initialization, the regularization-path phase transition has a clear dynamical analogue.

We verify this invariance numerically in Fig. 9. Starting from block-constant Hadamard coordinates, the entries of each vector block remain indistinguishable throughout training, up to numerical precision. Equivalently, the within-block variance of both $\mathbf{a}_{+,t}$ and $\mathbf{a}_{-,t}$ stays at machine precision along the trajectory.