

# Provable Generalization of Dataset Condensation for Classification via Signal–Noise Dynamics

Anonymous authors

Paper under double-blind review

## Abstract

Dataset condensation, particularly via gradient matching, distills massive datasets into compact synthetic sets, making it an important tool for training under severe storage or computation constraints. However, despite strong empirical performance on classification tasks, existing theory largely relies on regression surrogates or static analyses and gives limited explanation of the underlying classification dynamics. We study gradient-matching condensation for regularized hinge-loss SVMs under an additive sub-Gaussian classification model. Our analysis shows that the learned condensed samples act as signal-concentrating representatives: they aggregate class-level structure while suppressing finite-sample noise and initialization residuals. This mechanism leads to population generalization guarantees for geometry-based evaluators and yields an explicit advantage over random one-shot coresets. The dynamics also identify an early-stopping tradeoff: informative structure is encoded early, whereas overly long inner loops can weaken certified signal accumulation. We further give a multiclass one-condensed-sample-per-class extension through a classwise one-vs-rest update and nearest-prototype evaluation, and simulations on synthetic data and KMNIST corroborate the predicted geometry, schedule sensitivity, and multiclass behavior.

## 1 Introduction

Machine learning is entering a new data-driven era as datasets grow at an unprecedented pace. In deep learning, large-scale datasets such as ImageNet Deng et al. (2009) and multimodal datasets like LAION-5B Schuhmann et al. (2022) and Infinity-MM Gu et al. (2024) have become indispensable for driving technological progress, but storing and training on such data incurs substantial computational and infrastructure costs that hinder deployment in resource-constrained settings.

Dataset distillation and condensation aim to mitigate this bottleneck by synthesizing compact datasets whose training behavior approximates that of the full data Wang et al. (2018); Sucholutsky & Schonlau (2021); Bohdal et al. (2020); Zhao et al. (2021); Zhao & Bilen (2021); Lee et al. (2022); Jiang et al. (2023b); Kim et al. (2022). In particular, gradient-matching condensation Zhao et al. (2021), which aligns the optimization trajectories of synthetic and real data, has achieved remarkable empirical success even under extreme compression ratios. However, a significant disparity remains between this empirical progress and our theoretical understanding of its underlying mechanisms.

Existing theoretical frameworks predominantly rely on regression surrogates (e.g., Kernel Ridge Regression) or infinite-width limits Izzo & Zou (2023); Maalouf et al. (2023); Chen et al. (2024). These analyses often overlook the distinct complexities of classification dynamics, which are inherently non-smooth and state-dependent. Consequently, mechanism-level guarantees for high-dimensional classification—especially for margin-based learners that rely on geometric decision boundaries rather than algebraic fitting—remain scarce. This leaves a fundamental question unanswered: *How do condensed datasets consistently retain the generalization ability required for classification tasks despite the loss of sample diversity?*

We address this gap by conducting a rigorous dynamical analysis of a practical bilevel gradient-matching procedure for regularized hinge-loss SVMs. Departing from static analysis, we explicitly track the optimization trajectory under a standard additive sub-Gaussian noise model. Our primary contribution is a structural

characterization of the terminal condensed samples in the one-sample-per-class regime. We establish that a terminal condensed sample is not merely a selected training point, but a synthesized representation: it emerges as a nonnegative structured combination of signed training samples with deterministic coefficient bounds. This derivation yields a transparent signal–noise decomposition, quantifying how the algorithm actively aggregates class-wise signal while averaging out noise, and elucidates the critical role of the inner-loop schedule in certifying this signal accumulation.

The choice of hinge-loss SVMs is deliberate. Dataset distillation and condensation are often evaluated through a two-stage protocol: synthetic data are first optimized under a chosen condensation model or matching objective, and the learned set is then reused for downstream training or evaluation (Wang et al., 2018; Zhao et al., 2021; Zhao & Bilén, 2021; 2023). Our setting follows this separation: the hinge-loss SVM supplies tractable margin-based gradient signals during condensation, while the terminal condensed set is evaluated through geometry-based downstream rules. Although the SVM decision rule is linear, its gradient-matching dynamics are state-dependent through hinge activity switching, capturing a canonical non-smooth feature of margin-based classification while remaining analytically tractable.

To convert the structure of a minimal condensed set into test performance, we adopt a directional evaluation viewpoint. In this setting, a broad family of symmetric evaluators depends primarily on the direction induced by the two condensed samples, allowing population risk to be controlled by how strongly that direction captures the underlying signal relative to residual components. Instantiating this geometric control for both a random one-shot coreset and the learned condensed set yields a provable generalization guarantee for the terminal output and reveals an explicit class-averaging advantage of condensation over random selection; moreover, the same analysis identifies a dynamical tradeoff whereby overly long inner loops can weaken conservative certified bounds, motivating early stopping. We also give a direct  $K$ -class extension for the one-condensed-sample-per-class regime: a classwise OvR gradient-matching update preserves the coefficient-expansion mechanism class by class, and a nearest-prototype evaluator converts the resulting classwise signal aggregation into a pairwise multiclass risk bound.

Our contributions are summarized as follows.

1. We identify a signal-aggregation mechanism in gradient-matching condensation for hinge-loss SVMs. In the one-sample-per-class regime, the learned condensed samples concentrate class signal while averaging finite-sample noise, and our coefficient bounds quantify how this mechanism depends on the training schedule.
2. We turn condensed-set geometry into population generalization guarantees. For directional prototype evaluators, the analysis yields an explicit class-averaging improvement over random one-shot selection; the same evaluation principle is further extended to a  $K$ -class nearest-prototype rule with one condensed sample per class.
3. We explain schedule sensitivity through the dynamics and validate the resulting predictions. The simulations confirm signal alignment, early-stopping behavior, cross-model transfer, and the corresponding classwise aggregation effect in the multiclass OvR extension.

To contextualize our contributions, Table 1 compares our work with representative theoretical studies on distillation/condensation. Our core dynamical analysis is developed for gradient-matching condensation with hinge-loss SVMs in the binary setting, while the evaluation viewpoint applies more broadly to downstream rules governed by condensed-set geometry. Appendix B gives the corresponding  $K$ -class one-condensed-sample-per-class extension through a classwise OvR update and a nearest-prototype multiclass risk bound.

## 1.1 Related Works

### 1.1.1 Dataset Condensation and Distillation

Pioneered by Wang et al. (2018), dataset distillation aims to compress a large dataset into a small synthetic one that preserves training utility. To improve scalability, Zhao et al. (2021) proposed *dataset condensation*, which optimizes synthetic samples by matching the gradients of models trained on real versus synthetic

Table 1: Comparison of representative theoretical studies on dataset distillation/condensation. For works that provide method-agnostic laws without committing to a specific learner, the model entry is marked as “–”. (GLM: Generalized Linear Model, KRR: Kernel Ridge Regression, SVM: Support Vector Machine)

Work	Task Type	Model	Key Focus	Dynamics?	Noise Model
Izzo et al. Izzo & Zou (2023)	Regression	GLM / KRR	Exact Recovery	✗	✗
Maalouf et al. Maalouf et al. (2023)	Regression	KRR	Size–Error Bound	✗	✗
Chen et al. Chen et al. (2024)	Regression	KRR	Recovery Conditions	✗	✗
Luo and Xu Luo & Xu (2025)	General	–	Utility Boundary	✓	✗
<b>Our Work</b>	<b>Classification</b>	<b>SVM</b>	<b>Signal Learning</b>	<b>✓</b>	<b>✓</b>

data. This gradient-matching paradigm significantly accelerates optimization compared to back-propagating through unrolled trajectories and has been further refined via differentiable augmentations (Zhao & Bilen, 2021) and multi-level optimization strategies (Jiang et al., 2023a). Beyond single-step gradient alignment, recent works have explored alternative matching objectives. Cazenavette et al. (2022) introduced *trajectory matching*, which aligns the long-horizon training paths of student networks to capture temporal training dynamics. Conversely, to bypass the computational cost of bi-level optimization entirely, Zhao & Bilen (2023) proposed *distribution matching*, which directly aligns the feature distributions of synthetic and real data. While these methods offer different trade-offs in computational cost and performance, our work specifically focuses on establishing a theoretical foundation for the widely used gradient-matching framework.

Despite substantial empirical progress in dataset condensation, theoretical research on the topic remains comparatively limited and fragmented, especially for classification-oriented condensation. Existing analyses largely concentrate on regression and kernelized settings, providing recovery or approximation guarantees for generalized linear models and kernel ridge regression Izzo & Zou (2023); Maalouf et al. (2023); Chen et al. (2024). Recently, Aoyama et al. proposed a duality-gap-based kernel distillation framework that extends kernel inducing point methods beyond squared loss (e.g., to hinge and cross-entropy) and provides bounds on parameter deviation and the resulting prediction/test errors Aoyama et al. (2025). Luo and Xu developed a configuration–dynamics–error framework that yields a scaling law in distilled set size for a fixed configuration and a configuration–coverage law characterizing how the required distilled size grows with configuration diversity Luo & Xu (2025). Nevertheless, a mechanism-level understanding of when and why condensation preserves classification generalization under explicit noise models—particularly for margin-based learners—remains insufficiently characterized. Our work bridges this gap through a theoretical framework that jointly characterizes optimization dynamics and generalization guarantees in classification-oriented dataset condensation.

### 1.1.2 Benign Overfitting and High-Dimensional Classification

The phenomenon of benign overfitting has been extensively analyzed in high-dimensional linear classification. Seminal works have characterized the generalization error of maximum margin classifiers (SVMs) under Gaussian mixture models Wang & Thrampoulidis (2020); Chatterji & Long (2021) and extended these results to multiclass settings Wang et al. (2021). Subsequent research has broadened this scope to accommodate label noise Wen et al. (2023); Frei et al. (2022) and heavy-tailed distributions Okudo & Kobayashi (2024). While our work focuses on dataset condensation rather than overfitting per se, we adopt similar high-dimensional

analysis techniques and additive noise models to characterize the signal-to-noise dynamics of the learned condensed samples.

## 2 Preliminaries

A complete notation table is provided in Appendix A (Table 2). This section introduces the data generation model, the SVM objective, the gradient-matching condensation algorithm, the directional evaluation framework, the key assumptions, and the multiclass extension studied in Appendix B.

### 2.1 Data generation model

We consider an additive class-conditional model in  $\mathbb{R}^d$ . The label space is  $\mathcal{Y} = \{-1, +1\}$  for binary classification and  $\mathcal{Y} = [K] := \{1, \dots, K\}$  for  $K$ -class classification. A sample  $(\mathbf{x}, y)$  is generated as

$$\mathbf{x} = \boldsymbol{\mu}_y + \boldsymbol{\xi}, \quad (1)$$

where  $\boldsymbol{\mu}_y \in \mathbb{R}^d$  is the class-dependent signal and  $\boldsymbol{\xi} \in \mathbb{R}^d$  is mean-zero noise independent of  $y$ .

The noise is sub-Gaussian with covariance  $\boldsymbol{\Sigma}_\xi \succeq \mathbf{0}$ . Concretely, we write  $\boldsymbol{\xi} = \boldsymbol{\Sigma}_\xi^{1/2} \boldsymbol{\zeta}$ , where  $\boldsymbol{\zeta} \in \mathbb{R}^d$  has independent, mean-zero, sub-Gaussian coordinates with unit variance. In particular,  $\text{Cov}(\boldsymbol{\xi}) = \boldsymbol{\Sigma}_\xi$  (see Definition C.1).

Let  $\mathcal{D}$  denote the population law of  $(\mathbf{x}, y)$  induced by equation 1. We observe an i.i.d. training set  $\mathcal{T} = \{(\mathbf{x}_i^{\mathcal{T}}, y_i^{\mathcal{T}})\}_{i=1}^{n_{\mathcal{T}}}$ , where each  $(\mathbf{x}_i^{\mathcal{T}}, y_i^{\mathcal{T}})$  is independently sampled from  $\mathcal{D}$ . For  $k \in \mathcal{Y}$ , define the index set and sample size  $\mathcal{I}_k(\mathcal{T}) = \{i \in [n_{\mathcal{T}}] : y_i^{\mathcal{T}} = k\}$  and  $n_{\mathcal{T},k} = |\mathcal{I}_k(\mathcal{T})|$ . In the binary case, we set  $\boldsymbol{\mu}_{+1} = \boldsymbol{\mu}$  and  $\boldsymbol{\mu}_{-1} = -\boldsymbol{\mu}$  for some  $\boldsymbol{\mu} \neq \mathbf{0}$ , so that  $\mathbf{x} = y \boldsymbol{\mu} + \boldsymbol{\xi}$  is a special case of equation 1.

### 2.2 Support Vector Machines

**Binary SVM.** For  $y \in \{-1, +1\}$  we study the regularized hinge-loss objective

$$L_{\mathcal{T}}(\mathbf{w}) := \frac{1}{n_{\mathcal{T}}} \sum_{i=1}^{n_{\mathcal{T}}} (1 - y_i^{\mathcal{T}} \mathbf{w}^{\top} \mathbf{x}_i^{\mathcal{T}})_+ + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (2)$$

where  $\lambda > 0$  and  $(a)_+ := \max\{a, 0\}$ . We omit an explicit bias term for simplicity (it can be incorporated by augmenting  $\mathbf{x}$  with a constant coordinate).

Define the margin  $u_i^{\mathcal{T}}(\mathbf{w}) := y_i^{\mathcal{T}} \mathbf{w}^{\top} \mathbf{x}_i^{\mathcal{T}}$ . At the hinge boundary we fix the deterministic subgradient selection

$$q_i^{\mathcal{T}}(\mathbf{w}) := \mathbb{1}\{u_i^{\mathcal{T}}(\mathbf{w}) < 1\}. \quad (3)$$

This convention agrees with the ordinary hinge derivative away from  $u_i^{\mathcal{T}}(\mathbf{w}) = 1$  and assigns the boundary case to the inactive side. A valid subgradient of equation 2 is therefore

$$\nabla_{\mathbf{w}} L_{\mathcal{T}}(\mathbf{w}) \in \lambda \mathbf{w} - \frac{1}{n_{\mathcal{T}}} \sum_{i=1}^{n_{\mathcal{T}}} q_i^{\mathcal{T}}(\mathbf{w}) y_i^{\mathcal{T}} \mathbf{x}_i^{\mathcal{T}}. \quad (4)$$

Thus the SVM gradient is governed by the current active set of margin-violating samples. This active-set dependence makes the dynamics state-dependent, while within each fixed active-set region the update is affine. We define  $L_{\mathcal{S}}(\mathbf{w})$  and  $\nabla_{\mathbf{w}} L_{\mathcal{S}}(\mathbf{w})$  analogously for a synthetic dataset  $\mathcal{S}$ . We also write the signed sample  $\mathbf{z}_i^{\mathcal{T}} := y_i^{\mathcal{T}} \mathbf{x}_i^{\mathcal{T}}$ , so that the data-dependent term in equation 4 becomes  $\frac{1}{n_{\mathcal{T}}} \sum_i q_i^{\mathcal{T}}(\mathbf{w}) \mathbf{z}_i^{\mathcal{T}}$ .

### 2.3 Gradient-Matching Dataset Condensation Algorithm

We specify the dataset condensation (DC) algorithm analyzed in the rest of the paper. The procedure is a hinge-loss SVM specialization of gradient matching Zhao et al. (2021). It compresses the training set  $\mathcal{T}$  into

a synthetic set  $\mathcal{S}$  by making the data-dependent SVM subgradient computed on  $\mathcal{S}$  mimic the corresponding subgradient computed on  $\mathcal{T}$  along repeatedly initialized training trajectories.

Algorithm 1 implements the *condensation stage* of our pipeline. The SVM supplies a low-cost proxy gradient that shapes the synthetic samples; after condensation, the learned set  $\mathcal{S}_{\text{dc}}$  may be passed unchanged to any downstream evaluator whose decision rule depends on the condensed-set geometry described in Subsection 2.4.

The main analysis focuses on the binary one-sample-per-class regime  $\mathcal{S} = \{(\mathbf{x}_+^{\mathcal{S}}, +1), (\mathbf{x}_-^{\mathcal{S}}, -1)\}$  with  $n_{\mathcal{S}} = 2$ . We write the algorithm in signed coordinates  $\mathbf{z}_y^{\mathcal{S}} := y \mathbf{x}_y^{\mathcal{S}}$  for  $y \in \{\pm 1\}$ ; under the symmetric binary model, both class signals then equal the common direction  $\boldsymbol{\mu}$ . For a class block  $\mathcal{A}_y$  of a dataset  $\mathcal{A}$ , define the activity-weighted signed mean at the current model state by

$$\mathbf{g}_{\mathcal{A}_y}(\mathbf{w}) := \frac{1}{n_{\mathcal{A}_y}} \sum_{i \in \mathcal{I}_y(\mathcal{A})} q_i^{\mathcal{A}}(\mathbf{w}) \mathbf{z}_i^{\mathcal{A}}, \quad (5)$$

with the convention  $\mathbf{g}_{\mathcal{S}_y}(\mathbf{w}) = q_{\mathcal{S},y}(\mathbf{w}) \mathbf{z}_y^{\mathcal{S}}$  for the singleton condensed class. The regularization term  $\lambda \mathbf{w}$  is common to the real and synthetic SVM subgradients, so the classwise matching step only needs to align the data-dependent terms in equation 5. The gradient-matching discrepancy for class  $y$  is therefore

$$D_y(\mathcal{S}; \mathbf{w}) := \|\mathbf{g}_{\mathcal{S}_y}(\mathbf{w}) - \mathbf{g}_{\mathcal{T}_y}(\mathbf{w})\|_2^2. \quad (6)$$

The algorithm has two nested loops. At the beginning of every outer loop  $t_{\text{out}}$ , the SVM weight vector is reinitialized from a small isotropic Gaussian distribution. The condensed samples are not reinitialized; they are inherited from the end of the previous outer loop. Within the inner loop, the algorithm alternates between a classwise synthetic-data update and an SVM weight update. Throughout the synthetic-data update, the hard hinge activity indicators are evaluated at the current iterate and their derivatives with respect to the synthetic samples are taken to be zero, which agrees with their almost-everywhere derivative and fixes the convention at hinge switching boundaries. The resulting classwise gradient step induced by equation 6 takes the explicit form

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}}+1)} = \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} + 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \left( \mathbf{g}_{\mathcal{T}_y}^{(t_{\text{out}}, t_{\text{in}})} - \mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}}, t_{\text{in}})} \right). \quad (7)$$

Thus, when the condensed sample is hinge-active, it is pulled toward the activity-weighted signed mean of the real samples in the same class; when it is inactive, the corresponding classwise synthetic update vanishes. The SVM weights are then updated on the real training objective:

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}}+1)} = (1 - \eta_{\mathbf{w}} \lambda) \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} + \eta_{\mathbf{w}} \mathbf{g}_{\mathcal{T}}^{(t_{\text{out}}, t_{\text{in}})}, \quad \mathbf{g}_{\mathcal{T}}^{(t_{\text{out}}, t_{\text{in}})} := \frac{1}{n_{\mathcal{T}}} \sum_{i=1}^{n_{\mathcal{T}}} q_i^{\mathcal{T},(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_i^{\mathcal{T}}. \quad (8)$$

This explicit form explains why the terminal condensed samples admit a tractable signal-noise expansion. During hinge-active periods, equation 7 is an affine contraction toward the classwise signed mean of the training data, with contraction factor  $\rho_{\mathcal{S}} = 1 - 2\eta_{\mathcal{S}}$ . Across outer loops, repeated inheritance accumulates classwise signal while averaging finite-sample noise. The main results below quantify this mechanism by unrolling equation 7 over the bilevel schedule and translating the resulting condensed-set geometry into population risk bounds.

## 2.4 Prototypical Evaluation Framework

Algorithm 1 describes the condensation stage: the SVM provides a tractable margin-based gradient signal for updating the synthetic data, while the learned set  $\mathcal{S}_{\text{dc}}$  is later evaluated through downstream rules determined by its geometry.

For a binary condensed set with one sample per class, the relevant geometric object is the direction connecting the two condensed samples. The following definition records the class of evaluators whose population risk is governed by this direction.

**Algorithm 1** Gradient-Matching DC for Hinge-Loss SVM (One Sample per Class)

---

**Input:** Training set  $\mathcal{T}$ , loop lengths  $T_{\text{out}}$  and  $T_{\text{in}}$ , step sizes  $\eta_{\mathbf{w}}$  and  $\eta_{\mathcal{S}}$   
**Output:** Condensed set  $\mathcal{S}_{\text{dc}}$   
Initialize signed condensed samples  $\mathbf{z}_y^{\mathcal{S},(0,0)} \sim \mathcal{N}(0, \sigma_{\mathcal{S}}^2 \mathbf{I}_d)$  for  $y \in \{\pm 1\}$   
**for**  $t_{\text{out}} = 0$  **to**  $T_{\text{out}} - 1$  **do**  
  Initialize  $\mathbf{w}^{(t_{\text{out}},0)} \sim \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \mathbf{I}_d)$   
  **if**  $t_{\text{out}} > 0$  **then**  
    Set  $\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},0)} \leftarrow \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}-1, T_{\text{in}})}$  for each  $y \in \{\pm 1\}$   
  **end if**  
  **for**  $t_{\text{in}} = 0$  **to**  $T_{\text{in}} - 1$  **do**  
    Compute the hinge indicators  $q_i^{\mathcal{T},(t_{\text{out}},t_{\text{in}})}$  and  $q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}$   
    Compute  $\mathbf{g}_{\mathcal{T}}^{(t_{\text{out}},t_{\text{in}})}$ ,  $\mathbf{g}_{\mathcal{T}_y}^{(t_{\text{out}},t_{\text{in}})}$ , and  $\mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}},t_{\text{in}})}$  from equation 5  
    **for**  $y \in \{\pm 1\}$  **do**  
      Update  $\mathbf{z}_y^{\mathcal{S}}$  by equation 7  
    **end for**  
    Update  $\mathbf{w}$  by equation 8  
  **end for**  
**end for**  
**Return:**  $\mathcal{S}_{\text{dc}} = \{(y\mathbf{z}_y^{\mathcal{S},(T_{\text{out}}-1, T_{\text{in}})}, y) : y \in \{\pm 1\}\}$

---

**Definition 2.1** (Directional prototype evaluator). Let  $\mathcal{S} = \{(\mathbf{x}_+, +1), (\mathbf{x}_-, -1)\}$  be a binary condensed dataset with one sample per class. Define the *two-sample difference vector*  $\mathbf{s}(\mathcal{S}) := \mathbf{x}_+ - \mathbf{x}_-$ . A learning algorithm is a *directional prototype evaluator* if:

1. **Directional Alignment:** The learned weight vector aligns with the two-sample difference, i.e.,  $\hat{\mathbf{w}}_{\mathcal{S}} := \mathbf{w}_{\mathcal{S}} / \|\mathbf{w}_{\mathcal{S}}\|_2 = \mathbf{s}(\mathcal{S}) / \|\mathbf{s}(\mathcal{S})\|_2$ .
2. **Alignment-Based Risk:** The generalization risk on distribution  $\mathcal{D}$  depends solely on this direction:  $\mathcal{R}(\mathcal{S}) := \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}(y \langle \hat{\mathbf{w}}_{\mathcal{S}}, \mathbf{x} \rangle < 0)$ .

In the one-sample-per-class setting, several standard evaluation rules reduce to a boundary whose normal direction is  $\mathbf{s}(\mathcal{S})$ . The following examples illustrate common cases:

- **SVM and logistic regression.** With one sample per class, symmetric regularized linear objectives include cases whose fitted direction is parallel to  $\mathbf{x}_+ - \mathbf{x}_-$ ; in particular, a linear SVM has maximum-margin normal direction  $\mathbf{x}_+ - \mathbf{x}_-$  up to scaling.
- **Nearest centroid classifier.** The distance rule  $\|\mathbf{x} - \mathbf{x}_+\|_2 < \|\mathbf{x} - \mathbf{x}_-\|_2$  is algebraically equivalent to the linear comparison  $\langle \mathbf{x}, \mathbf{x}_+ - \mathbf{x}_- \rangle > C$ , where  $C = (\|\mathbf{x}_+\|_2^2 - \|\mathbf{x}_-\|_2^2)/2$ .
- **Feature-space nearest-centroid rules.** In representation spaces where neural-collapse geometry is observed (Papayan et al., 2020), class features concentrate around their means and the classifier behaves analogously to a nearest-centroid rule in feature space.

Our risk bounds therefore apply to any downstream rule satisfying Definition 2.1; the examples above identify familiar settings where this directional condition is natural.

## 2.5 Assumptions

This subsection collects the assumptions used in the binary analysis. They specify the signal scale, the concentration regime, the stable step-size range, and the initialization scale.

**Assumption 2.2.** Fix a tail parameter  $\delta \in (0, 1)$  and constant  $\kappa > 0$ . Let  $C_{\kappa} > 0$  denote the concentration constant determined by the choice of  $\kappa$ . We condition our analysis on the intersection of the high-probability

events defined in Lemmas C.2–C.5 (Section C.1), ensuring that the total failure probability is at most  $\delta$ . We assume the following conditions hold:

**A1. Signal strength.** The signal-to-noise ratio satisfies

$$\frac{\|\boldsymbol{\mu}\|_2}{\sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \sqrt{d}}} \geq 2C_\kappa.$$

**A2. High dimension.** The input dimension  $d$  is sufficiently large relative to the number of outer loops  $T_{\text{out}}$ , the condensed-set size  $n_S$ , the training-set size  $n_T$ , and the failure probability  $\delta$ :

$$d \geq \kappa \log\left(\frac{24 T_{\text{out}} n_S n_T}{\delta}\right).$$

**A3. Small learning rate.** The step sizes are chosen to ensure stability, where  $\eta_{\mathbf{w}}$  and  $\eta_S$  denote the inner-loop step sizes for updating  $\mathbf{w}$  and  $S$ , respectively:

$$0 < \eta_{\mathbf{w}} \lambda \leq 1, \quad 0 < 2\eta_S \leq 1.$$

**A4. Small initialization.** Define  $\rho_S := 1 - 2\eta_S$ . The weight initialization satisfies

$$\sigma_{\mathbf{w}} < \frac{1}{(2C_\kappa + 1)\|\boldsymbol{\mu}\|_2 \sqrt{d}}.$$

The synthetic initialization satisfies

$$\sigma_S < \min\left\{\frac{1}{C_\kappa \sigma_{\mathbf{w}} d}, \frac{\lambda}{(2C_\kappa + 1)\|\boldsymbol{\mu}\|_2 \sqrt{d}}, \rho_S^{T_{\text{in}}} \sqrt{\frac{2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}{n_T}}\right\}.$$

Assumptions A1 and A2 are the high-dimensional signal and concentration conditions used to separate the class signal from the sub-Gaussian noise Chatterji & Long (2021); Cao et al. (2022); Kou et al. (2023); Okudo & Kobayashi (2024). Assumption A3 keeps the weight and condensed-sample updates in the stable contraction regime. Assumption A4 imposes signal-scale small random initialization on both the model weights and the synthetic samples, so that the initial hinge margins lie in the active region and the initialization residual remains controlled by the noise scale.

## 2.6 Multiclass Classwise OvR Extension

The main results below focus on the binary one-sample-per-class setting, where signed coordinates reduce the two condensed samples to a common signal direction. Appendix B introduces the corresponding  $K$ -class setup. In that extension, the condensed object is a single multiclass set  $\mathcal{P} = \{(\mathbf{p}_k, k) : k \in [K]\}$  with one condensed sample per class, rather than  $K$  separate two-point condensed sets obtained from independent one-vs-rest reductions.

The multiclass condensation step uses  $K$  binary OvR heads. For head  $h$ , class  $h$  is positive and every other class is negative, with sign  $\tau_{h,k} = 2\mathbb{1}\{h = k\} - 1$ . The class- $k$  condensed sample is updated by matching, averaged over heads, its OvR signed gradient term against the same-class real-data OvR gradient term. This defines a classwise OvR gradient-matching objective while keeping a single shared multiclass condensed set.

The corresponding multiclass evaluator is the nearest-prototype rule induced by the condensed set,

$$h_{\mathcal{P}}(\mathbf{x}) := \arg \min_{k \in [K]} \|\mathbf{x} - \mathbf{p}_k\|_2^2 = \arg \max_{k \in [K]} \left\{ \langle \mathbf{x}, \mathbf{p}_k \rangle - \frac{1}{2} \|\mathbf{p}_k\|_2^2 \right\}.$$

Its population risk is  $\mathcal{R}_{\text{mc}}(\mathcal{P}) := \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}(h_{\mathcal{P}}(\mathbf{x}) \neq y)$ . This evaluator is the multiclass counterpart of the two-sample geometric evaluator used in the binary analysis: it compares the test point with one representative per class and includes the class-dependent quadratic bias in the equivalent score form. The detailed guarantees for this classwise OvR construction are stated after the binary results and in Appendix B.

### 3 Main Results

This section states the main guarantees. The binary one-sample-per-class result is the technical core: the terminal condensed samples achieve a population risk bound controlled by a certified signal-to-residual ratio. We then give the two ingredients behind this theorem: a dynamical expansion showing that the terminal condensed samples are nonnegative combinations of training samples (Section 3.2), and a geometric reduction converting that expansion into signal-to-residual lower bounds (Section 3.3). The corresponding  $K$ -class statement is summarized at the end of the section and stated in full in Appendix B.

#### 3.1 Provable generalization for gradient-matching condensation

The main theorem shows that the terminal output of the gradient-matching procedure retains enough signal to generalize under the population model. Let  $\mathbf{t}^* := (T_{\text{out}} - 1, T_{\text{in}})$  and write  $\mathbf{z}_y^S := \mathbf{z}_y^{S,(\mathbf{t}^*)}$  for the terminal signed condensed sample produced by Algorithm 1 for class  $y \in \{\pm 1\}$ . Since our dynamics are written in signed coordinates  $\mathbf{z} = y\mathbf{x}$ , the terminal condensed dataset in the original input space is

$$\mathcal{S}_{\text{dc}} := \{(y\mathbf{z}_y^S, y) : y \in \{\pm 1\}\} = \{(\mathbf{z}_+^S, +1), (-\mathbf{z}_-^S, -1)\}. \quad (9)$$

The bound applies to the directional evaluator induced by this two-point output. Its exponent is governed by the certified signal-to-residual ratio  $\underline{\text{SR}}_{\text{dc}}$ , derived in Lemma 3.7.

**Theorem 3.1** (Provable generalization for the terminal condensed samples). *Assume the additive model in equation 1 and Assumption 2.2. Run Algorithm 1 with one signed condensed sample per class, and let  $\mathcal{S}_{\text{dc}}$  be its terminal output defined in equation 9. Then, with probability at least  $1 - \delta$ , the population risk of the induced directional evaluator satisfies*

$$\mathcal{R}(\mathcal{S}_{\text{dc}}) \leq \exp\left(-\frac{\|\boldsymbol{\mu}\|_2^2}{2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}\left(\left(\frac{\underline{\text{SR}}_{\text{dc}} - 1}{\underline{\text{SR}}_{\text{dc}} + 1}\right)_+\right)^2\right), \quad (10)$$

where  $(u)_+ := \max\{u, 0\}$  and  $\underline{\text{SR}}_{\text{dc}}$  is the certified lower bound on the signal-to-residual ratio specified in equation 20.

Under the signal-strength condition in Assumption A1, the certified ratio satisfies

$$\underline{\text{SR}}_{\text{dc}} \geq 2\sqrt{C_\kappa} L\sqrt{n\mathcal{T}}, \quad (11)$$

where  $L$  is the dimensionless factor defined in Lemma 3.7. It summarizes the coefficient mass accumulated during certified classwise-averaging windows, the effect of class balance, and the contribution of coefficient mass outside those windows. The same signal-strength condition yields a direct dimension-dependent form of the population-risk exponent. Since  $\|\boldsymbol{\mu}\|_2^2/(2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}) \geq 2C_\kappa^2 d$ , equation 10 implies

$$\mathcal{R}(\mathcal{S}_{\text{dc}}) \leq \exp\left(-2C_\kappa^2 d \left(\left(\frac{\underline{\text{SR}}_{\text{dc}} - 1}{\underline{\text{SR}}_{\text{dc}} + 1}\right)_+\right)^2\right). \quad (12)$$

Thus the guarantee is nontrivial whenever  $\underline{\text{SR}}_{\text{dc}} > 1$ , and the exponent scales linearly with the ambient dimension when this ratio is bounded away from one. Section 3.3 explains the two components of the certified ratio: the prefactor  $\sqrt{n\mathcal{T}}$  is the class-averaging scale, while  $L$  records how much of that scale remains after accounting for coefficient mass outside the certified averaging windows.

#### 3.2 Dynamical structure: Signal–noise expansion

To prove Theorem 3.1, we first characterize the terminal condensed samples produced by the bilevel dynamics. We write  $\mathbf{t}^* := (T_{\text{out}} - 1, T_{\text{in}})$  for the terminal time index. The lemma below connects the gradient-matching recursion to condensed-set geometry. Part (i) gives *linear-window* lengths, namely initial inner prefixes on which the training-set and condensed-set hinge indicators remain active under the high-probability controls

used in the analysis. On these prefixes the update is affine and can be unrolled. Parts (ii)–(iii) then express each terminal signed condensed sample as a nonnegative combination of initialization and training data, with coefficient bounds determined by  $\rho_S$  and by the joint window length  $T_{y,\text{lin}}(r) = \min\{T_{\mathcal{T},\text{lin}}, T_{S,\text{lin}}(r)\}$ .

The condensed-set window  $T_{S,\text{lin}}(t_{\text{out}})$  is indexed by the outer round. Starting from the reset weights at  $(t_{\text{out}}, 0)$ , it is the length of the initial inner prefix on which both condensed samples remain strictly inside the  $\pm 1$  margin band, using the hinge indicators in equation 3. This mirrors the training-set window  $T_{\mathcal{T},\text{lin}}$ . Their minimum is the interval on which a classwise condensed update is certified to be an affine averaging step, and this is the interval that controls the coefficient lower bounds.

**Lemma 3.2** (Terminal condensed-sample expansion). *Assume the additive model in equation 1 and Assumption 2.2, and consider the one-sample-per-class output of Algorithm 1. Set  $\rho_S := 1 - 2\eta_S$ . Fix  $y \in \{\pm 1\}$ . Recall that we work with signed samples  $\mathbf{z} := y\mathbf{x}$ , and the synthetic initialization satisfies  $\mathbf{z}_y^{S,(0,0)} \sim \mathcal{N}(0, \sigma_S^2 \mathbf{I}_d)$ .*

(i) *There exist  $(T_{\mathcal{T},\text{lin}}, \{T_{S,\text{lin}}(t_{\text{out}})\}_{t_{\text{out}}=0}^{T_{\text{out}}-1}) \in \{0, \dots, T_{\text{in}}\} \times \{0, \dots, T_{\text{in}}\}^{T_{\text{out}}}$  such that  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t)}) = 1$  for all  $i \in \mathcal{I}(\mathcal{T})$  and all  $t \leq T_{\mathcal{T},\text{lin}}$ , and  $q_{S,y}^{(t_{\text{out}},t)} = 1$  for all  $t \leq T_{S,\text{lin}}(t_{\text{out}})$ .*

(ii) *There exist coefficients  $\tilde{u}_y^{(\mathbf{t}^*)} \geq 0$  and  $\{\tilde{c}_{y,j}^{(\mathbf{t}^*)} \geq 0\}_{j \in \mathcal{I}(\mathcal{T}_y)}$  such that, with  $\tilde{c}_{y,\mu}^{(\mathbf{t}^*)} := \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(\mathbf{t}^*)}$  and  $\tilde{\mathbf{c}}_y^{(\mathbf{t}^*)} := (\tilde{c}_{y,j}^{(\mathbf{t}^*)})_{j \in \mathcal{I}(\mathcal{T}_y)}$ , the terminal signed condensed sample admits the expansion*

$$\mathbf{z}_y^{S,(\mathbf{t}^*)} = \tilde{u}_y^{(\mathbf{t}^*)} \mathbf{z}_y^{S,(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(\mathbf{t}^*)} \mathbf{z}_j^{\mathcal{T}}. \quad (13)$$

(iii) *The initialization coefficient satisfies*

$$\rho_S^{T_{\text{out}}T_{\text{in}}} \leq \tilde{u}_y^{(\mathbf{t}^*)} \leq \rho_S^{\sum_{k=0}^{T_{\text{out}}-1} T_{S,\text{lin}}(k)}. \quad (14)$$

Define  $T_{y,\text{lin}}(r) := \min\{T_{\mathcal{T},\text{lin}}, T_{S,\text{lin}}(r)\}$ . There exist deterministic quantities  $B_{\text{out}} \geq 1$ ,  $u_{\text{out}} = (1 - \rho_S^{T_{\text{in}}})B_{\text{out}}$ , and  $\ell_y \in [0, u_{\text{out}}]$ , depending only on  $\rho_S$  and the certified window lengths, such that the following bounds hold. The explicit construction of these quantities is recorded in Lemma D.8 in the appendix. For every  $j \in \mathcal{I}(\mathcal{T}_y)$ ,

$$\frac{\ell_y}{n_{\mathcal{T}_y}} \leq \tilde{c}_{y,j}^{(\mathbf{t}^*)} \leq \frac{u_{\text{out}}}{n_{\mathcal{T}_y}}. \quad (15)$$

Consequently,

$$\ell_y \leq \tilde{c}_{y,\mu}^{(\mathbf{t}^*)} \leq u_{\text{out}}, \quad \frac{\ell_y}{\sqrt{n_{\mathcal{T}_y}}} \leq \|\tilde{\mathbf{c}}_y^{(\mathbf{t}^*)}\|_2 \leq \frac{u_{\text{out}}}{\sqrt{n_{\mathcal{T}_y}}}. \quad (16)$$

(iv) *Moreover, let  $G_{\text{max}}$  denote the high-probability upper bound on pairwise signed-sample inner products, so that  $\langle \mathbf{z}_i^{\mathcal{T}}, \mathbf{z}_{i'}^{\mathcal{T}} \rangle \leq G_{\text{max}}$  for all  $i, i' \in \mathcal{I}(\mathcal{T})$  on the event in Assumption 2.2. If  $G_{\text{max}} < \lambda$ , then  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t_{\text{in}})}) = 1$  and  $q_{S,y}^{(t_{\text{out}},t_{\text{in}})} = 1$  for all indices, and the coefficients simplify to*

$$\tilde{u}_y^{(\mathbf{t}^*)} = \rho_S^{T_{\text{out}}T_{\text{in}}} \quad \tilde{c}_{y,j}^{(\mathbf{t}^*)} = \frac{1 - \rho_S^{T_{\text{out}}T_{\text{in}}}}{n_{\mathcal{T}_y}}. \quad (17)$$

*Remark 3.3.* Part (i) provides initial linear windows for the training and condensed dynamics. On these windows, the hinge activity indicators stay at one throughout the corresponding initial inner loop. Thus the subgradient pattern is fixed on the certified prefix, and the inner-loop map reduces to an affine contraction that can be unrolled explicitly.

*Remark 3.4.* Part (ii) supplies nonnegative coefficients  $\tilde{u}_y^{(\mathbf{t}^*)}$  and  $\{\tilde{c}_{y,j}^{(\mathbf{t}^*)}\}$  for the initialization and the class- $y$  training samples. It specifies how the terminal signed condensed sample can be expressed as a nonnegative combination of the initialization and the signed training points. This is the basic “signal-noise” template:  $\tilde{u}_y^{(\mathbf{t}^*)}$  weights the random initialization, while  $\{\tilde{c}_{y,j}^{(\mathbf{t}^*)}\}$  allocates mass across the real samples in class  $y$ .

*Remark 3.5.* Part (iii) turns the abstract expansion in equation 13 into concrete two-sided bounds on each coefficient. For the initialization weight, the bound  $\rho_S^{T_{\text{out}}T_{\text{in}}} \leq \tilde{u}_y^{(\mathbf{t}^*)} \leq \rho_S^{\sum_k T_{S,\text{lin}}(k)}$  says that the random synthetic draw is attenuated by the active portions of the dynamics: the more time per outer round the condensed sample spends in the hinge-active regime, the more strongly the initialization is contracted. The quantity  $\ell_y$  is the amount of terminal training coefficient mass contributed by steps whose classwise updates are certified to be uniform averages, while  $u_{\text{out}}$  is the corresponding upper envelope for the total terminal training mass. The coordinatewise bound equation 15 separates this averaged contribution from the remaining coefficient mass not covered by that lower bound; this remainder is accounted for in Lemma 3.7 through  $\Delta_y = u_{\text{out}} - \ell_y$ .

*Remark 3.6.* Part (iv) gives a closed-form specialization in the strong-regularization regime  $G_{\text{max}} < \lambda$ : full activity persists throughout training and condensation, the initialization weight decays exponentially, and the training coefficients become uniform within each class, making explicit the averaging mechanism on the classwise noise.

### 3.3 Geometric analysis: From structure to risk

We next translate the structural controls from Lemma 3.2 into the population classification guarantee. The next lemma gives the geometric risk reduction for directional evaluators and lower bounds the signal-to-residual ratio  $\widehat{\text{SR}}(\mathcal{S})$  for both a random one-shot coresets and the terminal condensed set. The latter lower bound is the certified constant  $\underline{\text{SR}}_{\text{dc}}$  used in Theorem 3.1.

**Lemma 3.7** (Directional risk bound and signal-to-residual lower bounds). *Assume the test distribution follows the additive model  $\mathbf{x} = y\boldsymbol{\mu} + \boldsymbol{\xi}$  with  $y \in \{\pm 1\}$  and mean-zero sub-Gaussian noise  $\boldsymbol{\xi}$  with proxy covariance  $\boldsymbol{\Sigma}_\xi$ , and let  $(u)_+ := \max\{u, 0\}$ . For any two-point set  $\mathcal{S} = \{(\mathbf{x}_+, +1), (\mathbf{x}_-, -1)\}$ , define  $\widehat{\mathbf{w}}_{\mathcal{S}} := (\mathbf{x}_+ - \mathbf{x}_-) / \|\mathbf{x}_+ - \mathbf{x}_-\|_2$  and  $\mathcal{R}(\mathcal{S}) := \mathbb{P}(y\langle \widehat{\mathbf{w}}_{\mathcal{S}}, \mathbf{x} \rangle < 0)$ . Suppose  $\mathbf{x}_+ - \mathbf{x}_- = A(\mathcal{S})\boldsymbol{\mu} + \mathbf{r}(\mathcal{S})$  for some  $A(\mathcal{S}) > 0$ , and define  $\widehat{\text{SR}}(\mathcal{S}) := A(\mathcal{S})\|\boldsymbol{\mu}\|_2 / \|\mathbf{r}(\mathcal{S})\|_2$  (with  $\widehat{\text{SR}}(\mathcal{S}) = +\infty$  if  $\mathbf{r}(\mathcal{S}) = \mathbf{0}$ ). Then the following hold.*

(i) *The population risk satisfies*

$$\mathcal{R}(\mathcal{S}) \leq \exp\left(-\frac{\|\boldsymbol{\mu}\|_2^2}{2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}\left(\frac{\widehat{\text{SR}}(\mathcal{S}) - 1}{\widehat{\text{SR}}(\mathcal{S}) + 1}\right)_+^2\right). \quad (18)$$

(ii) *Assume Assumption 2.2 and the training model  $\mathbf{x}_i^T = y_i^T \boldsymbol{\mu} + \boldsymbol{\xi}_i$ . Pick any  $i_+ \in \mathcal{I}(\mathcal{T}_+)$  and  $i_- \in \mathcal{I}(\mathcal{T}_-)$  and set  $\mathcal{S}_{\text{core}} = \{(\mathbf{x}_{i_+}^T, +1), (\mathbf{x}_{i_-}^T, -1)\}$ . On the high-probability event in Assumption 2.2, letting  $C_\kappa$  denote the corresponding constant, we have*

$$\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) \geq \frac{2\|\boldsymbol{\mu}\|_2}{\sqrt{2(1 + 2C_\kappa)}\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}d}. \quad (19)$$

(iii) *Assume the hypotheses of Lemma 3.2 and consider the terminal condensed set  $\mathcal{S}_{\text{dc}}$  in equation 9. Let  $B_{\text{out}}$ ,  $u_{\text{out}}$ , and  $\ell_y$ , for  $y \in \{\pm 1\}$ , be the deterministic coefficient quantities from Lemma 3.2. Set  $p_y := n_{\mathcal{T}_y} / n_{\mathcal{T}}$ ,  $\beta_p := p_+^{-1/2} + p_-^{-1/2}$ , and  $\Delta_y := u_{\text{out}} - \ell_y$ . Define the dimensionless factor*

$$L := \frac{\ell_+ + \ell_-}{B_{\text{out}}\beta_p + \sum_{y \in \{\pm 1\}} \Delta_y (\sqrt{n_{\mathcal{T}}} - p_y^{-1/2})}.$$

*Then, on the same high-probability event,*

$$\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) \geq \underline{\text{SR}}_{\text{dc}} = \frac{\sqrt{n_{\mathcal{T}}}L\|\boldsymbol{\mu}\|_2}{\sqrt{C_\kappa}\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}d}. \quad (20)$$

*Remark 3.8.* The risk upper bound equation 18 is governed by the signal-to-residual ratio  $\widehat{\text{SR}}(\mathcal{S})$ . Larger values yield a more negative exponent and a tighter bound, while  $\widehat{\text{SR}}(\mathcal{S}) \leq 1$  makes the exponent term equal to zero and gives only the vacuous bound  $\mathcal{R}(\mathcal{S}) \leq 1$ .

*Remark 3.9.* Comparing equation 19 and equation 20, the condensed lower bound displays the class-averaging scale through the explicit prefactor  $\sqrt{n_{\mathcal{T}}}$ . The dimensionless factor  $L$  determines how much of that scale is certified by the dynamics. Its numerator  $\ell_+ + \ell_-$  is the coefficient mass that enters through certified classwise averaging, while the denominator combines the class-balance and averaged-residual term  $B_{\text{out}}\beta_p$  with the cost of the non-averaged remainder  $\sum_y \Delta_y(\sqrt{n_{\mathcal{T}}} - p_y^{-1/2})$ .

*Remark 3.10.* The factor  $L$  determines whether the explicit prefactor  $\sqrt{n_{\mathcal{T}}}$  in equation 20 is retained after accounting for activity switching. Lemma D.12 in the appendix quantifies the loss incurred when the certified active prefix is shorter than the full inner loop. With  $T := T_{\text{in}}$ ,  $\rho := \rho_{\mathcal{S}}$ , and

$$M_{\mathcal{S}} := \max_{0 \leq r \leq T_{\text{out}} - 1} \sum_{k=r+1}^{T_{\text{out}} - 1} (T - T_{\mathcal{S}, \text{lin}}(k)),$$

this deficit is

$$\mathcal{W}_{\text{win}} := (1 - \rho^T)(1 - \rho^{M_{\mathcal{S}}}) + \max_{\substack{y \in \{\pm 1\} \\ 0 \leq r \leq T_{\text{out}} - 1}} (1 - \rho^{T - T_{y, \text{lin}}(r)}).$$

The first term measures the effect of later condensed active prefixes that are shorter than  $T$ : if a contribution is injected in an early outer round, then subsequent outer rounds with  $T_{\mathcal{S}, \text{lin}}(k) < T$  leave less certified contraction on that contribution. The quantity  $M_{\mathcal{S}}$  records the largest cumulative future gap  $\sum_k (T - T_{\mathcal{S}, \text{lin}}(k))$ . The second term measures the gap within the current outer round: classwise averaging is certified only up to  $T_{y, \text{lin}}(r)$ , so the remaining  $T - T_{y, \text{lin}}(r)$  inner steps are not treated as averaged updates. Hence  $\mathcal{W}_{\text{win}} = 0$  when all relevant active prefixes have full length  $T$ , while larger values indicate that more terminal coefficient mass may lie outside the certified classwise-averaging part. Lemma D.12 proves the deterministic comparisons between  $\mathcal{W}_{\text{win}}$ , the remainders  $\Delta_y$ , and the factor  $L$ , yielding the following regimes.

- (i) In the full-activity regime  $G_{\text{max}} < \lambda$ , all windows have full length and  $\mathcal{W}_{\text{win}} = 0$ . The terminal training coefficients are exactly uniform within each class, so the non-averaged remainder vanishes. In this case

$$L = \frac{2(1 - \rho^T)}{\beta_p},$$

so the bound retains the full  $\sqrt{n_{\mathcal{T}}}$  class-averaging scale up to the class-balance factor.

- (ii) In the small-remainder regime, if  $\mathcal{W}_{\text{win}} = o(1 - \rho^T)$  and  $\mathcal{W}_{\text{win}} = O(\beta_p/\sqrt{n_{\mathcal{T}}})$ , then the coefficient mass outside the certified averaging part is at most comparable to the empirical averaging scale. Lemma D.12 gives  $L = \Omega((1 - \rho^T)/\beta_p)$ , and if the second condition is strengthened to  $\mathcal{W}_{\text{win}} = o(\beta_p/\sqrt{n_{\mathcal{T}}})$ , then  $L = (1 + o(1))2(1 - \rho^T)/\beta_p$ .
- (iii) In the switching-dominated regime, the total remainder  $\sum_y \Delta_y$  is large enough that the non-averaged part of the residual dominates the averaged-residual term in the denominator of  $L$ . For example, if  $n_{\text{min}} \rightarrow \infty$  and  $\sum_y \Delta_y = \Omega(B_{\text{out}}(1 - \rho^T))$ , Lemma D.12 gives  $L = O(n_{\mathcal{T}}^{-1/2})$ . The proof still yields a valid lower bound, but it no longer certifies a growing  $\sqrt{n_{\mathcal{T}}}$  class-averaging gain.

Theorem 3.1 follows by combining the certified signal-to-residual lower bound for the terminal condensed set with the directional-evaluator risk reduction.

*Proof of Theorem 3.1.* Lemma 3.7(iii) gives  $\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) \geq \underline{\text{SR}}_{\text{dc}}$ , and Lemma 3.7(i) bounds  $\mathcal{R}(\mathcal{S}_{\text{dc}})$  by the right-hand side of equation 18. Since  $t \mapsto \left( \left( \frac{t-1}{t+1} \right)_+ \right)^2$  is nondecreasing on  $[0, +\infty)$ , substituting  $\underline{\text{SR}}_{\text{dc}}$  yields equation 10.  $\square$

### 3.4 Multiclass consequences

The same coefficient-expansion mechanism extends to the  $K$ -class classwise OvR setting introduced in Subsection 2.6. Let  $\mathcal{P}_{\text{dc}} = \{(\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)}, k) : k \in [K]\}$  be the terminal multiclass condensed set. Theorem B.3 shows that, for every class  $k$ , the terminal condensed sample admits a nonnegative same-class expansion  $\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)} = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{\mathcal{S},(0,0)} + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \mathbf{x}_i^{\mathcal{T}}$ , with coefficient bounds controlled by the classwise OvR active windows. In the full-activity specialization, this expansion reduces to the contraction  $\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)} = \rho_S^{T_{\text{out}}T_{\text{in}}} \mathbf{p}_k^{\mathcal{S},(0,0)} + (1 - \rho_S^{T_{\text{out}}T_{\text{in}}}) \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}}$  toward the empirical class mean.

For the nearest-prototype evaluator introduced in Subsection 2.6, Theorem B.4 gives a pairwise multiclass risk bound controlled by the class separations  $\gamma_{k\ell} = \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_\ell\|_2$  and the perturbations of the condensed samples from their class means. Corollary B.5 combines this geometric bound with the OvR coefficient expansion to obtain a certified  $K$ -class risk guarantee for  $\mathcal{P}_{\text{dc}}$ . Thus the multiclass result preserves the binary proof strategy at the level of classwise affine recursions, while replacing the binary directional evaluator by a nearest-prototype pairwise comparison.

## 4 Simulation

This section empirically examines the mechanisms predicted by the theory. The binary experiments are the primary tests of the single-sample-per-class setting  $\mathcal{S} = \{(\mathbf{x}_+, +1), (\mathbf{x}_-, -1)\}$ : they measure signal aggregation, the effect of the inner-loop length, and the transferability of the learned condensed-set geometry across downstream evaluators. We further include a focused  $K$ -class OvR experiment to evaluate whether the same classwise aggregation effect persists in the one-condensed-sample-per-class extension.

**Setup.** We implement the bilevel gradient matching routine in Algorithm 1 with the deterministic hinge subgradient selection  $q(\mathbf{w}) = \mathbb{1}\{u(\mathbf{w}) < 1\}$  and the classwise update rule specialized to  $n_{\mathcal{S}} = 2$ , and we denote by  $\mathcal{S}^*$  the terminal condensed set after the prescribed  $(T_{\text{out}}, T_{\text{in}})$  bilevel iterations. In the synthetic experiments, we sample a balanced training set  $\mathcal{T}$  from the additive model  $\mathbf{x} = y\boldsymbol{\mu} + \boldsymbol{\xi}$  and report the geometric alignment  $\cos(\mathbf{s}(\mathcal{S}^*), \boldsymbol{\mu})$ , where  $\mathbf{s}(\mathcal{S}) := \mathbf{x}_+ - \mathbf{x}_-$ , together with the downstream test accuracy of a linear SVM trained on  $\mathcal{S}^*$  and evaluated on a held-out test set. On KMNIST, we use a two-class subset and evaluate the same condensed set under multiple downstream models, following the evaluation viewpoint in Subsection 2.4. All reported quantities are averaged over  $n_{\text{seeds}} = 5$  runs and shown as mean  $\pm$  one standard deviation. Full hyperparameters, reproducibility details, and additional numerical tables are reported in Appendix E. The main baseline is a random one-shot coreset that selects one real training sample per class and uses the same downstream training pipeline as  $\mathcal{S}^*$ .

**Condensed-sample aggregation and effectiveness on synthetic data.** To visualize aggregation, Figure 1a projects samples onto the signal coordinate  $\langle \mathbf{x}, \boldsymbol{\mu} / \|\boldsymbol{\mu}\|_2 \rangle$  and an orthogonal coordinate  $\langle \mathbf{x}, \mathbf{v} \rangle$  for a fixed unit vector  $\mathbf{v} \perp \boldsymbol{\mu}$ , and overlays the two learned condensed samples  $\mathbf{x}_+$  and  $\mathbf{x}_-$ ; they concentrate near the signal axis with small orthogonal components and separate along the signal coordinate, consistent with our terminal expansion and coefficient analysis. For downstream performance, Figure 1b compares the test accuracy of a linear SVM trained on the condensed set against two baselines: a random one-shot coreset and training on the full dataset. The condensed set reaches  $0.772 \pm 0.013$  accuracy, close to full-data training ( $0.768 \pm 0.021$ ) and well above the random one-shot baseline ( $0.518 \pm 0.029$ ). This pattern matches the theoretical picture that gradient matching turns the two condensed samples into signal-aligned aggregations of the training set rather than noisy representatives of individual examples.

**Early stopping via inner-loop length.** Our theory isolates activity windows and coefficient bounds that depend on the interplay between the training-set and condensed-set linear windows, suggesting that overly long inner loops may enter regimes where informative gradient directions weaken. We assess this schedule effect by sweeping  $(T_{\text{in}}, T_{\text{out}})$  on synthetic data and reporting the alignment heatmap in Figure 2a; the corresponding test-accuracy heatmap appears in Appendix E. The heatmap shows that alignment is not governed only by the total number of bilevel updates: increasing  $T_{\text{out}}$  is beneficial when the inner loop is

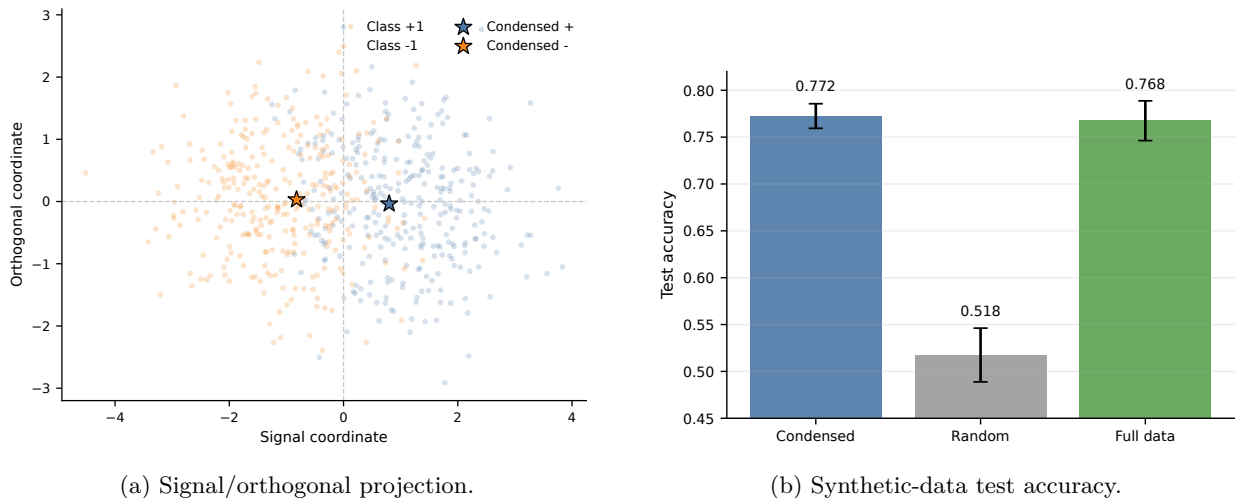


Figure 1: Condensed-sample aggregation and downstream effectiveness on synthetic data. Left: training samples and learned condensed samples in a signal/orthogonal projection. Right: condensed vs. random one-shot vs. full data, reported as mean  $\pm$  one standard deviation over seeds.

short, but the gain saturates and then weakens once  $T_{\text{in}}$  becomes large. The strongest alignments occur for relatively short inner loops combined with moderate outer-loop repetition, whereas large  $T_{\text{in}}$  lowers the learned signal alignment and is accompanied by a modest decrease in test accuracy. This pattern is consistent with the coefficient bounds in Lemma 3.7: after the certified active prefix, additional inner steps contribute less uniform classwise averaging and can discount the previously accumulated signal mass through the contraction factors. Appendix E also reports stepwise diagnostics; in the default synthetic setting the condensed samples remain active, while the training-set activity and gradient norms decrease along longer inner loops, indicating that the observed schedule effect is driven primarily by weakening training-side gradients rather than by a loss of condensed-sample activity.

**Cross-model transfer on KMNIST.** The KMNIST experiment examines whether the geometry learned by SVM-based condensation remains useful beyond the proxy learner used during condensation. We learn the condensed set once and then evaluate it, without re-condensation, under four downstream models: a linear SVM, logistic regression, a nearest-centroid classifier induced by  $\mathbf{s}(\mathcal{S}^*)$ , and a small two-layer ReLU MLP. Figure 2b shows that the same two SVM-condensed samples outperform a random one-shot coreset across all evaluators. The improvement appears not only for the linear SVM used as the proxy learner, but also for logistic regression, nearest-centroid classification, and the two-layer MLP, indicating that the learned samples encode a transferable class-separating geometry rather than only the details of the hinge-loss update. The nearest-centroid performance further confirms that this geometry is already present at the prototype level. This supports the condensed-set geometry perspective: once  $\mathcal{S}^*$  captures a stable class-separating direction, evaluators that depend on that geometry can benefit even when they are not identical to the proxy model used during condensation. Additional numerical tables are reported in Appendix E.

**Multiclass OvR evaluation.** Finally, we run the classwise OvR condensation rule from Appendix B on a balanced  $K = 5$  synthetic additive model with one condensed sample per class, and we evaluate the resulting multiclass set using the nearest-prototype rule in Subsection 2.6. The learned OvR condensed set obtains  $0.624 \pm 0.008$  test accuracy, matching the full-data nearest-prototype reference based on empirical class means within Monte Carlo variability ( $0.624 \pm 0.018$ ) and substantially exceeding a random one-shot multiclass coreset ( $0.231 \pm 0.030$ ). This result provides empirical evidence for the classwise aggregation mechanism behind the multiclass expansion; the corresponding table and figure are given in Appendix E.5.

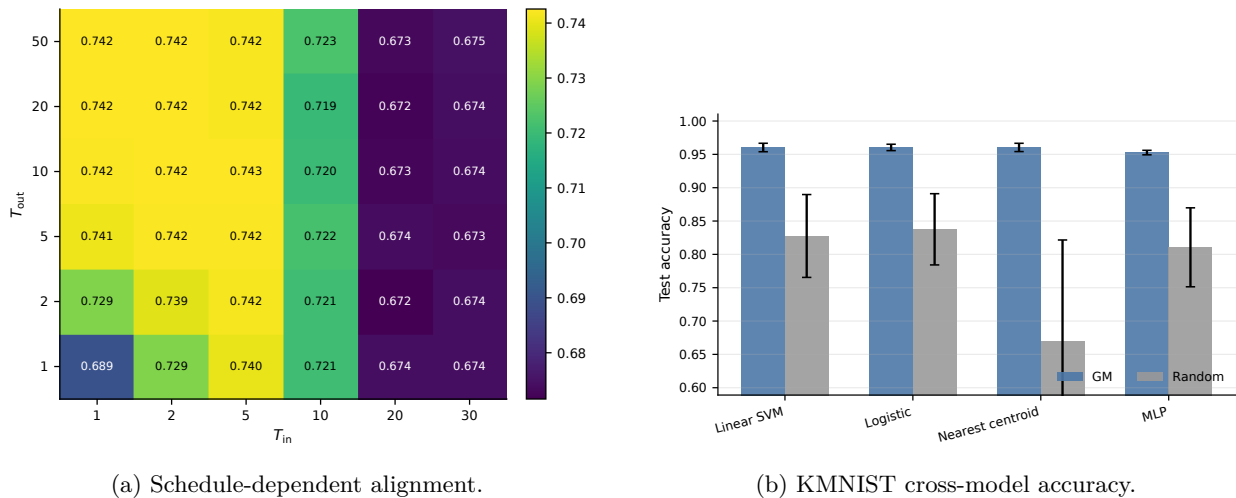


Figure 2: Additional empirical diagnostics. Left: alignment  $\cos(\mathbf{s}(\mathcal{S}^*), \boldsymbol{\mu})$  over the  $(T_{in}, T_{out})$  grid on synthetic data, showing that short-to-moderate inner loops combined with repeated outer-loop restarts yield the strongest signal alignment. Right: KMNIST cross-model evaluation of GM and random one-shot coresets, reported as mean  $\pm$  one standard deviation over seeds; the same condensed set transfers across evaluators beyond the hinge-loss SVM used during condensation.

## 5 Discussion

This work establishes a dynamical framework for dataset condensation in classification, moving beyond regression surrogates to analyze the non-smooth optimization of margin-based learners directly. By characterizing condensed samples as structured signal–noise combinations, we bridge the gap between the algorithmic mechanics of gradient matching and the geometric requirements of generalization.

While anchored in SVMs, this dynamical framework holds potential for characterizing optimization trajectories in broader margin-based classifiers, where gradient signals often change as examples enter or leave high-loss/high-margin regimes. Crucially, our analysis reveals a mechanism where excessive inner-loop updates can attenuate certified signal accumulation, suggesting that early stopping reflects a broader signal-accumulation versus activity-switching tradeoff rather than a peculiarity of the SVM objective.

Our guarantees are proved for the binary one-condensed-sample-per-class setting under an additive sub-Gaussian model, together with a classwise OvR extension for the  $K$ -class case. Extending the same level of guarantee to general softmax or Crammer–Singer multiclass training would require controlling cross-class competition, switching competitors, and coupled gradient-matching dynamics, which are absent from the affine classwise recursions analyzed here. Likewise, relaxing the activity-window and linear-window conditions would require a substantially finer analysis of when samples enter and leave the active hinge regime. These limitations identify concrete technical barriers for extending the analysis beyond affine classwise recursions, and clarify the role of the present tractable margin-based setting as a first step toward classification-specific condensation theory.

Practically, the results suggest that the inner-loop length should be monitored as a mechanism controlling the persistence of informative gradient activity, rather than treated as a purely computational hyperparameter. The analysis points to concrete diagnostics for schedule selection, including alignment, hinge activity, and validation performance.

Furthermore, our theoretical toolkit offers a blueprint for other condensation paradigms. For *distribution matching* (Zhao & Bilen, 2023), our structural decomposition and geometric risk bounds can be adapted to simplify analysis by bypassing non-smooth unrolling. Conversely, for *trajectory matching* (Cazenavette et al., 2022), which shares our dynamical nature, our coefficient expansion and linear-window analysis provide a

foundational approach to disentangle complex cross-time dependencies. Extending these tools to capture such long-horizon structures remains a promising avenue for future research.

## References

- Tatsuya Aoyama, Hanting Yang, Hiroyuki Hanada, Satoshi Akahane, Tomonari Tanaka, Yoshito Okura, Yu Inatsu, Noriaki Hashimoto, Taro Murayama, Hanju Lee, et al. Generalized kernel inducing points by duality gap for dataset distillation. *arXiv preprint arXiv:2502.12607*, 2025.
- Ondrej Bohdal, Yongxin Yang, and Timothy Hospedales. Flexible dataset distillation: Learn labels instead of images. *arXiv preprint arXiv:2006.08572*, 2020.
- Yuan Cao, Zixiang Chen, Misha Belkin, and Quanquan Gu. Benign overfitting in two-layer convolutional neural networks. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 25237–25250. Curran Associates, Inc., 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/a12c999be280372b157294e72a4bbc8b-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/a12c999be280372b157294e72a4bbc8b-Paper-Conference.pdf).
- George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu. Dataset distillation by matching training trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4750–4759, 2022.
- Niladri S Chatterji and Philip M Long. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *Journal of Machine Learning Research*, 22(129):1–30, 2021.
- Yilan Chen, Wei Huang, and Tsui-Wei Weng. Provable and efficient dataset distillation for kernel ridge regression. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 88739–88771. Curran Associates, Inc., 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/file/a1c716638d9b618a1a40a96f473c8250-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/a1c716638d9b618a1a40a96f473c8250-Paper-Conference.pdf).
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Spencer Frei, Niladri S Chatterji, and Peter Bartlett. Benign overfitting without linearity: Neural network classifiers trained by gradient descent for noisy linear data. In *Conference on Learning Theory*, pp. 2668–2703. PMLR, 2022.
- Shuhao Gu, Jialing Zhang, Siyuan Zhou, Kevin Yu, Zhaohu Xing, Liangdong Wang, Zhou Cao, Jintao Jia, Zhuoyi Zhang, Yixuan Wang, et al. Infinity-mm: Scaling multimodal performance with large-scale and high-quality instruction data. *arXiv preprint arXiv:2410.18558*, 2024.
- Zachary Izzo and James Zou. A theoretical study of dataset distillation. In *NeurIPS 2023 Workshop on Mathematics of Modern Machine Learning*, 2023. URL <https://openreview.net/forum?id=dq5QGxGxoJ>.
- Zixuan Jiang, Jiaqi Gu, Mingjie Liu, and David Z. Pan. Delving into effective gradient matching for dataset condensation. In *2023 IEEE International Conference on Omni-layer Intelligent Systems (COINS)*, pp. 1–6, 2023a. doi: 10.1109/COINS57856.2023.10189244.
- Zixuan Jiang, Jiaqi Gu, Mingjie Liu, and David Z Pan. Delving into effective gradient matching for dataset condensation. In *2023 IEEE International Conference on Omni-layer Intelligent Systems (COINS)*, pp. 1–6. IEEE, 2023b.
- Jang-Hyun Kim, Jinuk Kim, Seong Joon Oh, Sangdoo Yun, Hwanjun Song, Joonhyun Jeong, Jung-Woo Ha, and Hyun Oh Song. Dataset condensation via efficient synthetic-data parameterization. In *International Conference on Machine Learning*, pp. 11102–11118. PMLR, 2022.

- Yiwen Kou, Zixiang Chen, Yuanzhou Chen, and Quanquan Gu. Benign overfitting in two-layer ReLU convolutional neural networks. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 17615–17659. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/kou23a.html>.
- Saehyung Lee, Sanghyuk Chun, Sangwon Jung, Sangdoon Yun, and Sungroh Yoon. Dataset condensation with contrastive signals. In *International Conference on Machine Learning*, pp. 12352–12364. PMLR, 2022.
- Zhengquan Luo and Zhiqiang Xu. Utility boundary of dataset distillation: Scaling and configuration-coverage laws. *arXiv preprint arXiv:2512.05817*, 2025.
- Alaa Maalouf, Murad Tukan, Noel Loo, Ramin Hasani, Mathias Lechner, and Daniela Rus. On the size and approximation error of distilled datasets. *Advances in Neural Information Processing Systems*, 36: 61085–61102, 2023.
- Kota Okudo and Kei Kobayashi. Benign overfitting under learning rate conditions for  $\alpha$  sub-exponential input. *arXiv preprint arXiv:2409.00733*, 2024.
- 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.
- Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electronic Communications in Probability*, 18(none), Jan 2013. doi: 10.1214/ecp.v18-2865.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- Iliia Sucholutsky and Matthias Schonlau. Soft-label dataset distillation and text dataset distillation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2021.
- Ke Wang and Christos Thrampoulidis. Benign overfitting in binary classification of gaussian mixtures. *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4030–4034, 2020. doi: 10.1109/ICASSP39728.2021.9413946.
- Ke Wang, Vidya Muthukumar, and Christos Thrampoulidis. Benign overfitting in multiclass classification: All roads lead to interpolation. *IEEE Transactions on Information Theory*, 69:7909–7952, 2021. doi: 10.1109/tit.2023.3320098.
- Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- Kaiyue Wen, Jiaye Teng, and Jingzhao Zhang. Benign overfitting in classification: Provably counter label noise with larger models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=UrEwJebCxx>.
- Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In *International Conference on Machine Learning*, pp. 12674–12685. PMLR, 2021.
- Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 6514–6523, January 2023.
- Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=mSAKhLYLSs1>.

## A Notation

Table 2: Summary of notation used throughout the paper.

Symbol	Meaning
<i>Data and distributions</i>	
$d$	Input dimension.
$\mathcal{Y}$	Label space: $\{-1, +1\}$ in the binary case, or $[K] = \{1, \dots, K\}$ for $K$ -class problems.
$\mathbf{x}, \mathbf{y}$	Feature vector and label.
$\boldsymbol{\mu}_y$	Class-dependent mean (signal); in the binary symmetric case $\boldsymbol{\mu}_{+1} = \boldsymbol{\mu}$ , $\boldsymbol{\mu}_{-1} = -\boldsymbol{\mu}$ and $\mathbf{x} = y\boldsymbol{\mu} + \boldsymbol{\xi}$ .
$\boldsymbol{\xi}, \boldsymbol{\Sigma}_\xi, \zeta$	Mean-zero noise; covariance $\boldsymbol{\Sigma}_\xi$ ; whitened coordinates $\boldsymbol{\xi} = \boldsymbol{\Sigma}_\xi^{-1/2}\zeta$ .
$\ \boldsymbol{\Sigma}_\xi\ _{\text{op}}$	Operator norm of the noise covariance (used in high-probability bounds).
$\mathcal{D}$	Population law of $(\mathbf{x}, y)$ .
$\mathcal{T}, n_{\mathcal{T}}$	Training set $\{(\mathbf{x}_i^T, y_i^T)\}_{i=1}^{n_{\mathcal{T}}}$ and its size.
$\mathcal{I}_k(\mathcal{T}), n_{\mathcal{T},k}$	Indices with label $k$ and the corresponding class count.
$\mathbf{z}_i^T := y_i^T \mathbf{x}_i^T$	Signed training sample.
<i>SVM objective and hinge structure</i>	
$\mathbf{w}$	SVM weight vector.
$\lambda$	$\ell_2$ regularization parameter in $L_{\mathcal{T}}(\mathbf{w})$ and $L_{\mathcal{S}}(\mathbf{w})$ .
$(a)_+$	Positive part $\max\{a, 0\}$ .
$L_{\mathcal{T}}(\mathbf{w}), L_{\mathcal{S}}(\mathbf{w})$	Regularized hinge loss on $\mathcal{T}$ and on the synthetic set $\mathcal{S}$ .
$u_i^T(\mathbf{w})$	Margin $y_i^T \mathbf{w}^T \mathbf{x}_i^T$ .
$q_i^T(\mathbf{w})$	Hinge activity indicator $\mathbb{1}\{u_i^T(\mathbf{w}) < 1\}$ (subgradient convention).
<i>Dataset condensation (bi-level schedule)</i>	
$\mathcal{S}, n_{\mathcal{S}}$	Condensed dataset; one-sample-per-class gives $n_{\mathcal{S}} = 2$ in the binary case.
$T_{\text{out}}, T_{\text{in}}$	Numbers of outer and inner iterations.
$t_{\text{out}}, t_{\text{in}}$	Outer and inner loop indices.
$\eta_{\mathbf{w}}, \eta_{\mathcal{S}}$	Inner-loop step sizes for updating $\mathbf{w}$ (on $\mathcal{T}$ ) and condensed samples.
$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{x}_y^{S, (t_{\text{out}}, t_{\text{in}})}$	Weight and condensed sample after the indicated inner steps.
<i>Directional evaluation (Definition 2.1)</i>	
$\mathbf{s}(\mathcal{S})$	Two-sample difference $\mathbf{x}_+ - \mathbf{x}_-$ for $\mathcal{S} = \{(\mathbf{x}_+, +1), (\mathbf{x}_-, -1)\}$ .
$\hat{\mathbf{w}}_{\mathcal{S}}$	Normalized classifier direction aligned with $\mathbf{s}(\mathcal{S})$ .
$\mathcal{R}(\mathcal{S})$	Population misclassification risk $\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}(y \langle \hat{\mathbf{w}}_{\mathcal{S}}, \mathbf{x} \rangle < 0)$ .
<i>Multiclass classwise OvR extension (Appendix B)</i>	
$\mathcal{P}(t)$	Multiclass condensed set at time $t$ : $\{(\mathbf{p}_k^{S, (t)}, k) : k \in [K]\}$ , with one condensed sample per class.
$\mathcal{P}_{\text{dc}}$	Terminal multiclass condensed set returned by the classwise OvR condensation procedure.
$\mathbf{p}_k^{S, (t)}, \mathbf{p}_k$	Class- $k$ condensed sample at time $t$ ; generic class- $k$ representative in a multiclass nearest-prototype evaluator.
$\mathbf{W}^{(t)}, \mathbf{w}_h^{(t)}$	Collection of $K$ OvR heads and the head- $h$ weight vector at time $t$ .
$t^+$	Next inner-loop time $(t_{\text{out}}, t_{\text{in}} + 1)$ , used when $t_{\text{in}} < T_{\text{in}}$ .
$\tau_{h,k}$	OvR sign $2\mathbb{1}\{h = k\} - 1$ , treating class $h$ as positive for head $h$ .
$\mathbf{z}_{i,h}^T, \bar{\mathbf{z}}_h^T$	Head-wise signed training sample $\tau_{h,y_i^T} \mathbf{x}_i^T$ and its empirical mean over $\mathcal{T}$ .
$\mathbf{z}_{k,h}^{S, (t)}$	Head-wise signed condensed sample $\tau_{h,k} \mathbf{p}_k^{S, (t)}$ .
$\tilde{S}_{k,h}^{(t)}, q_{i,h}^T, q_{i,h}^{\mathcal{T}, (t)}$	Synthetic and training hinge activity indicators for class $k$ / sample $i$ under OvR head $h$ .
$\mathbf{g}_{k,h}^{S, (t)}, \mathbf{g}_{k,h}^{\mathcal{T}, (t)}$	Classwise synthetic and training OvR gradient terms for class $k$ under head $h$ .
$D_{\text{OvR}}^k(\mathcal{P}(t); \mathbf{W}^{(t)})$	Scalar classwise OvR gradient-matching discrepancy for class $k$ .
$\bar{q}_k^{S, (t)}, \bar{b}_{k,i}^{(t)}$	Averaged synthetic activity and averaged synthetic-training activity product over the $K$ OvR heads.
$h_{\mathcal{P}}, \mathcal{R}_{\text{mc}}(\mathcal{P})$	Nearest-prototype multiclass evaluator induced by $\mathcal{P}$ and its population misclassification risk.
$\pi_k, \gamma_{k\ell}, \epsilon_{\text{mc}}(\mathcal{P})$	Class prior, pairwise class separation $\ \boldsymbol{\mu}_k - \boldsymbol{\mu}_\ell\ _2$ , and maximum multiclass condensed-sample perturbation.
<i>Multiclass OvR coefficient and risk quantities</i>	
$u_k^{(\mathbf{t}^*)}, c_{k,i}^{(\mathbf{t}^*)}, c_{k,\mu}^{(\mathbf{t}^*)}$	Terminal OvR expansion coefficients for class $k$ , with $c_{k,\mu}^{(\mathbf{t}^*)} = \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)}$ .
$\mathbf{c}_k^{(\mathbf{t}^*)}$	Vector of terminal same-class coefficients $(c_{k,i}^{(\mathbf{t}^*)})_{i \in \mathcal{I}_k(\mathcal{T})}$ .

Continued on next page

Table 2: *Table 2 continued*

Symbol	Meaning
$T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r)$	Class- $k$ OvR training and condensed-sample linear-window lengths in outer loop $r$ .
$T_{k,\text{lin}}^{\text{OvR}}(r)$	Joint class- $k$ OvR window length $\min\{T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r)\}$ .
$B_{k,\text{out}}, \ell_k, u_{k,\text{out}}, \Delta_k$	Class- $k$ OvR coefficient upper budget, certified averaging mass, upper envelope, and non-averaged remainder.
$\bar{u}_k$	Upper bound on the terminal initialization coefficient induced by the class- $k$ OvR condensed-sample windows.
$\mathcal{W}_k^{\text{OvR}}$	Deterministic window-deficit diagnostic for class $k$ in the multiclass OvR active-window comparison.
$\widehat{\text{SR}}_{k\ell}^{\text{mc}}, \underline{\text{SR}}_{k\ell}^{\text{mc}}$	Pairwise multiclass signal-to-residual ratio and its certified lower bound for classes $k$ and $\ell$ .
$L_{k\ell}^{\text{mc}}$	Pairwise multiclass factor combining class-averaging mass, initialization remainder, signal-recovery error, and active-window deficit.
$\Psi_{\text{mc}}$	Monotone function mapping pairwise multiclass signal-to-residual ratios into the exponential risk bound.
$\gamma_{\min}, \underline{\text{SR}}_{\min}^{\text{mc}}$	Minimum pairwise class separation and minimum certified pairwise multiclass signal-to-residual ratio.
<i>Main theoretical objects</i>	
$t^*$	Terminal time index ( $T_{\text{out}} - 1, T_{\text{in}}$ ).
$\mathbf{z}_y^{\mathcal{S},(t^*)}$	Terminal signed condensed sample.
$\rho_{\mathcal{S}} := 1 - 2\eta_{\mathcal{S}}$	Contraction factor in condensed-sample dynamics (stability: $0 < 2\eta_{\mathcal{S}} \leq 1$ ).
$T_{\mathcal{T},\text{lin}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})$	Linear-window lengths where training / condensed hinge indicators stay active.
$T_{y,\text{lin}}(r)$	$\min\{T_{\mathcal{T},\text{lin}}, T_{\mathcal{S},\text{lin}}(r)\}$ .
$\tilde{u}_y^{(t^*)}, \tilde{c}_{y,j}^{(t^*)}$	Coefficients in the terminal expansion of $\mathbf{z}_y^{\mathcal{S}}$ over initialization and $\{\mathbf{z}_j^{\mathcal{T}}\}$ .
$p_y, \beta_p$	Class proportion $n_{\mathcal{T}_y}/n_{\mathcal{T}}$ and balance factor $p_+^{-1/2} + p_-^{-1/2}$ .
$B_{\text{out}}$	Outer-loop accumulation budget for terminal coefficient upper bounds; its construction is given in Lemma D.8.
$\ell_y, u_{\text{out}}, \Delta_y$	Coefficient mass certified by classwise averaging, coefficient upper envelope, and the non-averaged remainder $u_{\text{out}} - \ell_y$ ; constructed in Lemma D.8.
$\widehat{\text{SR}}(\mathcal{S})$	Signal-to-residual ratio for $\mathbf{x}_+ - \mathbf{x}_- = A(\mathcal{S})\boldsymbol{\mu} + \mathbf{r}(\mathcal{S})$ .
$\underline{\text{SR}}_{\text{dc}}$	Certified lower bound on $\widehat{\text{SR}}$ for the terminal condensed set.
$L$	Dimensionless factor combining coefficient mass from certified classwise averaging, class balance, and the non-averaged coefficient remainder.
$\mathcal{W}_{\text{win}}$	Deterministic diagnostic for the loss caused when certified active prefixes are shorter than the full inner loop, used in Lemma D.12 to compare regimes for $L$ .
$\mathcal{S}_{\text{dc}}, \mathcal{S}_{\text{core}}$	Terminal condensed set vs. random one-shot coresets from two training points.
$G_{\text{max}}$	Upper bound on inner products $\langle \mathbf{z}_i^{\mathcal{T}}, \mathbf{z}_{i'}^{\mathcal{T}} \rangle$ (appendix); regime $G_{\text{max}} < \lambda$ gives full hinge activity.
<i>Assumption 2.2</i>	
$\delta, \kappa, C_{\kappa}$	Tail probability, dimension-log constant, and associated concentration constant.
$\sigma_{\mathbf{w}}, \sigma_{\mathcal{S}}$	Initialization scales for $\mathbf{w}$ and condensed samples.
<i>Inner-loop weight dynamics (appendix)</i>	
$\rho := 1 - \lambda\eta_{\mathbf{w}}$	Contraction factor for weight updates under fixed hinge activity.
$\mathcal{I}(\mathcal{A}), n_{\mathcal{A}}$	Index set and size of a generic dataset $\mathcal{A}$ .
$\bar{\mathbf{z}}_{\mathcal{T}}$	Mean of signed training samples $\frac{1}{n_{\mathcal{T}}} \sum_{i \in \mathcal{I}(\mathcal{T})} \mathbf{z}_i^{\mathcal{T}}$ .
<i>Sub-Gaussian noise.</i> Vectors satisfying Definition C.1 (Appendix C.1).	

## B Multiclass One-Condensed-Sample-Per-Class Extension

This appendix extends the coefficient-expansion mechanism of the binary analysis to a  $K$ -class one-condensed-sample-per-class setting. Rather than applying a standard one-vs-rest reduction that runs  $K$  separate binary condensation problems and produces a two-point condensed set for each induced subproblem, we maintain a single multiclass condensed set with one synthetic sample per class. The extension uses a classwise OvR gradient-matching objective and preserves the classwise affine-contraction structure that drives the binary hinge-loss analysis.

We work under the  $K$ -class additive model in Section 2.1:  $y \in [K]$  has class prior  $\pi_k = \mathbb{P}(y = k)$ , and  $\mathbf{x} = \boldsymbol{\mu}_y + \boldsymbol{\xi}$  with mean-zero sub-Gaussian noise proxy covariance  $\boldsymbol{\Sigma}_\xi$ , with independent noise across samples. The results below are stated in terms of the empirical class sizes  $n_{\mathcal{T},k}$  and the pairwise separations  $\gamma_{k\ell} = \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_\ell\|_2$ , without imposing simplex symmetry on the class means.

### B.1 Classwise OvR Gradient Matching

We use the same two-dimensional time convention as in the main text:  $\mathbf{t} = (t_{\text{out}}, t_{\text{in}})$ , where  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and  $t_{\text{in}} \in \{0, \dots, T_{\text{in}}\}$ . For an inner-loop update with  $t_{\text{in}} < T_{\text{in}}$ , write  $\mathbf{t}^+ := (t_{\text{out}}, t_{\text{in}} + 1)$ . For  $t_{\text{out}} = 0, \dots, T_{\text{out}} - 2$ , the condensed samples are inherited across outer loops as  $\mathbf{p}_k^{\mathcal{S},(t_{\text{out}}+1,0)} = \mathbf{p}_k^{\mathcal{S},(t_{\text{out}},T_{\text{in}})}$ .

For  $K$ -class classification, the condensed object at time  $\mathbf{t}$  is the multiclass condensed set  $\mathcal{P}^{(\mathbf{t})} := \{(\mathbf{p}_k^{\mathcal{S},(\mathbf{t})}, k) : k \in [K]\}$ , where  $\mathbf{p}_k^{\mathcal{S},(\mathbf{t})} \in \mathbb{R}^d$  is the single condensed sample assigned to class  $k$ . We write  $\mathcal{P}^{(\mathbf{t})}$  to distinguish this multiclass object from the binary condensed-set notation  $\mathcal{S}$ . The condensed samples are initialized independently as  $\mathbf{p}_k^{\mathcal{S},(0,0)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathcal{S}}^2 \mathbf{I}_d)$  for  $k \in [K]$ . We collect the  $K$  OvR heads as  $\mathbf{W}^{(\mathbf{t})} := (\mathbf{w}_h^{(\mathbf{t})})_{h=1}^K$ , where  $\mathbf{w}_h^{(\mathbf{t})} \in \mathbb{R}^d$ . At the beginning of each outer loop, every head is reinitialized as  $\mathbf{w}_h^{(\mathbf{t}_{\text{out}},0)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{w}}^2 \mathbf{I}_d)$ .

For head  $h$ , class  $h$  is treated as positive and all other classes as negative. Define the OvR sign  $\tau_{h,k} := 2\mathbb{1}\{h = k\} - 1 \in \{+1, -1\}$  for  $h, k \in [K]$ . For head-wise signed comparisons, set  $\mathbf{z}_{i,h}^{\mathcal{T}} := \tau_{h,y_i^{\mathcal{T}}} \mathbf{x}_i^{\mathcal{T}}$ ,  $\bar{\mathbf{z}}_h^{\mathcal{T}} := n_{\mathcal{T}}^{-1} \sum_{i=1}^{n_{\mathcal{T}}} \mathbf{z}_{i,h}^{\mathcal{T}}$ , and  $\mathbf{z}_{k,h}^{\mathcal{S},(\mathbf{t})} := \tau_{h,k} \mathbf{p}_k^{\mathcal{S},(\mathbf{t})}$ . These signed coordinates express each OvR head in the same hinge-active form as the binary analysis.

**Assumption B.1.** For every head  $h$ , the signed training samples  $\{\mathbf{z}_{i,h}^{\mathcal{T}}\}_{i=1}^{n_{\mathcal{T}}}$  and signed condensed samples  $\{\mathbf{z}_{k,h}^{\mathcal{S},(\mathbf{t})}\}_{k=1}^K$  satisfy the analogues of Assumption D.1 and Lemmas C.2–C.5, with the binary signed variables replaced by their head-wise OvR counterparts and with the corresponding suprema taken over heads, classes, condensed samples, and training samples. Thus the initialization bounds, signed-sample concentration bounds, and active-window margin envelopes required for the binary coefficient estimates hold uniformly over the OvR signed coordinates.

The synthetic hinge activity of the class- $k$  condensed sample under head  $h$  is  $q_{k,h}^{\mathcal{S},(\mathbf{t})} := \mathbb{1}\{1 - \tau_{h,k} \langle \mathbf{w}_h^{(\mathbf{t})}, \mathbf{p}_k^{\mathcal{S},(\mathbf{t})} \rangle > 0\}$ . For a training example  $(\mathbf{x}_i^{\mathcal{T}}, y_i^{\mathcal{T}})$ , define  $q_{i,h}^{\mathcal{T},(\mathbf{t})} := \mathbb{1}\{1 - \tau_{h,y_i^{\mathcal{T}}} \langle \mathbf{w}_h^{(\mathbf{t})}, \mathbf{x}_i^{\mathcal{T}} \rangle > 0\}$ ; in particular, if  $i \in \mathcal{I}_k(\mathcal{T})$ , then  $y_i^{\mathcal{T}} = k$  and the sign is  $\tau_{h,k}$ .

For class  $k$  and head  $h$ , define the classwise OvR gradient terms  $\mathbf{g}_{k,h}^{\mathcal{S},(\mathbf{t})} := q_{k,h}^{\mathcal{S},(\mathbf{t})} \tau_{h,k} \mathbf{p}_k^{\mathcal{S},(\mathbf{t})}$  and  $\mathbf{g}_{k,h}^{\mathcal{T},(\mathbf{t})} := n_{\mathcal{T},k}^{-1} \sum_{i \in \mathcal{I}_k(\mathcal{T})} q_{i,h}^{\mathcal{T},(\mathbf{t})} \tau_{h,k} \mathbf{x}_i^{\mathcal{T}}$ . The scalar classwise OvR gradient-matching discrepancy for class  $k$  is

$$D_k^{\text{OvR}}(\mathcal{P}^{(\mathbf{t})}; \mathbf{W}^{(\mathbf{t})}) := \frac{1}{K} \sum_{h=1}^K \left\| \mathbf{g}_{k,h}^{\mathcal{S},(\mathbf{t})} - \mathbf{g}_{k,h}^{\mathcal{T},(\mathbf{t})} \right\|_2^2.$$

Since the same sign  $\tau_{h,k}$  multiplies the synthetic and training terms within a fixed class block, the head-wise squared norm is unchanged if this common sign is removed.

Define the averaged synthetic activity  $\bar{q}_k^{\mathcal{S},(\mathbf{t})} := K^{-1} \sum_{h=1}^K q_{k,h}^{\mathcal{S},(\mathbf{t})}$  and, for  $i \in \mathcal{I}_k(\mathcal{T})$ , the averaged activity product  $\bar{b}_{k,i}^{(\mathbf{t})} := K^{-1} \sum_{h=1}^K q_{k,h}^{\mathcal{S},(\mathbf{t})} q_{i,h}^{\mathcal{T},(\mathbf{t})}$ . By construction,  $0 \leq \bar{b}_{k,i}^{(\mathbf{t})} \leq \bar{q}_k^{\mathcal{S},(\mathbf{t})} \leq 1$ .

**Algorithm 2** Classwise OvR Gradient Matching for Multiclass Condensed Samples

---

**Input:** Training set  $\mathcal{T}$ , number of classes  $K$ , loop lengths  $T_{\text{out}}$  and  $T_{\text{in}}$ , step sizes  $\eta_{\mathbf{w}}$  and  $\eta_S$   
**Output:** Multiclass condensed set  $\mathcal{P}_{\text{dc}}$   
Initialize  $\mathbf{p}_k^{S,(0,0)} \sim \mathcal{N}(\mathbf{0}, \sigma_S^2 \mathbf{I}_d)$  for  $k \in [K]$   
Set  $\tau_{h,k} = 2\mathbb{1}\{h = k\} - 1$  for  $h, k \in [K]$   
**for**  $t_{\text{out}} = 0$  **to**  $T_{\text{out}} - 1$  **do**  
  Initialize  $\mathbf{w}_h^{(t_{\text{out}},0)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{w}}^2 \mathbf{I}_d)$  for  $h \in [K]$   
  **if**  $t_{\text{out}} > 0$  **then**  
    Set  $\mathbf{p}_k^{S,(t_{\text{out}},0)} \leftarrow \mathbf{p}_k^{S,(t_{\text{out}}-1,T_{\text{in}})}$  for each  $k \in [K]$   
  **end if**  
  **for**  $t_{\text{in}} = 0$  **to**  $T_{\text{in}} - 1$  **do**  
    Set  $\mathbf{t} = (t_{\text{out}}, t_{\text{in}})$  and  $\mathbf{t}^+ = (t_{\text{out}}, t_{\text{in}} + 1)$   
    Compute  $q_{k,h}^{S,\mathbf{t}}$  and  $q_{i,h}^{\mathcal{T},\mathbf{t}}$  for all  $k, h \in [K]$  and  $i \in \{1, \dots, n_{\mathcal{T}}\}$   
    Compute  $\bar{q}_k^{S,\mathbf{t}}$  and  $\bar{b}_{k,i}^{(\mathbf{t})}$  for every  $k \in [K]$  and  $i \in \mathcal{I}_k(\mathcal{T})$   
    **for**  $k \in [K]$  **do**  
      Update  $\mathbf{p}_k^{S,\mathbf{t}^+}$  by equation 21  
    **end for**  
    **for**  $h \in [K]$  **do**  
      Compute  $\mathbf{g}_h^{\mathcal{T},\mathbf{t}}$  and update  $\mathbf{w}_h^{(\mathbf{t}^+)}$  by equation 22  
    **end for**  
  **end for**  
**end for**  
**Return:**  $\mathcal{P}_{\text{dc}} = \{(\mathbf{p}_k^{S,(T_{\text{out}}-1,T_{\text{in}})}, k) : k \in [K]\}$

---

For  $t_{\text{in}} < T_{\text{in}}$ , the multiclass condensed-sample update is the classwise gradient step  $\mathbf{p}_k^{S,\mathbf{t}^+} = \mathbf{p}_k^{S,\mathbf{t}} - \eta_S \nabla_{\mathbf{p}_k^S} D_k^{\text{OvR}}(\mathcal{P}(\mathbf{t}); \mathbf{W}(\mathbf{t}))$ , where the hinge activities are evaluated at time  $\mathbf{t}$  and are not differentiated through, as in the binary analysis. Under this convention, the update has the explicit affine form

$$\mathbf{p}_k^{S,\mathbf{t}^+} = \left(1 - 2\eta_S \bar{q}_k^{S,\mathbf{t}}\right) \mathbf{p}_k^{S,\mathbf{t}} + \frac{2\eta_S}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \bar{b}_{k,i}^{(\mathbf{t})} \mathbf{x}_i^{\mathcal{T}}. \quad (21)$$

For  $t_{\text{in}} < T_{\text{in}}$  and each head  $h$ , the full training OvR gradient term is  $\mathbf{g}_h^{\mathcal{T},\mathbf{t}} := n_{\mathcal{T}}^{-1} \sum_{i=1}^{n_{\mathcal{T}}} q_{i,h}^{\mathcal{T},\mathbf{t}} \tau_{h,y_i^{\mathcal{T}}} \mathbf{x}_i^{\mathcal{T}}$ , and the real-data SVM update is

$$\mathbf{w}_h^{(\mathbf{t}^+)} = (1 - \eta_{\mathbf{w}} \lambda) \mathbf{w}_h^{(\mathbf{t})} + \eta_{\mathbf{w}} \mathbf{g}_h^{\mathcal{T},\mathbf{t}}. \quad (22)$$

Algorithm 2 summarizes the resulting classwise OvR gradient-matching procedure. It maintains a single multiclass condensed set across outer loops, while the  $K$  OvR heads are reinitialized at every outer-loop restart.

## B.2 Multiclass Nearest-Prototype Evaluation

The multiclass evaluator below turns a condensed set into a prediction rule and provides the target of the population risk bounds.

**Definition B.2** (Multiclass nearest-prototype evaluator). Let  $\mathcal{P} = \{(\mathbf{p}_k, k) : k \in [K]\}$  be a generic multiclass condensed set with one representative per class, not necessarily tied to a particular condensation time. In this evaluator, these representatives act as class prototypes, and the nearest-prototype classifier induced by  $\mathcal{P}$  is

$$h_{\mathcal{P}}(\mathbf{x}) := \arg \min_{k \in [K]} \|\mathbf{x} - \mathbf{p}_k\|_2^2 = \arg \max_{k \in [K]} \left\{ \langle \mathbf{x}, \mathbf{p}_k \rangle - \frac{1}{2} \|\mathbf{p}_k\|_2^2 \right\},$$

with deterministic tie-breaking. The corresponding multiclass population risk is

$$\mathcal{R}_{\text{mc}}(\mathcal{P}) := \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} (h_{\mathcal{P}}(\mathbf{x}) \neq y).$$

**Examples of multiclass nearest-prototype evaluators.** The nearest-prototype evaluator is a standard multiclass prediction family, and it appears in several familiar classifiers, either exactly or after a fixed feature transformation.

- **Nearest class mean.** If  $\hat{\boldsymbol{\mu}}_k^T := n_{\mathcal{T},k}^{-1} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \mathbf{x}_i^T$  and  $\mathcal{P}_{\text{NCM}} := \{(\hat{\boldsymbol{\mu}}_k^T, k) : k \in [K]\}$ , then  $h_{\mathcal{P}_{\text{NCM}}}(\mathbf{x}) = \arg \min_{k \in [K]} \|\mathbf{x} - \hat{\boldsymbol{\mu}}_k^T\|_2^2$ , which is the standard nearest-class-mean classifier. It is the evaluator most closely connected to the full-activity specialization, in which each condensed sample contracts toward the corresponding class mean.
- **Shared-covariance Gaussian Bayes and LDA.** If  $\mathbf{x} \mid y = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$  with equal class priors, the Bayes rule is  $\arg \min_{k \in [K]} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)$ . After whitening,  $\tilde{\mathbf{x}} := \boldsymbol{\Sigma}^{-1/2} \mathbf{x}$  and  $\tilde{\mathbf{p}}_k := \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\mu}_k$ , this becomes the nearest-prototype rule  $\arg \min_{k \in [K]} \|\tilde{\mathbf{x}} - \tilde{\mathbf{p}}_k\|_2^2$ . Unequal priors add class-dependent bias terms and therefore correspond to a biased nearest-prototype score rather than the unweighted Euclidean form unless one augments the feature space.
- **Affine linear heads.** Any affine multiclass linear prediction rule  $h_{\text{lin}}(\mathbf{x}) = \arg \max_{k \in [K]} \{\langle \mathbf{a}_k, \mathbf{x} \rangle + b_k\}$  is exactly a nearest-prototype rule when  $\mathbf{p}_k = \mathbf{a}_k$  and  $b_k = -\frac{1}{2} \|\mathbf{a}_k\|_2^2 + C$  for a class-independent constant  $C$ . More generally, every affine linear rule can be represented exactly as a nearest-prototype rule after adding one dummy feature dimension: choose  $C \geq \max_k \{b_k + \frac{1}{2} \|\mathbf{a}_k\|_2^2\}$ , set  $r_k^2 := 2(C - b_k) - \|\mathbf{a}_k\|_2^2 \geq 0$ , and use  $\bar{\mathbf{x}} := (\mathbf{x}, 0)$  and  $\bar{\mathbf{p}}_k := (\mathbf{a}_k, r_k)$ . Then  $\langle \bar{\mathbf{x}}, \bar{\mathbf{p}}_k \rangle - \frac{1}{2} \|\bar{\mathbf{p}}_k\|_2^2 = \langle \mathbf{x}, \mathbf{a}_k \rangle + b_k - C$ , so the class-independent shift does not change the argmax. This representation includes the prediction rules of multiclass logistic regression, linear softmax heads, and multiclass linear SVMs.
- **Prototypical and cosine classifiers.** In a fixed representation space  $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^r$ , the rule  $\arg \min_{k \in [K]} \|\phi(\mathbf{x}) - \mathbf{p}_k\|_2^2$  is the usual prototypical-network or embedding-space nearest-class-mean prediction rule when  $\mathbf{p}_k = n_{\mathcal{T},k}^{-1} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \phi(\mathbf{x}_i^T)$ . If the features and prototypes are normalized, or more generally if all prototypes have the same norm, the nearest-prototype score is equivalent to maximizing the cosine or inner-product score  $\langle \phi(\mathbf{x}), \mathbf{p}_k \rangle$ . This connects the evaluator to normalized linear heads in a fixed feature space.

For the  $K$ -class additive model, write  $\pi_k := \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}(y = k)$  and  $\gamma_{k\ell} := \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_\ell\|_2$  for  $k, \ell \in [K]$ . For a multiclass condensed set  $\mathcal{P} = \{(\mathbf{p}_k, k) : k \in [K]\}$ , define the maximum condensed-sample perturbation  $\epsilon_{\text{mc}}(\mathcal{P}) := \max_{k \in [K]} \|\mathbf{p}_k - \boldsymbol{\mu}_k\|_2$ . When the condensed set is clear from context, write  $\epsilon_{\text{mc}}$ .

### B.3 Main Multiclass Guarantees

The following results give the multiclass consequences of the classwise OvR update. The first result is the direct analogue of the binary coefficient expansion: each terminal condensed sample is a nonnegative combination of its initialization and same-class training samples, with coefficient bounds controlled by the certified active windows. Lemma D.13 gives sufficient margin-envelope conditions for these windows.

**Theorem B.3** (Classwise OvR coefficient expansion). *Assume  $0 < 2\eta_S \leq 1$  and  $n_{\mathcal{T},k} \geq 1$  for every  $k \in [K]$ . Set  $\rho_S := 1 - 2\eta_S$  and  $\mathbf{t}^* := (T_{\text{out}} - 1, T_{\text{in}})$ . Under classwise OvR gradient matching, for every class  $k \in [K]$  there exist coefficients  $u_k^{(\mathbf{t}^*)} \geq 0$  and  $c_{k,i}^{(\mathbf{t}^*)} \geq 0$ ,  $i \in \mathcal{I}_k(\mathcal{T})$ , such that*

$$\mathbf{p}_k^{S,(\mathbf{t}^*)} = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{S,(0,0)} + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \mathbf{x}_i^T.$$

Consequently, under the  $K$ -class additive model, with  $c_{k,\mu}^{(\mathbf{t}^*)} := \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)}$ ,

$$\mathbf{p}_k^{S,(\mathbf{t}^*)} = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{S,(0,0)} + c_{k,\mu}^{(\mathbf{t}^*)} \boldsymbol{\mu}_k + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \boldsymbol{\xi}_i.$$

Assume further that the class- $k$  training activities and condensed-sample activities are active on windows of lengths  $T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$  and  $T_{S,k,\text{lin}}^{\text{OvR}}(r)$ , respectively, in every outer loop  $r \in \{0, \dots, T_{\text{out}} - 1\}$ . Define  $T_{k,\text{lin}}^{\text{OvR}}(r) :=$

$$\min\{T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, T_{S,k,\text{lin}}^{\text{OvR}}(r)\},$$

$$B_{k,\text{out}} := 1 + \sum_{r=0}^{T_{\text{out}}-2} \rho_S^{\sum_{s=r+1}^{T_{\text{out}}-1}} T_{S,k,\text{lin}}^{\text{OvR}}(s), \quad u_{k,\text{out}} := (1 - \rho_S^{T_{\text{in}}})B_{k,\text{out}},$$

and

$$\ell_k := \sum_{r=0}^{T_{\text{out}}-1} \rho_S^{(T_{\text{out}}-r)T_{\text{in}}-T_{k,\text{lin}}^{\text{OvR}}(r)} (1 - \rho_S^{T_{k,\text{lin}}^{\text{OvR}}(r)}).$$

Then

$$\rho_S^{T_{\text{out}}T_{\text{in}}} \leq u_k^{(\mathbf{t}^*)} \leq \rho_S^{\sum_{r=0}^{T_{\text{out}}-1} T_{S,k,\text{lin}}^{\text{OvR}}(r)}, \quad \frac{\ell_k}{n_{\mathcal{T},k}} \leq c_{k,i}^{(\mathbf{t}^*)} \leq \frac{u_{k,\text{out}}}{n_{\mathcal{T},k}}.$$

In particular,  $\ell_k \leq c_{k,\mu}^{(\mathbf{t}^*)} \leq u_{k,\text{out}}$  and  $\ell_k/\sqrt{n_{\mathcal{T},k}} \leq \|\mathbf{c}_k^{(\mathbf{t}^*)}\|_2 \leq u_{k,\text{out}}/\sqrt{n_{\mathcal{T},k}}$ , where  $\mathbf{c}_k^{(\mathbf{t}^*)} := (c_{k,i}^{(\mathbf{t}^*)})_{i \in \mathcal{I}_k(\mathcal{T})}$ . If all OvR activities are active throughout the bilevel procedure, then

$$\mathbf{p}_k^{S,(\mathbf{t}^*)} = \rho_S^{T_{\text{out}}T_{\text{in}}} \mathbf{p}_k^{S,(0,0)} + (1 - \rho_S^{T_{\text{out}}T_{\text{in}}}) \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}}, \quad \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}} := \frac{1}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \mathbf{x}_i^{\mathcal{T}}.$$

*Proof.* The active-window certification is Lemma D.13. The coefficient expansion, coefficient bounds, and full-activity contraction are Lemma D.14.  $\square$

The next result is independent of the OvR dynamics. It states that nearest-prototype evaluation has a pairwise risk bound controlled by class separation and condensed-sample perturbation.

**Theorem B.4** (Multiclass nearest-prototype risk). *Assume the  $K$ -class additive test model  $\mathbf{x} = \boldsymbol{\mu}_y + \boldsymbol{\xi}$ , where  $\boldsymbol{\xi}$  is mean-zero sub-Gaussian with proxy covariance  $\boldsymbol{\Sigma}_\xi$ . Let  $\mathcal{P} = \{(\mathbf{p}_k, k) : k \in [K]\}$  be any multiclass condensed set used with the nearest-prototype evaluator. Define  $\epsilon_k := \|\mathbf{p}_k - \boldsymbol{\mu}_k\|_2$ ,  $R_{k\ell} := \epsilon_k + \epsilon_\ell$ , and  $\widehat{\text{SR}}_{k\ell}^{\text{mc}} := \gamma_{k\ell}/R_{k\ell}$ , with value  $+\infty$  when  $R_{k\ell} = 0$ . For  $s \in (0, +\infty)$ , set*

$$\Psi_{\text{mc}}(s) := \frac{((1-s^{-1})_+)^4}{(1+s^{-1})^2}, \quad \Psi_{\text{mc}}(0) := 0, \quad \Psi_{\text{mc}}(+\infty) := 1.$$

Then

$$\mathcal{R}_{\text{mc}}(\mathcal{P}) \leq \sum_{k=1}^K \pi_k \sum_{\ell \neq k} \exp\left(-\frac{\gamma_{k\ell}^2}{8\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \Psi_{\text{mc}}\left(\widehat{\text{SR}}_{k\ell}^{\text{mc}}\right)\right).$$

*Proof.* The statement is proved in Lemma D.15.  $\square$

Combining the OvR coefficient expansion with the nearest-prototype risk bound gives a certified multiclass guarantee for the terminal condensed set. The bound is written through pairwise signal-to-residual ratios; larger pairwise ratios yield stronger exponential risk decay.

**Corollary B.5** (Certified multiclass risk for classwise OvR condensation). *Assume the hypotheses of Theorem B.3 and the high-probability event in Assumption B.1. Let  $C_\kappa$  denote the corresponding uniform concentration constant and define the terminal multiclass condensed set  $\mathcal{P}_{\text{dc}} := \{(\mathbf{p}_k^{S,(\mathbf{t}^*)}, k) : k \in [K]\}$ . Set  $p_k := n_{\mathcal{T},k}/n_{\mathcal{T}}$ ,*

$$\alpha_S := \frac{\sigma_S \sqrt{n_{\mathcal{T}}}}{\sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}}, \quad \alpha_{\mu,k} := \frac{\|\boldsymbol{\mu}_k\|_2 \sqrt{n_{\mathcal{T}}}}{\sqrt{C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}}, \quad \bar{u}_k := \rho_S^{\sum_{r=0}^{T_{\text{out}}-1} T_{S,k,\text{lin}}^{\text{OvR}}(r)}.$$

For  $k \neq \ell$ , define

$$L_{k\ell}^{\text{mc}} := \left[ \alpha_S (\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k} (1 - \ell_k) + \alpha_{\mu,\ell} (1 - \ell_\ell) + \ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2} + \sqrt{n_{\mathcal{T}}} (\Delta_k + \Delta_\ell) \right]^{-1},$$

where  $\Delta_k := u_{k,\text{out}} - \ell_k$ , and set

$$\underline{\text{SR}}_{k\ell}^{\text{mc}} := \frac{\sqrt{n_{\mathcal{T}}} L_{k\ell}^{\text{mc}} \gamma_{k\ell}}{\sqrt{C_{\kappa} \|\Sigma_{\xi}\|_{\text{op}} d}}.$$

Then

$$\mathcal{R}_{\text{mc}}(\mathcal{P}_{\text{dc}}) \leq \sum_{k=1}^K \pi_k \sum_{\ell \neq k} \exp\left(-\frac{\gamma_{k\ell}^2}{8\|\Sigma_{\xi}\|_{\text{op}}} \Psi_{\text{mc}}(\underline{\text{SR}}_{k\ell}^{\text{mc}})\right).$$

In particular, with  $\gamma_{\min} := \min_{k \neq \ell} \gamma_{k\ell}$  and  $\underline{\text{SR}}_{\min}^{\text{mc}} := \min_{k \neq \ell} \underline{\text{SR}}_{k\ell}^{\text{mc}}$ ,

$$\mathcal{R}_{\text{mc}}(\mathcal{P}_{\text{dc}}) \leq (K-1) \exp\left(-\frac{\gamma_{\min}^2}{8\|\Sigma_{\xi}\|_{\text{op}}} \Psi_{\text{mc}}(\underline{\text{SR}}_{\min}^{\text{mc}})\right).$$

If all OvR activities are active throughout training, then  $L_{k\ell}^{\text{mc}}$  can be replaced by

$$L_{k\ell,\text{full}}^{\text{mc}} := \left[ \rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}} (2\alpha_{\mathcal{S}} + \alpha_{\mu,k} + \alpha_{\mu,\ell}) + (1 - \rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}}) (p_k^{-1/2} + p_{\ell}^{-1/2}) \right]^{-1}$$

in the definition of  $\underline{\text{SR}}_{k\ell}^{\text{mc}}$ .

*Proof.* The statement is proved in Corollary D.16. □

*Remark B.6 (Active-window interpretation).* The factor  $L_{k\ell}^{\text{mc}}$  separates the class-averaging contribution from the remainder induced by incomplete active windows. Lemma D.17 gives the deterministic comparison: when initialization and signal remainders are no larger than the class-averaging term and the active-window deficits satisfy  $\Delta_k + \Delta_{\ell} = O((\ell_k p_k^{-1/2} + \ell_{\ell} p_{\ell}^{-1/2})/\sqrt{n_{\mathcal{T}}})$ , the certified ratio has the same  $\sqrt{n_{\mathcal{T}}}$  scale as class averaging. If the deficits are instead of constant order, the pairwise factor satisfies  $L_{k\ell}^{\text{mc}} = O(n_{\mathcal{T}}^{-1/2})$ , so the averaging gain is no longer reflected in the certified ratio. In the full-window case, the transition is governed by the comparison between the initialization term  $\rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}}(2\alpha_{\mathcal{S}} + \alpha_{\mu,k} + \alpha_{\mu,\ell})$  and the class-averaging term  $(1 - \rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}})(p_k^{-1/2} + p_{\ell}^{-1/2})$ .

## C Preliminary Lemmas

This section introduces several preliminary lemmas that serve as the foundation for the subsequent proofs and results.

### C.1 Bounds for Random Variables and Aggregates

Our theoretical analysis relies on controlling the behavior of various random quantities that arise from the data generation process and random initializations. This subsection establishes the foundational high-probability bounds for these quantities, starting with the core definition of sub-Gaussian vectors which characterizes our noise model.

**Definition C.1** (Sub-Gaussian with proxy covariance). A random vector  $\boldsymbol{\xi} \in \mathbb{R}^d$  is sub-Gaussian with proxy covariance  $\boldsymbol{\Sigma} \succeq \mathbf{0}$  if  $\mathbb{E} \exp(\lambda \langle \mathbf{v}, \boldsymbol{\xi} \rangle) \leq \exp(\frac{\lambda^2}{2} \mathbf{v}^\top \boldsymbol{\Sigma} \mathbf{v})$  for all  $\lambda \in \mathbb{R}$  and  $\mathbf{v} \in \mathbb{R}^d$ .

The following lemmas establish high-probability bounds for various random variables and their aggregates, crucial for controlling noise and initialization effects throughout our analysis.

**Lemma C.2.** *Assume the training set follows the additive model of Section 2.1 with sub-Gaussian noise parameter  $\boldsymbol{\Sigma}_\xi$ . Fix  $\delta \in (0, 1)$  and suppose the consolidated dimension condition  $d \geq \kappa \log(6T_{\text{out}}n_{\mathcal{T}}/\delta)$  for some  $\kappa > 0$ . There exist absolute constants  $C > 0$  and  $C_\kappa > 0$  such that, with probability at least  $1 - \delta$ , simultaneously for all  $i, i' \in [n_{\mathcal{T}}]$ ,*

$$|\langle \boldsymbol{\xi}_i, \boldsymbol{\mu} \rangle| \leq \sqrt{2 \log(6n_{\mathcal{T}}/\delta)} \|\boldsymbol{\Sigma}_\xi^{1/2} \boldsymbol{\mu}\|_2, \leq C_\kappa \sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \|\boldsymbol{\mu}\|_2} \sqrt{d}. \quad (23)$$

$$|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle| \leq C \left( \|\boldsymbol{\Sigma}_\xi\|_F \sqrt{\log(6n_{\mathcal{T}}^2/\delta)} + \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \log(6n_{\mathcal{T}}^2/\delta) \right), \quad i \neq i', \leq C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d. \quad (24)$$

$$\left| \|\boldsymbol{\xi}_i\|_2^2 - \text{tr}(\boldsymbol{\Sigma}_\xi) \right| \leq C \left( \|\boldsymbol{\Sigma}_\xi\|_F \sqrt{\log(6n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \log(6n_{\mathcal{T}}/\delta) \right), \leq C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d. \quad (25)$$

*Proof.* For equation 23, since  $\langle \boldsymbol{\mu}, \boldsymbol{\xi}_i \rangle$  is sub-Gaussian with variance proxy  $\boldsymbol{\mu}^\top \boldsymbol{\Sigma}_\xi \boldsymbol{\mu}$ , for any  $t > 0$ ,

$$\mathbb{P}(|\langle \boldsymbol{\mu}, \boldsymbol{\xi}_i \rangle| \geq t) \leq 2 \exp\left(-\frac{t^2}{2 \boldsymbol{\mu}^\top \boldsymbol{\Sigma}_\xi \boldsymbol{\mu}}\right).$$

Taking  $t = \sqrt{2 \boldsymbol{\mu}^\top \boldsymbol{\Sigma}_\xi \boldsymbol{\mu} \log(6n_{\mathcal{T}}/\delta)}$  gives  $\mathbb{P}(|\langle \boldsymbol{\mu}, \boldsymbol{\xi}_i \rangle| \geq t) \leq \delta/(3n_{\mathcal{T}})$ ; a union bound over  $i$  yields the first inequality. Using  $\boldsymbol{\mu}^\top \boldsymbol{\Sigma}_\xi \boldsymbol{\mu} \leq \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \|\boldsymbol{\mu}\|_2^2$  and  $\sqrt{\log(6n_{\mathcal{T}}/\delta)} \leq \sqrt{\log(6T_{\text{out}}n_{\mathcal{T}}/\delta)} \leq \kappa^{-1/2} \sqrt{d}$ , we obtain the second line of equation 23 after choosing  $C_\kappa \geq \kappa^{-1/2}$ .

For equation 24, write  $\boldsymbol{\xi}_i = \boldsymbol{\Sigma}_\xi^{1/2} \mathbf{z}_i$  where  $\mathbf{z}_i$  has independent, mean-zero, sub-Gaussian coordinates with unit variance. For  $i \neq i'$  we have  $\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle = \mathbf{z}_i^\top \boldsymbol{\Sigma}_\xi \mathbf{z}_{i'}$ . A decoupled Hanson–Wright type inequality for bilinear forms of independent sub-Gaussian vectors (e.g. Rudelson & Vershynin, 2013) yields a universal constant  $c > 0$  such that, for all  $t > 0$ ,

$$\mathbb{P}(|\mathbf{z}_i^\top \boldsymbol{\Sigma}_\xi \mathbf{z}_{i'}| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{\|\boldsymbol{\Sigma}_\xi\|_F^2}, \frac{t}{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}\right\}\right).$$

Choosing

$$t = C \left( \|\boldsymbol{\Sigma}_\xi\|_F \sqrt{\log(6n_{\mathcal{T}}^2/\delta)} + \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \log(6n_{\mathcal{T}}^2/\delta) \right)$$

with  $C$  large enough (depending only on  $c$ ) gives  $\mathbb{P}(|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle| \geq t) \leq \delta/(3n_{\mathcal{T}}^2)$ . A union bound over ordered pairs  $(i, i')$  with  $i \neq i'$  yields the first line of equation 24. For the coarse bound, use  $\|\boldsymbol{\Sigma}_\xi\|_F \leq \sqrt{d} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}$  and

$$\sqrt{\log(6n_{\mathcal{T}}^2/\delta)} \leq \sqrt{\log(6T_{\text{out}}n_{\mathcal{T}}^2/\delta)} \leq \sqrt{2 \log(6T_{\text{out}}n_{\mathcal{T}}/\delta)} \leq \sqrt{2/\kappa} \sqrt{d},$$

together with  $\log(6n_{\mathcal{T}}/\delta) \leq (2/\kappa)d$ , which follows from the dimension condition. Enlarging constants gives the second line of equation 24.

For equation 25, we again write  $\boldsymbol{\xi}_i = \boldsymbol{\Sigma}_{\xi}^{1/2} \mathbf{z}_i$  so that  $\|\boldsymbol{\xi}_i\|_2^2 = \boldsymbol{\zeta}_i^{\top} \boldsymbol{\Sigma}_{\xi} \boldsymbol{\zeta}_i$  and  $\mathbb{E}[\boldsymbol{\zeta}_i^{\top} \boldsymbol{\Sigma}_{\xi} \boldsymbol{\zeta}_i] = \text{tr}(\boldsymbol{\Sigma}_{\xi})$ . The Hanson–Wright inequality implies that for all  $t > 0$ ,

$$\mathbb{P}(|\boldsymbol{\zeta}_i^{\top} \boldsymbol{\Sigma}_{\xi} \boldsymbol{\zeta}_i - \text{tr}(\boldsymbol{\Sigma}_{\xi})| \geq t) \leq 2 \exp\left(-c \min\left\{\frac{t^2}{\|\boldsymbol{\Sigma}_{\xi}\|_F^2}, \frac{t}{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}}\right\}\right).$$

Taking

$$t = C\left(\|\boldsymbol{\Sigma}_{\xi}\|_F \sqrt{\log(6n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} \log(6n_{\mathcal{T}}/\delta)\right)$$

gives a tail probability at most  $\delta/(3n_{\mathcal{T}})$  for  $C$  sufficiently large. A union bound over  $i$  establishes the first line of equation 25. The second line follows from the same simplifications as above:  $\|\boldsymbol{\Sigma}_{\xi}\|_F \leq \sqrt{d} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}$  and  $\sqrt{\log(6n_{\mathcal{T}}/\delta)} \leq \kappa^{-1/2} \sqrt{d}$ ,  $\log(6n_{\mathcal{T}}/\delta) \leq \kappa^{-1}d$ .

Finally, the three parts are proved with failure probabilities at most  $\delta/3$  each, so a union bound yields a joint event of probability at least  $1 - \delta$  on which equation 23–equation 25 hold simultaneously.  $\square$

**Lemma C.3.** *Assume the training set follows the additive model of Section 2.1 with sub-Gaussian noise parameter  $\boldsymbol{\Sigma}_{\xi}$ . Let  $\mathbf{w}^{(t_{\text{out}},0)} \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{w}}^2 \mathbf{I})$  be independent of  $\{\boldsymbol{\xi}_i\}$ . Fix  $\delta \in (0, 1)$  and suppose  $d \geq \kappa \log(6T_{\text{out}}n_{\mathcal{T}}/\delta)$  for some  $\kappa > 0$ . There exist absolute constants  $C_{\kappa} > 0$ , such that with probability at least  $1 - \delta$ , simultaneously for all  $t_{\text{out}} \leq T_{\text{out}} - 1$  and  $i \in [n_{\mathcal{T}}]$ ,*

$$\|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2 \leq \sigma_{\mathbf{w}}^2 \left(d + 2\sqrt{d \log(6T_{\text{out}}/\delta)} + 2 \log(6T_{\text{out}}/\delta)\right), \quad \leq C_{\kappa} \sigma_{\mathbf{w}}^2 d. \quad (26)$$

$$|\langle \mathbf{w}^{(t_{\text{out}},0)}, \boldsymbol{\mu} \rangle| \leq \sigma_{\mathbf{w}} \|\boldsymbol{\mu}\|_2 \sqrt{2 \log(6T_{\text{out}}/\delta)}, \quad \leq C_{\kappa} \sigma_{\mathbf{w}} \|\boldsymbol{\mu}\|_2 \sqrt{d}. \quad (27)$$

$$\begin{aligned} |\langle \mathbf{w}^{(t_{\text{out}},0)}, \boldsymbol{\xi}_i \rangle| &\leq \sigma_{\mathbf{w}} \sqrt{2 \log(6T_{\text{out}}n_{\mathcal{T}}/\delta)} \\ &\quad \times \sqrt{\text{tr}(\boldsymbol{\Sigma}_{\xi}) + C\left(\|\boldsymbol{\Sigma}_{\xi}\|_F \sqrt{\log(6n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} \log(6n_{\mathcal{T}}/\delta)\right)}, \\ &\leq C_{\kappa} \sigma_{\mathbf{w}} \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d}. \end{aligned} \quad (28)$$

*Proof.* For equation 28, condition on  $\boldsymbol{\xi}_i$ :

$$\langle \mathbf{w}^{(t_{\text{out}},0)}, \boldsymbol{\xi}_i \rangle \mid \boldsymbol{\xi}_i \sim \mathcal{N}(0, \sigma_{\mathbf{w}}^2 \|\boldsymbol{\xi}_i\|_2^2),$$

so for any  $\tau > 0$ ,

$$\mathbb{P}\left(|\langle \mathbf{w}^{(t_{\text{out}},0)}, \boldsymbol{\xi}_i \rangle| \geq \tau \mid \boldsymbol{\xi}_i\right) \leq 2 \exp\left(-\frac{\tau^2}{2\sigma_{\mathbf{w}}^2 \|\boldsymbol{\xi}_i\|_2^2}\right).$$

Take

$$\tau = \sigma_{\mathbf{w}} \sqrt{2 \log(6T_{\text{out}}n_{\mathcal{T}}/\delta)} \sqrt{\text{tr}(\boldsymbol{\Sigma}_{\xi}) + C\left(\|\boldsymbol{\Sigma}_{\xi}\|_F \sqrt{\log(6n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} \log(6n_{\mathcal{T}}/\delta)\right)}.$$

By Lemma C.2, inequality equation 25 holds simultaneously for all  $i$ , hence

$$\|\boldsymbol{\xi}_i\|_2^2 \leq \text{tr}(\boldsymbol{\Sigma}_{\xi}) + C\left(\|\boldsymbol{\Sigma}_{\xi}\|_F \sqrt{\log(6n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} \log(6n_{\mathcal{T}}/\delta)\right).$$

A union bound over  $i$  and  $t_{\text{out}}$  gives the first inequality in equation 28. Finally, using  $\|\boldsymbol{\Sigma}_{\xi}\|_F \leq \sqrt{d} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}$  and the dimension condition to bound  $\sqrt{\log(6T_{\text{out}}n_{\mathcal{T}}/\delta)} \leq C_{\kappa} \sqrt{d}$  and  $\log(6n_{\mathcal{T}}/\delta) \leq C_{\kappa}d$ , we obtain the coarse bound in the second line of equation 28.  $\square$

**Lemma C.4.** Assume the training set follows the additive model of Section 2.1 with sub-Gaussian noise parameter  $\Sigma_\xi$ . Let  $\{\mathbf{x}_i^{S,(0,0)}\}_{i \in \mathcal{I}(S)}$  be initialized independently as  $\mathbf{x}_i^{S,(0,0)} \sim \mathcal{N}(\mathbf{0}, \sigma_S^2 \mathbf{I})$  and let  $\mathbf{w}^{(t_{\text{out}},0)} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$  be independent of everything else for  $t_{\text{out}} = 0, \dots, T_{\text{out}} - 1$ . Fix  $\delta \in (0, 1)$  and suppose the consolidated dimension condition  $d \geq \kappa \log(24T_{\text{out}} n_S n_T / \delta)$  for some  $\kappa > 0$ . There exist absolute constants  $C > 0$  and  $C_\kappa > 0$  such that, with probability at least  $1 - \delta$ , simultaneously for all  $i \in \mathcal{I}(S)$ ,  $j \in \mathcal{I}(T)$ , and  $t_{\text{out}} \leq T_{\text{out}} - 1$ ,

$$|\langle \mathbf{x}_i^{S,(0,0)}, \boldsymbol{\mu} \rangle| \leq \sigma_S \|\boldsymbol{\mu}\|_2 \sqrt{2 \log(24n_S/\delta)}, \quad \leq C_\kappa \sigma_S \|\boldsymbol{\mu}\|_2 \sqrt{d}. \quad (29)$$

$$\begin{aligned} |\langle \mathbf{x}_i^{S,(0,0)}, \boldsymbol{\xi}_j \rangle| &\leq \sigma_S \sqrt{2 \log(24n_S n_T / \delta)} \\ &\quad \times \sqrt{\text{tr}(\Sigma_\xi) + C \left( \|\Sigma_\xi\|_F \sqrt{\log(24n_T/\delta)} + \|\Sigma_\xi\|_{\text{op}} \log(24n_T/\delta) \right)} \\ &\leq C_\kappa \sigma_S \sqrt{\|\Sigma_\xi\|_{\text{op}} d}. \end{aligned} \quad (30)$$

$$\begin{aligned} |\langle \mathbf{x}_i^{S,(0,0)}, \mathbf{w}^{(t_{\text{out}},0)} \rangle| &\leq \sigma_S \sqrt{2 \log(24T_{\text{out}} n_S / \delta)} \\ &\quad \times \sqrt{\sigma_w^2 \left( d + 2\sqrt{d \log(24T_{\text{out}}/\delta)} + 2 \log(24T_{\text{out}}/\delta) \right)} \\ &\leq C_\kappa \sigma_S \sigma_w d. \end{aligned} \quad (31)$$

$$\|\mathbf{x}_i^{S,(0,0)}\|_2^2 \leq \sigma_S^2 \left( d + 2\sqrt{d \log(24n_S/\delta)} + 2 \log(24n_S/\delta) \right), \quad \leq C_\kappa \sigma_S^2 d. \quad (32)$$

*Proof.* We control each displayed inequality on a high-probability event and then take a union bound.

For equation 29, since  $\mathbf{x}_i^{S,(0,0)} \sim \mathcal{N}(\mathbf{0}, \sigma_S^2 \mathbf{I})$ , we have  $\langle \mathbf{x}_i^{S,(0,0)}, \boldsymbol{\mu} \rangle \sim \mathcal{N}(0, \sigma_S^2 \|\boldsymbol{\mu}\|_2^2)$  and hence, for any  $t > 0$ ,

$$\mathbb{P}\left(|\langle \mathbf{x}_i^{S,(0,0)}, \boldsymbol{\mu} \rangle| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2\sigma_S^2 \|\boldsymbol{\mu}\|_2^2}\right).$$

Taking  $t = \sigma_S \|\boldsymbol{\mu}\|_2 \sqrt{2 \log(24n_S/\delta)}$  gives

$$\mathbb{P}\left(|\langle \mathbf{x}_i^{S,(0,0)}, \boldsymbol{\mu} \rangle| \geq t\right) \leq \frac{\delta}{12n_S}.$$

A union bound over  $i \in \mathcal{I}(S)$  yields the first inequality in equation 29. Moreover,

$$\sqrt{\log(24n_S/\delta)} \leq \sqrt{\log(24T_{\text{out}} n_S n_T / \delta)} \leq \kappa^{-1/2} \sqrt{d},$$

so enlarging constants gives the second line of equation 29.

For equation 32, note that  $\|\mathbf{x}_i^{S,(0,0)}\|_2^2 / \sigma_S^2 \sim \chi_d^2$ . The Laurent–Massart inequality implies that for any  $u > 0$ ,

$$\mathbb{P}\left(\|\mathbf{x}_i^{S,(0,0)}\|_2^2 > \sigma_S^2 (d + 2\sqrt{d}u + 2u)\right) \leq e^{-u}.$$

Taking  $u = \log(24n_S/\delta)$  gives

$$\mathbb{P}\left(\|\mathbf{x}_i^{S,(0,0)}\|_2^2 > \sigma_S^2 \left( d + 2\sqrt{d \log(24n_S/\delta)} + 2 \log(24n_S/\delta) \right)\right) \leq \frac{\delta}{24n_S}.$$

A union bound over  $i \in \mathcal{I}(S)$  yields the first inequality in equation 32. Using  $\log(24n_S/\delta) \leq \kappa^{-1} d$  and  $\sqrt{\log(24n_S/\delta)} \leq \kappa^{-1/2} \sqrt{d}$  from the dimension condition, we have

$$d + 2\sqrt{d \log(24n_S/\delta)} + 2 \log(24n_S/\delta) \leq \left(1 + 2\kappa^{-1/2} + 2\kappa^{-1}\right) d,$$

which gives the second line of equation 32 after enlarging constants.

For equation 30, we first invoke Lemma C.2 with failure probability  $\delta/4$ , which ensures that with probability at least  $1 - \delta/4$ , simultaneously for all  $j \in \mathcal{I}(\mathcal{T})$ ,

$$\|\boldsymbol{\xi}_j\|_2^2 \leq \text{tr}(\boldsymbol{\Sigma}_\xi) + C \left( \|\boldsymbol{\Sigma}_\xi\|_F \sqrt{\log(24n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \log(24n_{\mathcal{T}}/\delta) \right).$$

Condition on  $\boldsymbol{\xi}_j$ . Given  $\boldsymbol{\xi}_j$ ,

$$\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \boldsymbol{\xi}_j \rangle \mid \boldsymbol{\xi}_j \sim \mathcal{N}(0, \sigma_{\mathcal{S}}^2 \|\boldsymbol{\xi}_j\|_2^2),$$

so for any  $\tau > 0$ ,

$$\mathbb{P} \left( |\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \boldsymbol{\xi}_j \rangle| \geq \tau \mid \boldsymbol{\xi}_j \right) \leq 2 \exp \left( -\frac{\tau^2}{2\sigma_{\mathcal{S}}^2 \|\boldsymbol{\xi}_j\|_2^2} \right).$$

On the event above, take

$$\tau = \sigma_{\mathcal{S}} \sqrt{2 \log(24n_{\mathcal{S}}n_{\mathcal{T}}/\delta)} \sqrt{\text{tr}(\boldsymbol{\Sigma}_\xi) + C \left( \|\boldsymbol{\Sigma}_\xi\|_F \sqrt{\log(24n_{\mathcal{T}}/\delta)} + \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \log(24n_{\mathcal{T}}/\delta) \right)}.$$

Then  $\tau^2/(2\sigma_{\mathcal{S}}^2 \|\boldsymbol{\xi}_j\|_2^2) \geq \log(24n_{\mathcal{S}}n_{\mathcal{T}}/\delta)$  and hence

$$\mathbb{P} \left( |\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \boldsymbol{\xi}_j \rangle| \geq \tau \mid \boldsymbol{\xi}_j \right) \leq \frac{\delta}{12n_{\mathcal{S}}n_{\mathcal{T}}}.$$

A union bound over  $(i, j) \in \mathcal{I}(\mathcal{S}) \times \mathcal{I}(\mathcal{T})$  gives the first inequality in equation 30. For the coarse bound, use  $\text{tr}(\boldsymbol{\Sigma}_\xi) \leq d \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}$ ,  $\|\boldsymbol{\Sigma}_\xi\|_F \leq \sqrt{d} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}$ , and  $\sqrt{\log(24n_{\mathcal{S}}n_{\mathcal{T}}/\delta)} \leq C_\kappa \sqrt{d}$ ,  $\log(24n_{\mathcal{T}}/\delta) \leq C_\kappa d$  (which follow from the dimension condition) to obtain the second line of equation 30 after enlarging constants.

For equation 31, first apply the Laurent–Massart inequality to each  $t_{\text{out}}$ :

$$\mathbb{P} \left( \|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2 > \sigma_{\mathbf{w}}^2 \left( d + 2\sqrt{d \log(24T_{\text{out}}/\delta)} + 2 \log(24T_{\text{out}}/\delta) \right) \right) \leq \frac{\delta}{24T_{\text{out}}}.$$

A union bound over  $t_{\text{out}} = 0, \dots, T_{\text{out}} - 1$  yields an event of probability at least  $1 - \delta/24$  on which equation 26–type bounds hold for all  $\mathbf{w}^{(t_{\text{out}},0)}$ . Condition on  $\mathbf{w}^{(t_{\text{out}},0)}$ . Given  $\mathbf{w}^{(t_{\text{out}},0)}$ ,

$$\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \mathbf{w}^{(t_{\text{out}},0)} \rangle \mid \mathbf{w}^{(t_{\text{out}},0)} \sim \mathcal{N}(0, \sigma_{\mathcal{S}}^2 \|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2),$$

hence for any  $\tau > 0$ ,

$$\mathbb{P} \left( |\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \mathbf{w}^{(t_{\text{out}},0)} \rangle| \geq \tau \mid \mathbf{w}^{(t_{\text{out}},0)} \right) \leq 2 \exp \left( -\frac{\tau^2}{2\sigma_{\mathcal{S}}^2 \|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2} \right).$$

On the event controlling  $\|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2$ , take

$$\tau = \sigma_{\mathcal{S}} \sqrt{2 \log(24T_{\text{out}}n_{\mathcal{S}}/\delta)} \sqrt{\sigma_{\mathbf{w}}^2 \left( d + 2\sqrt{d \log(24T_{\text{out}}/\delta)} + 2 \log(24T_{\text{out}}/\delta) \right)}.$$

Then  $\tau^2/(2\sigma_{\mathcal{S}}^2 \|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2) \geq \log(24T_{\text{out}}n_{\mathcal{S}}/\delta)$  and thus

$$\mathbb{P} \left( |\langle \mathbf{x}_i^{\mathcal{S},(0,0)}, \mathbf{w}^{(t_{\text{out}},0)} \rangle| \geq \tau \mid \mathbf{w}^{(t_{\text{out}},0)} \right) \leq \frac{\delta}{12T_{\text{out}}n_{\mathcal{S}}}.$$

A union bound over  $(i, t_{\text{out}}) \in \mathcal{I}(\mathcal{S}) \times \{0, \dots, T_{\text{out}} - 1\}$  yields the first inequality in equation 31. The coarse bound follows from the dimension condition, which implies  $\sqrt{\log(24T_{\text{out}}n_{\mathcal{S}}/\delta)} \leq C_\kappa \sqrt{d}$  and

$$\sqrt{d + 2\sqrt{d \log(24T_{\text{out}}/\delta)} + 2 \log(24T_{\text{out}}/\delta)} \leq C_\kappa \sqrt{d},$$

after enlarging constants.

Finally, combining the above high-probability events with a union bound yields a joint event of probability at least  $1 - \delta$  on which equation 29–equation 32 all hold simultaneously.  $\square$

**Lemma C.5.** Assume the training set follows the additive model of Section 2.1 with sub-Gaussian noise parameter  $\Sigma_\xi$ . Let  $\mathbf{w}^{(t_{\text{out}},0)} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I})$  be independent of  $\{\xi_i\}$  for  $t_{\text{out}} = 0, \dots, T_{\text{out}} - 1$ . Fix  $\delta \in (0, 1)$  and suppose the dimension condition

$$d \geq \kappa \log(4T_{\text{out}}/\delta) \quad \text{for some } \kappa > 0.$$

Fix, in addition, a deterministic vector  $\mathbf{a} = (a_1, \dots, a_{n_T}) \in \mathbb{R}^{n_T}$ . There exist absolute constants  $C > 0$  and  $C_\kappa > 0$  such that, with probability at least  $1 - \delta$ , the following hold simultaneously for all  $t_{\text{out}} \leq T_{\text{out}} - 1$ :

$$\begin{aligned} \left| \sum_{i \in \mathcal{I}(T)} a_i \langle \xi_i, \boldsymbol{\mu} \rangle \right| &\leq \|\mathbf{a}\|_2 \|\Sigma_\xi^{1/2} \boldsymbol{\mu}\|_2 \sqrt{2 \log(4/\delta)} \\ &\leq C_\kappa \|\mathbf{a}\|_2 \sqrt{\|\Sigma_\xi\|_{\text{op}}} \|\boldsymbol{\mu}\|_2 \sqrt{d}. \end{aligned} \quad (33)$$

$$\begin{aligned} \left| \sum_{i \in \mathcal{I}(T)} a_i \langle \xi_i, \mathbf{w}^{(t_{\text{out}},0)} \rangle \right| &\leq \|\mathbf{a}\|_2 \sigma_w \sqrt{\|\Sigma_\xi\|_{\text{op}}} \left( d + 2\sqrt{d \log(4T_{\text{out}}/\delta)} + 2 \log(4T_{\text{out}}/\delta) \right) \\ &\quad \times \sqrt{2 \log(4T_{\text{out}}/\delta)} \\ &\leq C_\kappa \|\mathbf{a}\|_2 \sigma_w \sqrt{\|\Sigma_\xi\|_{\text{op}}} d. \end{aligned} \quad (34)$$

$$\begin{aligned} \left\| \sum_{i \in \mathcal{I}(T)} a_i \xi_i \right\|_2^2 &\leq \|\mathbf{a}\|_2^2 \text{tr}(\Sigma_\xi) + C \|\mathbf{a}\|_2^2 \left( \|\Sigma_\xi\|_F \sqrt{\log(4/\delta)} + \|\Sigma_\xi\|_{\text{op}} \log(4/\delta) \right) \\ &\leq C_\kappa \|\mathbf{a}\|_2^2 \|\Sigma_\xi\|_{\text{op}} d. \end{aligned} \quad (35)$$

*Proof.* For equation 33, the variables  $\langle \xi_i, \boldsymbol{\mu} \rangle$  are independent, mean-zero, sub-Gaussian with variance proxy  $\boldsymbol{\mu}^\top \Sigma_\xi \boldsymbol{\mu}$ . Hence, for any  $t > 0$ ,

$$\mathbb{P} \left( \left| \sum_i a_i \langle \xi_i, \boldsymbol{\mu} \rangle \right| \geq t \right) \leq 2 \exp \left( -\frac{t^2}{2 \|\mathbf{a}\|_2^2 \boldsymbol{\mu}^\top \Sigma_\xi \boldsymbol{\mu}} \right).$$

Taking  $t = \|\mathbf{a}\|_2 \|\Sigma_\xi^{1/2} \boldsymbol{\mu}\|_2 \sqrt{2 \log(4/\delta)}$  gives the first line with failure probability at most  $\delta/4$ . Using  $\boldsymbol{\mu}^\top \Sigma_\xi \boldsymbol{\mu} \leq \|\Sigma_\xi\|_{\text{op}} \|\boldsymbol{\mu}\|_2^2$  and  $\sqrt{\log(4/\delta)} \leq \sqrt{\log(4T_{\text{out}}/\delta)} \leq C_\kappa \sqrt{d}$  by the dimension condition yields the second line.

For equation 34, condition on  $\mathbf{w}^{(t_{\text{out}},0)}$  and apply a sub-Gaussian tail bound to the weighted sum. This introduces a random variance term  $\mathbf{w}^{(t_{\text{out}},0)\top} \Sigma_\xi \mathbf{w}^{(t_{\text{out}},0)}$ , which we bound deterministically. First, note that  $\mathbf{w}^\top \Sigma_\xi \mathbf{w} \leq \|\Sigma_\xi\|_{\text{op}} \|\mathbf{w}\|_2^2$ . The norm term  $\|\mathbf{w}^{(t_{\text{out}},0)}\|_2^2$  is in turn bounded with high probability for all  $t_{\text{out}}$  by inequality equation 26 in Lemma C.3. By combining these two bounds and applying a union bound over  $t_{\text{out}}$  for the conditional tail event, we establish the first line of equation 34 with failure probability at most  $\delta/4$ . The second line follows from the dimension condition, which implies  $\sqrt{\log(4T_{\text{out}}/\delta)} \leq C_\kappa \sqrt{d}$ .

For equation 35, write  $\xi_i = \Sigma_\xi^{1/2} \mathbf{z}_i$  with independent mean-zero sub-Gaussian  $\mathbf{z}_i$  and set  $\mathbf{Z} := \sum_i a_i \mathbf{z}_i$ . Then  $\|\sum_i a_i \xi_i\|_2^2 = \mathbf{Z}^\top \Sigma_\xi \mathbf{Z}$ . Hanson–Wright implies that, for any  $u > 0$ ,

$$\mathbb{P} \left( \mathbf{Z}^\top \Sigma_\xi \mathbf{Z} > \|\mathbf{a}\|_2^2 \text{tr}(\Sigma_\xi) + C \|\mathbf{a}\|_2^2 \left( \|\Sigma_\xi\|_F \sqrt{u} + \|\Sigma_\xi\|_{\text{op}} u \right) \right) \leq e^{-u}.$$

Taking  $u = \log(4/\delta)$  yields the first line of equation 35 with failure probability at most  $\delta/4$ . Using  $\text{tr}(\Sigma_\xi) \leq d \|\Sigma_\xi\|_{\text{op}}$ ,  $\|\Sigma_\xi\|_F \leq \sqrt{d} \|\Sigma_\xi\|_{\text{op}}$ , and  $\log(4/\delta) \leq C_\kappa d$  concludes the second line.

The three displayed estimates hold simultaneously with probability at least  $1 - \delta$  for the fixed vector  $\mathbf{a}$ .  $\square$

## D Missing Proof

This appendix collects the proofs deferred from the main text and isolates the deterministic identities that govern the bilevel dynamics of dataset condensation.

### D.1 Notation

Throughout, the two-dimensional time index is  $\mathbf{t} = (t_{\text{out}}, t_{\text{in}})$  with  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and  $t_{\text{in}} \in \{0, \dots, T_{\text{in}}\}$ . In this appendix we index the condensed samples by their class labels and write  $\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}$  and  $\mathbf{x}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})}$  for  $y \in \{\pm 1\}$ , while  $\mathbf{w}^{\mathbf{t}}$  and  $\mathbf{x}_y^{\mathcal{S}, \mathbf{t}}$  may be used as shorthand when the tuple index is the main object.

For a dataset  $\mathcal{A} = \{(\mathbf{x}_i^{\mathcal{A}}, y_i^{\mathcal{A}})\}_{i \in \mathcal{I}(\mathcal{A})}$  with  $y_i^{\mathcal{A}} \in \{-1, +1\}$ , we write  $n_{\mathcal{A}} := |\mathcal{I}(\mathcal{A})|$  and define the signed samples

$$\mathbf{z}_i^{\mathcal{A}} := y_i^{\mathcal{A}} \mathbf{x}_i^{\mathcal{A}}, \quad i \in \mathcal{I}(\mathcal{A}).$$

As in the main text, for  $y \in \{\pm 1\}$  we write  $\mathcal{A}_y$  for the classwise block, with index set  $\mathcal{I}(\mathcal{A}_y)$  and size  $n_{\mathcal{A}_y} := |\mathcal{I}(\mathcal{A}_y)|$ .

For the training set, we denote the class proportions and the associated balance factor by

$$p_y := \frac{n_{\mathcal{T}_y}}{n_{\mathcal{T}}}, \quad \beta_p := p_+^{-1/2} + p_-^{-1/2}, \quad y \in \{\pm 1\}.$$

For the condensed set, we adopt the setting that each classwise block contains exactly one sample, i.e.,  $n_{\mathcal{S}_y} = 1$  for  $y \in \pm 1$  (hence  $n_{\mathcal{S}} = 2$ ). We index the condensed samples by their labels and define the time-dependent signed samples

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} := y \mathbf{x}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})}, \quad y \in \{\pm 1\}.$$

For the hinge loss, we use the activity indicator

$$q_{\mathcal{A}, i}(\mathbf{w}) := \mathbb{1}\{1 - \langle \mathbf{w}, \mathbf{z}_i^{\mathcal{A}} \rangle > 0\}, \quad i \in \mathcal{I}(\mathcal{A}),$$

and the regularized hinge objective

$$L_{\mathcal{A}}(\mathbf{w}) := \frac{1}{n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} \max(0, 1 - \langle \mathbf{w}, \mathbf{z}_i^{\mathcal{A}} \rangle) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2.$$

Along an inner-loop trajectory  $\{\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}\}_{t_{\text{in}} \geq 0}$ , for a fixed dataset  $\mathcal{A}$  we write

$$q_{\mathcal{A}, i}^{(t_{\text{out}}, t_{\text{in}})} := q_{\mathcal{A}, i}(\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}), \quad \bar{q}_{\mathcal{A}}^{(t_{\text{out}}, t_{\text{in}})} := \frac{1}{n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_{\mathcal{A}, i}^{(t_{\text{out}}, t_{\text{in}})} \in [0, 1],$$

and define the activation-weighted signed mean

$$\mathbf{g}_{\mathcal{A}}^{(t_{\text{out}}, t_{\text{in}})} := \frac{1}{n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_{\mathcal{A}, i}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_i^{\mathcal{A}}.$$

For the condensed set, the signed samples are time dependent and each class contains a single condensed sample. We therefore define, for each  $y \in \{\pm 1\}$ ,

$$q_{\mathcal{S}, y}^{(t_{\text{out}}, t_{\text{in}})} := \mathbb{1}\left\{1 - \langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \rangle > 0\right\},$$

and set

$$\bar{q}_{\mathcal{S}_y}^{(t_{\text{out}}, t_{\text{in}})} := q_{\mathcal{S}, y}^{(t_{\text{out}}, t_{\text{in}})} \in [0, 1], \quad \mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}}, t_{\text{in}})} := q_{\mathcal{S}, y}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})}.$$

We also denote the empirical signed mean of the training set by

$$\bar{\mathbf{z}}_{\mathcal{T}} := \frac{1}{n_{\mathcal{T}}} \sum_{i \in \mathcal{I}(\mathcal{T})} \mathbf{z}_i^{\mathcal{T}} = \frac{1}{n_{\mathcal{T}}} \sum_{i \in \mathcal{I}(\mathcal{T})} y_i^{\mathcal{T}} \mathbf{x}_i^{\mathcal{T}}.$$

## D.2 Assumptions

The following assumptions are used to prove the lemmas and the main theorems. We use the time index  $\mathbf{t} = (t_{\text{out}}, t_{\text{in}})$  from Section D.1.

**Assumption D.1.** Fix a tail parameter  $\delta \in (0, 1)$  and constants  $\kappa, C_\kappa > 0$ . We work on the intersection of the high-probability events of Lemmas C.2–C.5 in Section C.1, so that the total failure probability is at most  $\delta$ . The following conditions are assumed:

**A1. Signal strength.**

$$\frac{\|\boldsymbol{\mu}\|_2}{\sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}} \sqrt{d}}} \geq 2C_\kappa.$$

**A2. High dimension.**

$$d \geq \kappa \log\left(\frac{24T_{\text{out}} n_S n_{\mathcal{T}}}{\delta}\right).$$

**A3. Small learning rate.** The step sizes satisfy

$$0 < \eta_{\mathbf{w}} \lambda \leq 1, \quad 0 < 2\eta_S \leq 1.$$

**A4. Small initialization.** Denote  $\rho_S := 1 - 2\eta_S$ . The initialization scales satisfy

$$\sigma_{\mathbf{w}} < \frac{1}{(2C_\kappa + 1)\|\boldsymbol{\mu}\|_2 \sqrt{d}}, \quad \sigma_S < \min\left\{\frac{1}{C_\kappa \sigma_{\mathbf{w}} d}, \frac{\lambda}{(2C_\kappa + 1)\|\boldsymbol{\mu}\|_2 \sqrt{d}}, \rho_S^{T_{\text{in}}} \sqrt{\frac{2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}{n_{\mathcal{T}}}}\right\}.$$

Assumption D.1A1 implies  $\sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} d \leq (2C_\kappa)^{-1} \|\boldsymbol{\mu}\|_2 \sqrt{d}$ , and hence  $C_\kappa(\|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} d) \leq (C_\kappa + \frac{1}{2})\|\boldsymbol{\mu}\|_2 \sqrt{d}$ . Thus the initialization scale in Assumption D.1A4 implies the mixed-scale small-margin bounds used below.

## D.3 Auxiliary Results for Lemma 3.2

The first auxiliary result gives algebraic identities for the inner-loop hinge recursion. Independently of any switching pattern of the activity indicators, the weight sequence admits an unconditional unrolling in terms of the activation-weighted signed means. On intervals where the activity indicators are fixed, the same identity reduces to an affine recursion with a closed-form solution. Fix  $\rho := 1 - \lambda\eta_{\mathbf{w}}$ .

**Lemma D.2.** Fix an outer-loop index  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$ . Let  $\mathcal{A} = \{(\mathbf{x}_i^A, y_i^A)\}_{i \in \mathcal{I}(\mathcal{A})}$  be a finite dataset. Consider the inner-loop update with step size  $\eta_{\mathbf{w}} > 0$ ,

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}}+1)} = (1 - \lambda\eta_{\mathbf{w}})\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} + \frac{\eta_{\mathbf{w}}}{n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_{\mathcal{A}, i}(\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}) \mathbf{z}_i^A. \quad (36)$$

Then the following hold.

(i) For every integer  $t_{\text{in}} \geq 0$ ,

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}}, 0)} + \eta_{\mathbf{w}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \mathbf{g}_{\mathcal{A}}^{(t_{\text{out}}, s)}, \quad (37)$$

where the sum is interpreted as 0 when  $t_{\text{in}} = 0$ . Moreover, if there exist constants  $\{q_i^*\}_{i \in \mathcal{I}(\mathcal{A})} \subset \{0, 1\}$  such that  $q_{\mathcal{A}, i}(\mathbf{w}^{(t_{\text{out}}, s)}) = q_i^*$  for all  $i \in \mathcal{I}(\mathcal{A})$  and all  $s \in \{0, 1, \dots, t_{\text{in}} - 1\}$ , then

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} = (1 - \lambda\eta_{\mathbf{w}})^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}}, 0)} + \frac{1 - (1 - \lambda\eta_{\mathbf{w}})^{t_{\text{in}}}}{\lambda n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_i^* \mathbf{z}_i^A. \quad (38)$$

*Proof.* For (i), rewrite equation 36 as

$$\mathbf{w}^{(t_{\text{out}},s+1)} = \rho \mathbf{w}^{(t_{\text{out}},s)} + \eta_{\mathbf{w}} \mathbf{g}_{\mathcal{A}}^{(t_{\text{out}},s)}, \quad s \geq 0.$$

If  $t_{\text{in}} = 0$  then equation 37 is immediate. For  $t_{\text{in}} \geq 1$ , multiply the display by  $\rho^{t_{\text{in}}-1-s}$  and sum over  $s = 0, \dots, t_{\text{in}} - 1$ :

$$\sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \mathbf{w}^{(t_{\text{out}},s+1)} = \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-s} \mathbf{w}^{(t_{\text{out}},s)} + \eta_{\mathbf{w}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \mathbf{g}_{\mathcal{A}}^{(t_{\text{out}},s)}.$$

Re-index the left-hand side with  $u := s + 1$ :

$$\sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \mathbf{w}^{(t_{\text{out}},s+1)} = \sum_{u=1}^{t_{\text{in}}} \rho^{t_{\text{in}}-u} \mathbf{w}^{(t_{\text{out}},u)}.$$

Similarly,

$$\sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-s} \mathbf{w}^{(t_{\text{out}},s)} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + \sum_{u=1}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-u} \mathbf{w}^{(t_{\text{out}},u)}.$$

Subtracting these two expressions yields

$$\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} - \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} = \eta_{\mathbf{w}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \mathbf{g}_{\mathcal{A}}^{(t_{\text{out}},s)},$$

which is equation 37.

Under the constancy assumption, we have

$$\mathbf{g}_{\mathcal{A}}^{(t_{\text{out}},s)} = \frac{1}{n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_i^* \mathbf{z}_i^{\mathcal{A}} \quad \text{for all } s \in \{0, 1, \dots, t_{\text{in}} - 1\}.$$

Substituting this into equation 37 gives

$$\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + \frac{\eta_{\mathbf{w}}}{n_{\mathcal{A}}} \left( \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \right) \sum_{i \in \mathcal{I}(\mathcal{A})} q_i^* \mathbf{z}_i^{\mathcal{A}}.$$

The geometric series satisfies

$$\sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} = \sum_{k=0}^{t_{\text{in}}-1} \rho^k = \frac{1 - \rho^{t_{\text{in}}}}{1 - \rho} = \frac{1 - \rho^{t_{\text{in}}}}{\lambda \eta_{\mathbf{w}}}.$$

Substituting this identity yields

$$\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + \frac{1 - \rho^{t_{\text{in}}}}{\lambda n_{\mathcal{A}}} \sum_{i \in \mathcal{I}(\mathcal{A})} q_i^* \mathbf{z}_i^{\mathcal{A}}.$$

Replacing  $\rho$  by  $1 - \lambda \eta_{\mathbf{w}}$  gives equation 38.  $\square$

The next result gives deterministic geometric bounds for the training set that will be used to certify a fully-active linear window. Under the high-probability events in Assumption D.1, it controls the initialization margins and the pairwise inner products of the signed training samples. Define

$$B_{\mu} := C_{\kappa} \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}} \|\boldsymbol{\mu}\|_2 \sqrt{d}, \quad B_{\text{init}} := C_{\kappa} \sigma_{\mathbf{w}} \left( \|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}} d \right), \quad B_{\xi\xi} := C_{\kappa} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d,$$

and set

$$G_{\min} := \|\boldsymbol{\mu}\|_2^2 - 2B_{\mu} - B_{\xi\xi}, \quad G_{\max} := \|\boldsymbol{\mu}\|_2^2 + 2B_{\mu} + B_{\xi\xi}.$$

**Lemma D.3.** Assume Assumption D.1 and the additive model  $\mathbf{x}_i^T = y_i^T \boldsymbol{\mu} + \boldsymbol{\xi}_i$ . Then the following hold for every  $(t_{\text{out}}, t_{\text{in}})$ .

(i) The initialization margins satisfy, for every  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and every  $i \in \mathcal{I}(\mathcal{T})$ ,

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_i^T \rangle| \leq B_{\text{init}}.$$

(ii) The signed training samples satisfy the pairwise inner-product bounds

$$G_{\text{min}} \leq \langle \mathbf{z}_i^T, \mathbf{z}_{i'}^T \rangle \leq G_{\text{max}}, \quad i, i' \in \mathcal{I}(\mathcal{T}).$$

In particular, for every  $i \in \mathcal{I}(\mathcal{T})$ ,

$$G_{\text{min}} \leq \langle \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_i^T \rangle \leq G_{\text{max}}.$$

*Proof.* For (i), the additive model gives

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{x}_i^T \rangle| \leq |\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \boldsymbol{\mu} \rangle| + |\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \boldsymbol{\xi}_i \rangle|.$$

Lemma C.3 equation 27 and equation 28 hold simultaneously for all  $t_{\text{out}}$  and  $i$ , hence

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{x}_i^T \rangle| \leq C_{\kappa} \sigma_{\mathbf{w}} \left( \|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d} \right).$$

Since  $\mathbf{z}_i^T = y_i^T \mathbf{x}_i^T$  and  $|y_i^T| = 1$ , the same bound holds with  $\mathbf{x}_i^T$  replaced by  $\mathbf{z}_i^T$ , which proves (i).

For (ii), fix  $i, i' \in \mathcal{I}(\mathcal{T})$ . Under the additive model,  $\mathbf{z}_i^T = \boldsymbol{\mu} + y_i^T \boldsymbol{\xi}_i$ , hence

$$\langle \mathbf{z}_i^T, \mathbf{z}_{i'}^T \rangle = \|\boldsymbol{\mu}\|_2^2 + \langle \boldsymbol{\mu}, y_i^T \boldsymbol{\xi}_i \rangle + \langle \boldsymbol{\mu}, y_{i'}^T \boldsymbol{\xi}_{i'} \rangle + y_i^T y_{i'}^T \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle.$$

Under Assumption D.1, Lemma C.2 equation 23 gives  $\max_j |\langle \boldsymbol{\mu}, \boldsymbol{\xi}_j \rangle| \leq B_{\mu}$ . Moreover, Lemma C.2 equation 24 bounds  $|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle|$  for  $i \neq i'$ , and Lemma C.2 equation 25 together with  $\text{tr}(\boldsymbol{\Sigma}_{\xi}) \leq \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d$  bounds  $\|\boldsymbol{\xi}_i\|_2^2 = \langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_i \rangle$ . Enlarging  $C_{\kappa}$  once if needed, we obtain  $\max_{i, i' \in \mathcal{I}(\mathcal{T})} |\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle| \leq B_{\xi\xi}$ . Combining the last displays yields

$$\langle \mathbf{z}_i^T, \mathbf{z}_{i'}^T \rangle \leq \|\boldsymbol{\mu}\|_2^2 + 2B_{\mu} + B_{\xi\xi} = G_{\text{max}}, \quad \langle \mathbf{z}_i^T, \mathbf{z}_{i'}^T \rangle \geq \|\boldsymbol{\mu}\|_2^2 - 2B_{\mu} - B_{\xi\xi} = G_{\text{min}}.$$

Finally, since  $\bar{\mathbf{z}}_{\mathcal{T}} = \frac{1}{n_{\mathcal{T}}} \sum_{j \in \mathcal{I}(\mathcal{T})} \mathbf{z}_j^T$ , we have

$$\langle \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_i^T \rangle = \frac{1}{n_{\mathcal{T}}} \sum_{j \in \mathcal{I}(\mathcal{T})} \langle \mathbf{z}_j^T, \mathbf{z}_i^T \rangle,$$

so the same two-sided bounds carry over to  $\langle \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_i^T \rangle$  by averaging.  $\square$

By combining the initialization and inner-product bounds in Lemma D.3 with the closed-form dynamics in Lemma D.2, we establish the existence of an initial linear window during which all training samples remain active.

**Corollary D.4.** Assume Assumption D.1 and the additive model  $\mathbf{x}_i^T = y_i^T \boldsymbol{\mu} + \boldsymbol{\xi}_i$ . Recall  $B_{\text{init}}$  and  $G_{\text{max}}$  from Lemma D.3. Define the linear window length

$$T_{\text{lin}} := \max \left\{ t \in \{0, 1, \dots, T_{\text{in}}\} \mid \rho^t B_{\text{init}} + \frac{1 - \rho^t}{\lambda} G_{\text{max}} < 1 \right\}, \quad (39)$$

with the convention  $\max \emptyset = -1$ . Then the following hold.

(i) For every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{lin}}\}$  and every  $i \in \mathcal{I}(\mathcal{T})$ ,

$$q_{\mathcal{T}, i}(\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}) = 1,$$

and consequently

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}}, 0)} + \frac{1 - \rho^{t_{\text{in}}}}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}. \quad (40)$$

(ii) If  $G_{\text{max}} < \lambda$ , then  $T_{\text{lin}} = T_{\text{in}}$ , implying that full activity persists throughout the entire inner loop.

*Proof.* For (i), we proceed by induction on  $t_{\text{in}}$  to show that  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t_{\text{in}})}) = 1$  for all  $i \in \mathcal{I}(\mathcal{T})$ . For the base case  $t_{\text{in}} = 0$ , Lemma D.3(i) ensures  $|\langle \mathbf{w}^{(t_{\text{out}},0)}, \mathbf{z}_i^{\mathcal{T}} \rangle| \leq B_{\text{init}}$ , which is strictly less than 1 by Assumptions D.1A1 and D.1A4, implying  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},0)}) = 1$ .

Assume inductively that  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},s)}) = 1$  for all  $i$  and all  $s < t_{\text{in}}$ , with  $t_{\text{in}} \leq T_{\text{lin}}$ . Lemma D.2(i) then yields the closed-form update  $\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + \frac{1-\rho^{t_{\text{in}}}}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}$ . Decomposing the margin  $\langle \mathbf{w}^{(t_{\text{out}},t_{\text{in}})}, \mathbf{z}_i^{\mathcal{T}} \rangle$  and applying Lemma D.3(i) and Lemma D.3(ii) yields

$$\langle \mathbf{w}^{(t_{\text{out}},t_{\text{in}})}, \mathbf{z}_i^{\mathcal{T}} \rangle \leq \rho^{t_{\text{in}}} B_{\text{init}} + \frac{1-\rho^{t_{\text{in}}}}{\lambda} G_{\text{max}}.$$

Define the affine map  $\psi(\alpha) := \alpha B_{\text{init}} + (1-\alpha)G_{\text{max}}/\lambda$ . The right-hand side equals  $\psi(\rho^{t_{\text{in}}})$ . Since  $\psi$  is monotone on  $[0, 1]$ , we have  $\psi(\rho^{t_{\text{in}}}) \leq \max\{\psi(1), \psi(\rho^{T_{\text{lin}}})\} < 1$  by Assumption D.1 and the definition equation 39. Therefore  $\langle \mathbf{w}^{(t_{\text{out}},t_{\text{in}})}, \mathbf{z}_i^{\mathcal{T}} \rangle < 1$ , so  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t_{\text{in}})}) = 1$ , and equation 40 holds.

For (ii), if  $G_{\text{max}} < \lambda$ , then  $\psi(0) = G_{\text{max}}/\lambda < 1$  and  $\psi(1) = B_{\text{init}} < 1$ , hence  $\psi(\alpha) < 1$  for all  $\alpha \in [0, 1]$ . Therefore the defining inequality in equation 39 holds for all  $t \in \{0, \dots, T_{\text{in}}\}$ , yielding  $T_{\text{lin}} = T_{\text{in}}$ .  $\square$

The next result expands the inner-loop weights as an explicit linear combination of signed training samples. The coefficients encode the cumulative support-vector activity and satisfy deterministic bounds controlled by the initial fully-active window.

Fix  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and set  $\rho := 1 - \lambda \eta_{\mathbf{w}}$ . For  $t_{\text{in}} \geq 0$ , define the coefficient vector  $\mathbf{a}^{(t_{\text{out}},t_{\text{in}})} \in \mathbb{R}^{n_{\mathcal{T}}}$  by

$$a_i^{(t_{\text{out}},t_{\text{in}})} := \frac{\eta_{\mathbf{w}}}{n_{\mathcal{T}}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} q_{\mathcal{T},i}^{(t_{\text{out}},s)}, \quad i \in \mathcal{I}(\mathcal{T}), \quad (41)$$

with the convention that the sum is 0 when  $t_{\text{in}} = 0$ , and set

$$a_{\mu}^{(t_{\text{out}},t_{\text{in}})} := \sum_{i \in \mathcal{I}(\mathcal{T})} a_i^{(t_{\text{out}},t_{\text{in}})} = \eta_{\mathbf{w}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} \bar{q}_{\mathcal{T}}^{(t_{\text{out}},s)}. \quad (42)$$

**Lemma D.5.** *Assume Assumption D.1 and the additive model  $\mathbf{x}_i^{\mathcal{T}} = y_i^{\mathcal{T}} \boldsymbol{\mu} + \boldsymbol{\xi}_i$ . Fix  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$ . Then the following hold.*

(i) *For every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ , the inner-loop weights admit the signed-sample expansion*

$$\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + \sum_{i \in \mathcal{I}(\mathcal{T})} a_i^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_i^{\mathcal{T}}, \quad (43)$$

and hence, under the additive model,

$$\mathbf{w}^{(t_{\text{out}},t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}},0)} + a_{\mu}^{(t_{\text{out}},t_{\text{in}})} \boldsymbol{\mu} + \sum_{i \in \mathcal{I}(\mathcal{T})} a_i^{(t_{\text{out}},t_{\text{in}})} y_i^{\mathcal{T}} \boldsymbol{\xi}_i. \quad (44)$$

(ii) *The coefficients are nonnegative. Let  $T_{\text{lin}} \in \{0, 1, \dots, T_{\text{in}}\}$  be as in Corollary D.4 and assume its fully-active conclusion holds on  $\{0, 1, \dots, T_{\text{lin}} - 1\}$ . Define  $(u)_+ := \max\{u, 0\}$ . Then for every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$  and every  $i \in \mathcal{I}(\mathcal{T})$ ,*

$$\frac{\rho^{(t_{\text{in}}-T_{\text{lin}})_+} (1 - \rho^{\min\{t_{\text{in}}, T_{\text{lin}}\}})}{\lambda n_{\mathcal{T}}} \leq a_i^{(t_{\text{out}},t_{\text{in}})} \leq \frac{1 - \rho^{t_{\text{in}}}}{\lambda n_{\mathcal{T}}}. \quad (45)$$

Moreover,

$$\frac{\rho^{(t_{\text{in}}-T_{\text{lin}})_+} (1 - \rho^{\min\{t_{\text{in}}, T_{\text{lin}}\}})}{\lambda} \leq a_{\mu}^{(t_{\text{out}},t_{\text{in}})} \leq \frac{1 - \rho^{t_{\text{in}}}}{\lambda}, \quad (46)$$

and

$$\frac{\rho^{(t_{\text{in}}-T_{\text{lin}})_+} (1 - \rho^{\min\{t_{\text{in}}, T_{\text{lin}}\}})}{\lambda \sqrt{n_{\mathcal{T}}}} \leq \|\mathbf{a}^{(t_{\text{out}},t_{\text{in}})}\|_2 \leq \frac{1 - \rho^{t_{\text{in}}}}{\lambda \sqrt{n_{\mathcal{T}}}}. \quad (47)$$

*Proof.* For (i), this is a direct consequence of Lemma D.2(i) applied with  $\mathcal{A} = \mathcal{T}$ : substitute  $\mathbf{g}_{\mathcal{T}}^{(t_{\text{out}},s)} = \frac{1}{n_{\mathcal{T}}} \sum_{i \in \mathcal{I}(\mathcal{T})} q_{\mathcal{T},i}^{(t_{\text{out}},s)} \mathbf{z}_i^{\mathcal{T}}$  into equation 37, exchange the order of summation, and identify the resulting coefficients with equation 41, which gives equation 43. Summing equation 41 over  $i \in \mathcal{I}(\mathcal{T})$  yields equation 42, and substituting  $\mathbf{z}_i^{\mathcal{T}} = \boldsymbol{\mu} + y_i^{\mathcal{T}} \boldsymbol{\xi}_i$  gives equation 44.

For (ii), since  $q_{\mathcal{T},i}^{(t_{\text{out}},s)} \in \{0, 1\}$ , we have  $a_i^{(t_{\text{out}},t_{\text{in}})} \geq 0$  for all  $i$ . Under the fully-active conclusion on  $\{0, \dots, T_{\text{lin}} - 1\}$ , we have  $q_{\mathcal{T},i}^{(t_{\text{out}},s)} = 1$  for all  $i$  and all  $s \leq \min\{t_{\text{in}}, T_{\text{lin}}\} - 1$ , hence

$$a_i^{(t_{\text{out}},t_{\text{in}})} \geq \frac{\eta_{\mathbf{w}}}{n_{\mathcal{T}}} \sum_{s=0}^{\min\{t_{\text{in}}, T_{\text{lin}}\}-1} \rho^{t_{\text{in}}-1-s} = \frac{\eta_{\mathbf{w}}}{n_{\mathcal{T}}} \rho^{(t_{\text{in}}-T_{\text{lin}})_+} \sum_{u=0}^{\min\{t_{\text{in}}, T_{\text{lin}}\}-1} \rho^u.$$

Using  $\sum_{u=0}^{m-1} \rho^u = (1 - \rho^m)/(1 - \rho)$  and  $1 - \rho = \lambda \eta_{\mathbf{w}}$  yields the lower bound in equation 45. For the upper bound, use  $q_{\mathcal{T},i}^{(t_{\text{out}},s)} \leq 1$  in equation 41 to obtain

$$a_i^{(t_{\text{out}},t_{\text{in}})} \leq \frac{\eta_{\mathbf{w}}}{n_{\mathcal{T}}} \sum_{s=0}^{t_{\text{in}}-1} \rho^{t_{\text{in}}-1-s} = \frac{1 - \rho^{t_{\text{in}}}}{\lambda n_{\mathcal{T}}}.$$

Summing equation 45 over  $i$  gives equation 46. The  $\ell_2$  bounds in equation 47 follow from  $\|\mathbf{a}\|_2 \geq \sqrt{n_{\mathcal{T}}} \min_i a_i$  and  $\|\mathbf{a}\|_2 \leq \sqrt{n_{\mathcal{T}}} \max_i a_i$  together with equation 45.  $\square$

The corresponding expansion for the condensed samples follows from the affine gradient-matching update. It represents each evolving condensed sample as a linear combination of its initialization and the accumulated training gradients, and it gives coefficient bounds determined by the training and condensed-set activity windows. Define  $\rho_{\mathcal{S}} := 1 - 2\eta_{\mathcal{S}}$ .

**Lemma D.6.** *Fix an outer-loop index  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and a class  $y \in \{\pm 1\}$ . Let  $\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t_{\text{in}})}$  denote the time-dependent signed condensed sample in class  $y$ , and write  $q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} \in \{0, 1\}$  for its activity indicator. For the training set, recall  $\mathbf{z}_j^{\mathcal{T}}$  and  $q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})}$  for  $j \in \mathcal{I}(\mathcal{T}_y)$ .*

(i) *Define the coefficient sequences  $\{u_y^{(t_{\text{out}},t_{\text{in}})}\}_{t_{\text{in}} \geq 0}$  and  $\{c_{y,j}^{(t_{\text{out}},t_{\text{in}})}\}_{t_{\text{in}} \geq 0}$  for each  $j \in \mathcal{I}(\mathcal{T}_y)$  by the recursions*

$$u_y^{(t_{\text{out}},0)} := 1, \quad u_y^{(t_{\text{out}},t_{\text{in}}+1)} := (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) u_y^{(t_{\text{out}},t_{\text{in}})}, \quad t_{\text{in}} \geq 0, \quad (48)$$

and

$$c_{y,j}^{(t_{\text{out}},0)} := 0, \quad c_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} := (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) c_{y,j}^{(t_{\text{out}},t_{\text{in}})} + \frac{2\eta_{\mathcal{S}}}{n_{\mathcal{T}_y}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})}, \quad t_{\text{in}} \geq 0. \quad (49)$$

Set

$$c_{y,\mu}^{(t_{\text{out}},t_{\text{in}})} := \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}},t_{\text{in}})}, \quad \mathbf{c}_y^{(t_{\text{out}},t_{\text{in}})} := (c_{y,j}^{(t_{\text{out}},t_{\text{in}})})_{j \in \mathcal{I}(\mathcal{T}_y)}. \quad (50)$$

Then for every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{lin}}\}$ ,

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t_{\text{in}})} = u_y^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}},0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}. \quad (51)$$

Moreover, if the additive model  $\mathbf{x}_j^{\mathcal{T}} = y_j^{\mathcal{T}} \boldsymbol{\mu} + \boldsymbol{\xi}_j$  holds, then for every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{lin}}\}$ ,

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t_{\text{in}})} = u_y^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}},0)} + c_{y,\mu}^{(t_{\text{out}},t_{\text{in}})} \boldsymbol{\mu} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}},t_{\text{in}})} y_j^{\mathcal{T}} \boldsymbol{\xi}_j. \quad (52)$$

(ii) *Let  $T_{\text{lin}} \in \{0, 1, \dots, T_{\text{in}}\}$  and  $T_{\mathcal{S},\text{lin}}(t_{\text{out}}) \in \{0, 1, \dots, T_{\text{in}}\}$  be such that*

$$q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})} = 1 \quad \text{for all } j \in \mathcal{I}(\mathcal{T}_y) \text{ and all } t_{\text{in}} \in \{0, 1, \dots, T_{\text{lin}}-1\}, \quad q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} = 1 \quad \text{for all } t_{\text{in}} \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}})-1\}.$$

Define  $T_{y,\text{lin}} := \min\{T_{\text{lin}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}$  and  $(u)_+ := \max\{u, 0\}$ . Then for every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ ,

$$\rho_{\mathcal{S}}^{t_{\text{in}}} \leq u_y^{(t_{\text{out}}, t_{\text{in}})} \leq \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}}, \quad u_y^{(t_{\text{out}}, t_{\text{in}})} = \rho_{\mathcal{S}}^{t_{\text{in}}} \text{ for all } t_{\text{in}} \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}, \quad (53)$$

and

$$\rho_{\mathcal{S}}^{(t_{\text{in}} - T_{y,\text{lin}})_+} (1 - \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{y,\text{lin}}\}}) \leq c_{y,\mu}^{(t_{\text{out}}, t_{\text{in}})} \leq 1 - \rho_{\mathcal{S}}^{t_{\text{in}}}. \quad (54)$$

Moreover, for every  $j \in \mathcal{I}(\mathcal{T}_y)$  and every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ ,

$$\frac{\rho_{\mathcal{S}}^{(t_{\text{in}} - T_{y,\text{lin}})_+} (1 - \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{y,\text{lin}}\}})}{n_{\mathcal{T}_y}} \leq c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \leq \frac{1 - \rho_{\mathcal{S}}^{t_{\text{in}}}}{n_{\mathcal{T}_y}}, \quad c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} = \frac{1 - \rho_{\mathcal{S}}^{t_{\text{in}}}}{n_{\mathcal{T}_y}} \text{ for all } t_{\text{in}} \in \{0, 1, \dots, T_{y,\text{lin}}\}. \quad (55)$$

In particular,

$$\frac{\rho_{\mathcal{S}}^{(t_{\text{in}} - T_{y,\text{lin}})_+} (1 - \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{y,\text{lin}}\}})}{\sqrt{n_{\mathcal{T}_y}}} \leq \|\mathbf{c}_y^{(t_{\text{out}}, t_{\text{in}})}\|_2 \leq \frac{1 - \rho_{\mathcal{S}}^{t_{\text{in}}}}{\sqrt{n_{\mathcal{T}_y}}}. \quad (56)$$

*Proof.* For (i), the signed condensed update in class  $y$  takes the form

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}}+1)} = \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} + 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \left( \mathbf{g}_{\mathcal{T}_y}^{(t_{\text{out}}, t_{\text{in}})} - \mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}}, t_{\text{in}})} \right).$$

In the one-sample-per-class setting,  $\mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}}, t_{\text{in}})} = q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})}$ , and by definition of the activation-weighted mean on the training block,

$$\mathbf{g}_{\mathcal{T}_y}^{(t_{\text{out}}, t_{\text{in}})} = \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} q_{\mathcal{T},j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}.$$

Substituting both identities and using  $(q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})})^2 = q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})}$  yields

$$\begin{aligned} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}}+1)} &= \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} + \frac{2\eta_{\mathcal{S}}}{n_{\mathcal{T}_y}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} q_{\mathcal{T},j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}} - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} \\ &= (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})}) \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} + \frac{2\eta_{\mathcal{S}}}{n_{\mathcal{T}_y}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} q_{\mathcal{T},j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}. \end{aligned}$$

Assume that equation 51 holds at time  $t_{\text{in}}$ , i.e.,  $\mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} = u_y^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, 0)} + \sum_j c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}$ . Substituting this into the previous display gives

$$\begin{aligned} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}}+1)} &= (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})}) u_y^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, 0)} \\ &\quad + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \left[ (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})}) c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} + \frac{2\eta_{\mathcal{S}}}{n_{\mathcal{T}_y}} q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} q_{\mathcal{T},j}^{(t_{\text{out}}, t_{\text{in}})} \right] \mathbf{z}_j^{\mathcal{T}}. \end{aligned}$$

The coefficient of  $\mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, 0)}$  is exactly  $u_y^{(t_{\text{out}}, t_{\text{in}}+1)}$  in equation 48, and the coefficient of each  $\mathbf{z}_j^{\mathcal{T}}$  is exactly  $c_{y,j}^{(t_{\text{out}}, t_{\text{in}}+1)}$  in equation 49. With the initialization  $u_y^{(t_{\text{out}}, 0)} = 1$  and  $c_{y,j}^{(t_{\text{out}}, 0)} = 0$ , this establishes equation 51 for all  $t_{\text{in}} \in \{0, \dots, T_{\text{in}}\}$ . Under the additive model,  $\mathbf{z}_j^{\mathcal{T}} = \boldsymbol{\mu} + y_j^{\mathcal{T}} \boldsymbol{\xi}_j$ , so

$$\sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}} = \left( \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \right) \boldsymbol{\mu} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} y_j^{\mathcal{T}} \boldsymbol{\xi}_j,$$

and equation 52 follows from equation 50.

For (ii), iterating equation 48 gives the product representation

$$u_y^{(t_{\text{out}}, t_{\text{in}})} = \prod_{s=0}^{t_{\text{in}}-1} (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}}, s)}), \quad t_{\text{in}} \geq 1,$$

with the empty product interpreted as 1 at  $t_{\text{in}} = 0$ . Since  $q_{\mathcal{S},y}^{(t_{\text{out}},s)} \in \{0, 1\}$  and  $\rho_{\mathcal{S}} = 1 - 2\eta_{\mathcal{S}}$ , each factor lies in  $[\rho_{\mathcal{S}}, 1]$ , so

$$u_y^{(t_{\text{out}},t_{\text{in}})} \geq \rho_{\mathcal{S}}^{t_{\text{in}}}.$$

Moreover,  $q_{\mathcal{S},y}^{(t_{\text{out}},s)} = 1$  for  $s \in \{0, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}}) - 1\}$  implies that the first  $\min\{t_{\text{in}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}$  factors equal  $\rho_{\mathcal{S}}$ , while the remaining factors are at most 1. Thus

$$u_y^{(t_{\text{out}},t_{\text{in}})} \leq \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}}, \quad u_y^{(t_{\text{out}},t_{\text{in}})} = \rho_{\mathcal{S}}^{t_{\text{in}}} \text{ for all } t_{\text{in}} \leq T_{\mathcal{S},\text{lin}}(t_{\text{out}}),$$

which is equation 53.

Next, equation 49 and the nonnegativity of the driving term imply  $c_{y,j}^{(t_{\text{out}},t_{\text{in}})} \geq 0$  for all  $j$  and  $t_{\text{in}}$ . On the joint linear window  $\{0, 1, \dots, T_{y,\text{lin}} - 1\}$ , we have  $q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} = 1$  and  $q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})} = 1$  for every  $j \in \mathcal{I}(\mathcal{T}_y)$ , so equation 49 reduces to

$$c_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} = \rho_{\mathcal{S}} c_{y,j}^{(t_{\text{out}},t_{\text{in}})} + \frac{1 - \rho_{\mathcal{S}}}{n_{\mathcal{T}_y}}, \quad c_{y,j}^{(t_{\text{out}},0)} = 0,$$

which yields  $c_{y,j}^{(t_{\text{out}},t_{\text{in}})} = (1 - \rho_{\mathcal{S}}^{t_{\text{in}}})/n_{\mathcal{T}_y}$  for all  $t_{\text{in}} \leq T_{y,\text{lin}}$ . Summing over  $j$  gives  $c_{y,\mu}^{(t_{\text{out}},t_{\text{in}})} = 1 - \rho_{\mathcal{S}}^{t_{\text{in}}}$  for  $t_{\text{in}} \leq T_{y,\text{lin}}$ .

For an upper bound valid for all  $t_{\text{in}}$ , fix  $j \in \mathcal{I}(\mathcal{T}_y)$  and define

$$d_{y,j}^{(t_{\text{out}},t_{\text{in}})} := u_y^{(t_{\text{out}},t_{\text{in}})} + n_{\mathcal{T}_y} c_{y,j}^{(t_{\text{out}},t_{\text{in}})}.$$

Using equation 48–equation 49, for each  $t_{\text{in}} \geq 0$  we have

$$\begin{aligned} d_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} &= (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) u_y^{(t_{\text{out}},t_{\text{in}})} + n_{\mathcal{T}_y} (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) c_{y,j}^{(t_{\text{out}},t_{\text{in}})} + 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})} \\ &= (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) d_{y,j}^{(t_{\text{out}},t_{\text{in}})} + 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} q_{\mathcal{T},j}^{(t_{\text{out}},t_{\text{in}})} \\ &\leq (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) d_{y,j}^{(t_{\text{out}},t_{\text{in}})} + 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}. \end{aligned}$$

If  $q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} = 0$ , the last line gives  $d_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} \leq d_{y,j}^{(t_{\text{out}},t_{\text{in}})}$ . If  $q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} = 1$ , then it reads  $d_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} \leq \rho_{\mathcal{S}} d_{y,j}^{(t_{\text{out}},t_{\text{in}})} + (1 - \rho_{\mathcal{S}})$ . Since  $d_{y,j}^{(t_{\text{out}},0)} = u_y^{(t_{\text{out}},0)} = 1$ , both cases imply  $d_{y,j}^{(t_{\text{out}},t_{\text{in}})} \leq 1$  for all  $t_{\text{in}} \in \{0, \dots, T_{\text{in}}\}$ . Therefore,

$$c_{y,j}^{(t_{\text{out}},t_{\text{in}})} \leq \frac{1 - u_y^{(t_{\text{out}},t_{\text{in}})}}{n_{\mathcal{T}_y}} \leq \frac{1 - \rho_{\mathcal{S}}^{t_{\text{in}}}}{n_{\mathcal{T}_y}},$$

which is the upper bound in equation 55. Summing over  $j$  yields  $c_{y,\mu}^{(t_{\text{out}},t_{\text{in}})} \leq 1 - u_y^{(t_{\text{out}},t_{\text{in}})} \leq 1 - \rho_{\mathcal{S}}^{t_{\text{in}}}$ , which is the upper bound in equation 54.

For the lower bounds beyond the joint linear window, note that equation 49 gives

$$c_{y,j}^{(t_{\text{out}},t_{\text{in}}+1)} \geq (1 - 2\eta_{\mathcal{S}} q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})}) c_{y,j}^{(t_{\text{out}},t_{\text{in}})} \geq \rho_{\mathcal{S}} c_{y,j}^{(t_{\text{out}},t_{\text{in}})}.$$

Iterating this inequality from  $T_{y,\text{lin}}$  to  $t_{\text{in}}$  yields, for  $t_{\text{in}} \geq T_{y,\text{lin}}$ ,

$$c_{y,j}^{(t_{\text{out}},t_{\text{in}})} \geq \rho_{\mathcal{S}}^{t_{\text{in}} - T_{y,\text{lin}}} c_{y,j}^{(t_{\text{out}},T_{y,\text{lin}})} = \frac{\rho_{\mathcal{S}}^{t_{\text{in}} - T_{y,\text{lin}}} (1 - \rho_{\mathcal{S}}^{T_{y,\text{lin}}})}{n_{\mathcal{T}_y}},$$

which matches the unified lower bound in equation 55. Summing over  $j$  gives the lower bound in equation 54.

Finally, since all coordinates of  $\mathbf{c}_y^{(t_{\text{out}},t_{\text{in}})}$  are nonnegative, we have

$$\|\mathbf{c}_y^{(t_{\text{out}},t_{\text{in}})}\|_2 \leq \sqrt{n_{\mathcal{T}_y}} \max_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}},t_{\text{in}})}, \quad \|\mathbf{c}_y^{(t_{\text{out}},t_{\text{in}})}\|_2 \geq \sqrt{n_{\mathcal{T}_y}} \min_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}},t_{\text{in}})}.$$

Combining these relations with the coordinatewise bounds in equation 55 yields equation 56.  $\square$

Building on the coefficient expansions for the inner-loop weights (Lemma D.5) and the condensed signed sample (Lemma D.6), we next upper bound the condensed hinge margins uniformly over the inner loop. The resulting deterministic envelope yields a Corollary D.4-type sufficient condition ensuring an initial fully-active window for the condensed sample at each outer-loop restart. To account for the inheritance of  $\mathcal{S}$  across outer loops, we also introduce cumulative contraction factors that control the outer-loop initial correlations against  $\mathbf{w}^{(t_{\text{out}},0)}$  and  $\{\mathbf{z}_j^T\}$ .

**Lemma D.7.** *Assume Assumption D.1 and the additive model  $\mathbf{x}_i^T = y_i^T \boldsymbol{\mu} + \boldsymbol{\xi}_i$ . Set  $\rho := 1 - \lambda \eta_{\mathbf{w}}$  and  $\rho_{\mathcal{S}} := 1 - 2\eta_{\mathcal{S}}$ . Define the deterministic constants*

$$\begin{aligned} B_{\text{wS},0} &:= C_{\kappa} \sigma_{\mathcal{S}} \sigma_{\mathbf{w}} d, \\ G_{\text{TS},0} &:= C_{\kappa} \sigma_{\mathcal{S}} \left( \|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d} \right), \end{aligned} \quad (57)$$

and recall

$$\begin{aligned} B_{\text{init}} &:= C_{\kappa} \sigma_{\mathbf{w}} \left( \|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d} \right), \\ G_{\text{max}} &:= \|\boldsymbol{\mu}\|_2^2 + 2C_{\kappa} \sqrt{\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}} \|\boldsymbol{\mu}\|_2 \sqrt{d} + C_{\kappa} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d. \end{aligned} \quad (58)$$

Fix  $t_{\text{out}} \in \{0, \dots, T_{\text{out}} - 1\}$  and assume that for every  $k \in \{0, \dots, t_{\text{out}} - 1\}$ ,

$$q_{\mathcal{S},y}^{(k,t_{\text{in}})} = 1 \quad \text{for all } y \in \{\pm 1\} \text{ and all } t_{\text{in}} \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(k) - 1\},$$

for some integers  $T_{\mathcal{S},\text{lin}}(k) \in \{0, 1, \dots, T_{\text{in}}\}$ . Define the cumulative exponents

$$\begin{aligned} \Omega_{t_{\text{out}}} &:= \sum_{k=0}^{t_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k), \\ \Gamma_{t_{\text{out}}} &:= \sum_{r=0}^{t_{\text{out}}-1} \rho_{\mathcal{S}}^{\sum_{k=r+1}^{t_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k)}, \end{aligned} \quad (59)$$

with the conventions  $\Omega_0 = 0$  and  $\Gamma_0 = 0$ . Set

$$\begin{aligned} B_{\text{wS}}^{(t_{\text{out}})} &:= \rho_{\mathcal{S}}^{\Omega_{t_{\text{out}}}} B_{\text{wS},0} + (1 - \rho_{\mathcal{S}}^{T_{\text{in}}}) \Gamma_{t_{\text{out}}} B_{\text{init}}, \\ G_{\text{TS}}^{(t_{\text{out}})} &:= \rho_{\mathcal{S}}^{\Omega_{t_{\text{out}}}} G_{\text{TS},0} + (1 - \rho_{\mathcal{S}}^{T_{\text{in}}}) \Gamma_{t_{\text{out}}} G_{\text{max}}. \end{aligned} \quad (60)$$

Then the following hold.

(i) For every  $y \in \{\pm 1\}$  and every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ ,

$$\langle \mathbf{w}^{(t_{\text{out}},t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t_{\text{in}})} \rangle \leq \Phi_{t_{\text{out}}}(t_{\text{in}}), \quad (61)$$

where

$$\begin{aligned} \Phi_{t_{\text{out}}}(t_{\text{in}}) &:= \rho^{t_{\text{in}}} B_{\text{wS}}^{(t_{\text{out}})} + \frac{1 - \rho^{t_{\text{in}}}}{\lambda} G_{\text{TS}}^{(t_{\text{out}})} + \rho^{t_{\text{in}}} (1 - \rho_{\mathcal{S}}^{t_{\text{in}}}) B_{\text{init}} \\ &\quad + \frac{(1 - \rho^{t_{\text{in}}})(1 - \rho_{\mathcal{S}}^{t_{\text{in}}})}{\lambda} G_{\text{max}}. \end{aligned} \quad (62)$$

(ii) Define

$$T_{\mathcal{S},\text{lin}}(t_{\text{out}}) := \max \left\{ t \in \{0, 1, \dots, T_{\text{in}}\} \mid \Phi_{t_{\text{out}}}(s) < 1 \text{ for all } s \in \{0, 1, \dots, t\} \right\}, \quad (63)$$

with the convention  $\max \emptyset = -1$ . Then for every  $y \in \{\pm 1\}$  and every  $t \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}$ ,  $q_{\mathcal{S},y}^{(t_{\text{out}},t)} = 1$ . Moreover, if the training-set linear-window conclusion of Corollary D.4(i) holds on  $\{0, 1, \dots, T_{\text{lin}} - 1\}$ , then with  $T_{y,\text{lin}} := \min\{T_{\text{lin}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}$ ,

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t)} = \rho_{\mathcal{S}}^t \mathbf{z}_y^{\mathcal{S},(t_{\text{out}},0)} + (1 - \rho_{\mathcal{S}}^t) \bar{\mathbf{z}}_{\mathcal{T}_y}, \quad \bar{\mathbf{z}}_{\mathcal{T}_y} := \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \mathbf{z}_j^T, \quad t \in \{0, 1, \dots, T_{y,\text{lin}}\}. \quad (64)$$

(iii) Define the outer-loop inheritance coefficients  $(\tilde{u}_y^{(k)}, \tilde{c}_{y,j}^{(k)})$  for  $k \in \{0, 1, \dots, t_{\text{out}}\}$  and  $j \in \mathcal{I}(\mathcal{T}_y)$  by

$$\mathbf{z}_y^{\mathcal{S},(k,0)} = \tilde{u}_y^{(k)} \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(k)} \mathbf{z}_j^{\mathcal{T}}, \quad \tilde{u}_y^{(0)} := 1, \quad \tilde{c}_{y,j}^{(0)} := 0, \quad (65)$$

and for  $k \geq 0$ ,

$$\tilde{u}_y^{(k+1)} = u_y^{(k, T_{\text{in}})} \tilde{u}_y^{(k)}, \quad \tilde{c}_{y,j}^{(k+1)} = u_y^{(k, T_{\text{in}})} \tilde{c}_{y,j}^{(k)} + c_{y,j}^{(k, T_{\text{in}})}. \quad (66)$$

For each  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ , set

$$\tilde{u}_y^{(t_{\text{out}}, t_{\text{in}})} := u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{u}_y^{(t_{\text{out}})}, \quad \tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} := u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{c}_{y,j}^{(t_{\text{out}})} + c_{y,j}^{(t_{\text{out}}, t_{\text{in}})}. \quad (67)$$

Then for every  $y \in \{\pm 1\}$  and every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ ,

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}}, t_{\text{in}})} = \tilde{u}_y^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}. \quad (68)$$

Moreover, all coefficients are nonnegative, and they satisfy the deterministic bounds

$$\rho_{\mathcal{S}}^{t_{\text{out}} T_{\text{in}} + t_{\text{in}}} \leq \tilde{u}_y^{(t_{\text{out}}, t_{\text{in}})} \leq \rho_{\mathcal{S}}^{\Omega_{t_{\text{out}}} + \min\{t_{\text{in}}, T_{\mathcal{S}, \text{lin}}(t_{\text{out}})\}}. \quad (69)$$

If, in addition, the training-set linear-window conclusion of Corollary D.4(i) holds on  $\{0, 1, \dots, T_{\text{lin}} - 1\}$  for every outer-loop restart  $k \in \{0, 1, \dots, t_{\text{out}}\}$ , and we set  $T_{y, \text{lin}}(k) := \min\{T_{\text{lin}}, T_{\mathcal{S}, \text{lin}}(k)\}$ , then with

$$\underline{\Gamma}_{t_{\text{out}}, y} := \sum_{r=0}^{t_{\text{out}}-1} \rho_{\mathcal{S}}^{(t_{\text{out}}-1-r)T_{\text{in}}} \rho_{\mathcal{S}}^{(T_{\text{in}}-T_{y, \text{lin}}(r))_+} \left(1 - \rho_{\mathcal{S}}^{\min\{T_{\text{in}}, T_{y, \text{lin}}(r)\}}\right), \quad (70)$$

we have for every  $j \in \mathcal{I}(\mathcal{T}_y)$  and every  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ ,

$$\begin{aligned} \tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} &\geq \frac{\rho_{\mathcal{S}}^{t_{\text{in}}} \underline{\Gamma}_{t_{\text{out}}, y} + \rho_{\mathcal{S}}^{(t_{\text{in}}-T_{y, \text{lin}}(t_{\text{out}}))_+} \left(1 - \rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{y, \text{lin}}(t_{\text{out}})\}}\right)}{n_{\mathcal{T}_y}}, \\ \tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} &\leq \frac{\rho_{\mathcal{S}}^{\min\{t_{\text{in}}, T_{\mathcal{S}, \text{lin}}(t_{\text{out}})\}} \left(1 - \rho_{\mathcal{S}}^{T_{\text{in}}}\right) \Gamma_{t_{\text{out}}} + \left(1 - \rho_{\mathcal{S}}^{t_{\text{in}}}\right)}{n_{\mathcal{T}_y}}. \end{aligned} \quad (71)$$

*Proof.* For every  $k \geq 0$ , the inheritance  $\mathbf{z}_y^{\mathcal{S},(k+1,0)} = \mathbf{z}_y^{\mathcal{S},(k, T_{\text{in}})}$  and Lemma D.6(i) at time  $T_{\text{in}}$  give

$$\mathbf{z}_y^{\mathcal{S},(k+1,0)} = u_y^{(k, T_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(k,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(k, T_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}.$$

Substituting equation 65 for  $\mathbf{z}_y^{\mathcal{S},(k,0)}$  yields

$$\mathbf{z}_y^{\mathcal{S},(k+1,0)} = \left(u_y^{(k, T_{\text{in}})} \tilde{u}_y^{(k)}\right) \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \left(u_y^{(k, T_{\text{in}})} \tilde{c}_{y,j}^{(k)} + c_{y,j}^{(k, T_{\text{in}})}\right) \mathbf{z}_j^{\mathcal{T}},$$

which matches equation 66 after identifying coefficients in equation 65 at  $k+1$ . Lemma D.6(ii) gives  $u_y^{(k, T_{\text{in}})} \geq 0$  and  $c_{y,j}^{(k, T_{\text{in}})} \geq 0$ , hence induction on  $k$  implies  $\tilde{u}_y^{(k)} \geq 0$  and  $\tilde{c}_{y,j}^{(k)} \geq 0$  for all  $k$  and  $j$ .

Iterating equation 66 gives the unrolled forms

$$\tilde{u}_y^{(t_{\text{out}})} = \prod_{k=0}^{t_{\text{out}}-1} u_y^{(k, T_{\text{in}})}, \quad \tilde{c}_{y,j}^{(t_{\text{out}})} = \sum_{r=0}^{t_{\text{out}}-1} \left( \prod_{k=r+1}^{t_{\text{out}}-1} u_y^{(k, T_{\text{in}})} \right) c_{y,j}^{(r, T_{\text{in}})}.$$

Since  $u_y^{(k, T_{\text{in}})} = \prod_{s=0}^{T_{\text{in}}-1} (1 - 2\eta_s q_{\mathcal{S}, y}^{(k, s)})$  and  $1 - 2\eta_s q_{\mathcal{S}, y}^{(k, s)} \in [\rho_{\mathcal{S}}, 1]$ , we have  $u_y^{(k, T_{\text{in}})} \geq \rho_{\mathcal{S}}^{T_{\text{in}}}$  and therefore

$$\tilde{u}_y^{(t_{\text{out}})} \geq \rho_{\mathcal{S}}^{t_{\text{out}} T_{\text{in}}}.$$

The hypothesis on the previous outer loops and Lemma D.6(ii) give  $u_y^{(k, T_{\text{in}})} \leq \rho_S^{T_{\text{S}, \text{lin}}(k)}$ , hence

$$\tilde{u}_y^{(t_{\text{out}})} \leq \prod_{k=0}^{t_{\text{out}}-1} \rho_S^{T_{\text{S}, \text{lin}}(k)} = \rho_S^{\Omega_{t_{\text{out}}}}.$$

Lemma D.6(ii) also gives  $c_{y,j}^{(r, T_{\text{in}})} \leq (1 - \rho_S^{T_{\text{in}}})/n_{\mathcal{T}_y}$  and  $\prod_{k=r+1}^{t_{\text{out}}-1} u_y^{(k, T_{\text{in}})} \leq \rho_S^{\sum_{k=r+1}^{t_{\text{out}}-1} T_{\text{S}, \text{lin}}(k)}$ , hence

$$\tilde{c}_{y,j}^{(t_{\text{out}})} \leq \frac{1 - \rho_S^{T_{\text{in}}}}{n_{\mathcal{T}_y}} \sum_{r=0}^{t_{\text{out}}-1} \rho_S^{\sum_{k=r+1}^{t_{\text{out}}-1} T_{\text{S}, \text{lin}}(k)} = \frac{(1 - \rho_S^{T_{\text{in}}}) \Gamma_{t_{\text{out}}}}{n_{\mathcal{T}_y}}, \quad \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}})} \leq (1 - \rho_S^{T_{\text{in}}}) \Gamma_{t_{\text{out}}}.$$

For (i), applying equation 65 at  $k = t_{\text{out}}$  and taking inner products with  $\mathbf{w}^{(t_{\text{out}}, 0)}$  gives

$$\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle = \tilde{u}_y^{(t_{\text{out}})} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}})} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_j^{\mathcal{T}} \rangle.$$

Taking absolute values and using Lemma C.4 and Lemma D.3(i) yields

$$\begin{aligned} \left| \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle \right| &\leq \tilde{u}_y^{(t_{\text{out}})} B_{\text{WS}, 0} + \left( \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}})} \right) B_{\text{init}} \\ &\leq \rho_S^{\Omega_{t_{\text{out}}}} B_{\text{WS}, 0} + (1 - \rho_S^{T_{\text{in}}}) \Gamma_{t_{\text{out}}} B_{\text{init}} = B_{\text{WS}}^{(t_{\text{out}})}. \end{aligned}$$

Similarly, for any fixed  $k \in \mathcal{I}(\mathcal{T})$ ,

$$\langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle = \tilde{u}_y^{(t_{\text{out}})} \langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}})} \langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_j^{\mathcal{T}} \rangle,$$

and Lemma C.4 together with Lemma D.3(ii) yields

$$\begin{aligned} \left| \langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle \right| &\leq \tilde{u}_y^{(t_{\text{out}})} G_{\text{TS}, 0} + \left( \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}})} \right) G_{\text{max}} \\ &\leq \rho_S^{\Omega_{t_{\text{out}}}} G_{\text{TS}, 0} + (1 - \rho_S^{T_{\text{in}}}) \Gamma_{t_{\text{out}}} G_{\text{max}} = G_{\text{TS}}^{(t_{\text{out}})}. \end{aligned}$$

Fix  $t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}$ . Lemma D.5(i) gives

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} = \rho^{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}}, 0)} + \sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_k^{\mathcal{T}}, \quad a_k^{(t_{\text{out}}, t_{\text{in}})} \geq 0, \quad \sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \leq \frac{1 - \rho^{t_{\text{in}}}}{\lambda},$$

and Lemma D.6(i) gives

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} = u_y^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}}, \quad 0 \leq u_y^{(t_{\text{out}}, t_{\text{in}})} \leq 1, \quad c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \geq 0.$$

Expanding the inner product yields

$$\begin{aligned} \langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \rangle &= \rho^{t_{\text{in}}} u_y^{(t_{\text{out}}, t_{\text{in}})} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle + \rho^{t_{\text{in}}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_j^{\mathcal{T}} \rangle \\ &\quad + u_y^{(t_{\text{out}}, t_{\text{in}})} \sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} \rangle + \sum_{k \in \mathcal{I}(\mathcal{T})} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} a_k^{(t_{\text{out}}, t_{\text{in}})} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \langle \mathbf{z}_k^{\mathcal{T}}, \mathbf{z}_j^{\mathcal{T}} \rangle. \end{aligned}$$

Using  $u_y^{(t_{\text{out}}, t_{\text{in}})} \leq 1$  together with  $a_k^{(t_{\text{out}}, t_{\text{in}})} \geq 0$  and  $c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \geq 0$  yields

$$\begin{aligned} \langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \rangle &\leq \rho^{t_{\text{in}}} B_{\text{WS}}^{(t_{\text{out}})} + \rho^{t_{\text{in}}} B_{\text{init}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} + G_{\text{TS}}^{(t_{\text{out}})} \sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \\ &\quad + G_{\text{max}} \left( \sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \right) \left( \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \right). \end{aligned}$$

Lemma D.6(ii) gives  $\sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} = c_{y,\mu}^{(t_{\text{out}}, t_{\text{in}})} \leq 1 - \rho_S^{t_{\text{in}}}$ , and Lemma D.5(i) gives  $\sum_{k \in \mathcal{I}(\mathcal{T})} a_k^{(t_{\text{out}}, t_{\text{in}})} \leq (1 - \rho^{t_{\text{in}}})/\lambda$ , hence

$$\langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \rangle \leq \rho^{t_{\text{in}}} B_{\text{WS}}^{(t_{\text{out}})} + \rho^{t_{\text{in}}} (1 - \rho_S^{t_{\text{in}}}) B_{\text{init}} + \frac{1 - \rho^{t_{\text{in}}}}{\lambda} G_{\text{TS}}^{(t_{\text{out}})} + \frac{(1 - \rho^{t_{\text{in}}})(1 - \rho_S^{t_{\text{in}}})}{\lambda} G_{\text{max}},$$

which is equation 61–equation 62.

For (ii), if  $t \leq T_{\mathcal{S}, \text{lin}}(t_{\text{out}})$  then equation 63 gives  $\Phi_{t_{\text{out}}}(t) < 1$ , and equation 61 implies  $\langle \mathbf{w}^{(t_{\text{out}}, t)}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t)} \rangle < 1$ , hence  $q_{\mathcal{S}, y}^{(t_{\text{out}}, t)} = 1$ . For  $t \in \{0, 1, \dots, T_{y, \text{lin}}\}$  and each  $s \in \{0, 1, \dots, t-1\}$ , the joint activity on the corresponding linear windows implies  $\mathbf{g}_{\mathcal{T}_y}^{(t_{\text{out}}, s)} = \bar{\mathbf{z}}_{\mathcal{T}_y}$  and  $\mathbf{g}_{\mathcal{S}_y}^{(t_{\text{out}}, s)} = \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, s)}$ , hence the condensed update reduces to

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, s+1)} = \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, s)} + 2\eta_S (\bar{\mathbf{z}}_{\mathcal{T}_y} - \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, s)}) = \rho_S \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, s)} + (1 - \rho_S) \bar{\mathbf{z}}_{\mathcal{T}_y}.$$

Iterating gives  $\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t)} - \bar{\mathbf{z}}_{\mathcal{T}_y} = \rho_S^t (\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} - \bar{\mathbf{z}}_{\mathcal{T}_y})$ , which is equation 64.

For (iii), Lemma D.6(i) gives

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} = u_y^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, 0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}},$$

and substituting equation 65 at  $k = t_{\text{out}}$  yields

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} = (u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{u}_y^{(t_{\text{out}})}) \mathbf{z}_y^{\mathcal{S}, (0, 0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} (u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{c}_{y,j}^{(t_{\text{out}})} + c_{y,j}^{(t_{\text{out}}, t_{\text{in}})}) \mathbf{z}_j^{\mathcal{T}},$$

which is equation 68 with coefficients defined in equation 67. The bounds in equation 69 follow by multiplying the bounds for  $u_y^{(t_{\text{out}}, t_{\text{in}})}$  from Lemma D.6(ii) with the bounds for  $\tilde{u}_y^{(t_{\text{out}})}$  established above. For the upper bound in equation 71, equation 67 gives  $\tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} = u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{c}_{y,j}^{(t_{\text{out}})} + c_{y,j}^{(t_{\text{out}}, t_{\text{in}})}$ , and Lemma D.6(ii) together with the bound on  $\tilde{c}_{y,j}^{(t_{\text{out}})}$  derived above yields

$$\tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \leq \rho_S^{\min\{t_{\text{in}}, T_{\mathcal{S}, \text{lin}}(t_{\text{out}})\}} \frac{(1 - \rho_S^{T_{\text{in}}}) \Gamma_{t_{\text{out}}}}{n_{\mathcal{T}_y}} + \frac{1 - \rho_S^{t_{\text{in}}}}{n_{\mathcal{T}_y}}.$$

For the lower bound in equation 71, the unrolling of  $\tilde{c}_{y,j}^{(t_{\text{out}})}$  together with Lemma D.6(ii) gives, for each  $r \in \{0, 1, \dots, t_{\text{out}} - 1\}$ ,

$$\prod_{k=r+1}^{t_{\text{out}}-1} u_y^{(k, T_{\text{in}})} \geq \rho_S^{(t_{\text{out}}-1-r)T_{\text{in}}}, \quad c_{y,j}^{(r, T_{\text{in}})} \geq \frac{(\rho_S^{(T_{\text{in}}-T_{y, \text{lin}}(r))+} (1 - \rho_S^{\min\{T_{\text{in}}, T_{y, \text{lin}}(r)\}}))}{n_{\mathcal{T}_y}},$$

hence  $\tilde{c}_{y,j}^{(t_{\text{out}})} \geq \underline{\Gamma}_{t_{\text{out}}, y} / n_{\mathcal{T}_y}$ . Combining  $\tilde{c}_{y,j}^{(t_{\text{out}}, t_{\text{in}})} = u_y^{(t_{\text{out}}, t_{\text{in}})} \tilde{c}_{y,j}^{(t_{\text{out}})} + c_{y,j}^{(t_{\text{out}}, t_{\text{in}})}$  with Lemma D.6(ii), namely  $u_y^{(t_{\text{out}}, t_{\text{in}})} \geq \rho_S^{t_{\text{in}}}$  and

$$c_{y,j}^{(t_{\text{out}}, t_{\text{in}})} \geq \frac{(\rho_S^{(t_{\text{in}}-T_{y, \text{lin}}(t_{\text{out}}))+} (1 - \rho_S^{\min\{t_{\text{in}}, T_{y, \text{lin}}(t_{\text{out}})\}}))}{n_{\mathcal{T}_y}},$$

yields the stated lower bound in equation 71.  $\square$

#### D.4 Proof of Lemma 3.2

The following restatement records the terminal expansion and coefficient bounds used in the main text.

**Lemma D.8.** *Assume Assumption D.1 and the additive model  $\mathbf{x}_i^{\mathcal{T}} = y_i^{\mathcal{T}} \boldsymbol{\mu} + \boldsymbol{\xi}_i$ , and set  $\rho_S := 1 - 2\eta_S$ . Fix  $y \in \{\pm 1\}$ .*

(i) *There exist  $T_{\text{lin}} \in \{0, 1, \dots, T_{\text{in}}\}$  and  $T_{\mathcal{S}, \text{lin}}(t_{\text{out}}) \in \{0, 1, \dots, T_{\text{in}}\}$  for each  $t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}$  such that*

$$q_{\mathcal{T}, i}(\mathbf{w}^{(t_{\text{out}}, t)}) = 1, \quad i \in \mathcal{I}(\mathcal{T}), \quad t \in \{0, 1, \dots, T_{\text{lin}}\}, \quad t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\},$$

and

$$q_{\mathcal{S},y}^{(t_{\text{out}},t)} = 1, \quad t \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}, \quad t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}.$$

(ii) Let  $(\tilde{u}_y^{(t_{\text{out}},t_{\text{in}})}, \tilde{c}_{y,j}^{(t_{\text{out}},t_{\text{in}})})$  be the coefficients from Lemma D.7(iii), and define

$$\tilde{c}_{y,\mu}^{(t_{\text{out}},t_{\text{in}})} := \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}},t_{\text{in}})}, \quad \tilde{\mathbf{c}}_y^{(t_{\text{out}},t_{\text{in}})} := (\tilde{c}_{y,j}^{(t_{\text{out}},t_{\text{in}})})_{j \in \mathcal{I}(\mathcal{T}_y)}.$$

Then

$$\mathbf{z}_y^{\mathcal{S},(T_{\text{out}}-1,T_{\text{in}})} = \tilde{u}_y^{(T_{\text{out}}-1,T_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(T_{\text{out}}-1,T_{\text{in}})} \mathbf{z}_j^{\mathcal{I}},$$

and

$$\tilde{u}_y^{(T_{\text{out}}-1,T_{\text{in}})} \geq 0, \quad \tilde{c}_{y,j}^{(T_{\text{out}}-1,T_{\text{in}})} \geq 0, \quad j \in \mathcal{I}(\mathcal{T}_y).$$

(iii) Define

$$T_{y,\text{lin}}(r) := \min\{T_{\text{lin}}, T_{\mathcal{S},\text{lin}}(r)\}, \quad r \in \{0, 1, \dots, T_{\text{out}} - 1\}.$$

Define

$$B_{\text{out}} := 1 + \sum_{r=0}^{T_{\text{out}}-2} \rho_{\mathcal{S}}^{\sum_{k=r+1}^{T_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k)}, \quad u_{\text{out}} := (1 - \rho_{\mathcal{S}}^{T_{\text{in}}}) B_{\text{out}},$$

and

$$\ell_y := \sum_{r=0}^{T_{\text{out}}-1} \rho_{\mathcal{S}}^{(T_{\text{out}}-r)T_{\text{in}}-T_{y,\text{lin}}(r)} (1 - \rho_{\mathcal{S}}^{T_{y,\text{lin}}(r)}).$$

Then

$$\rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}} \leq \tilde{u}_y^{(T_{\text{out}}-1,T_{\text{in}})} \leq \rho_{\mathcal{S}}^{\sum_{k=0}^{T_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k)}.$$

Moreover, for every  $j \in \mathcal{I}(\mathcal{T}_y)$ ,

$$\frac{\ell_y}{n_{\mathcal{T}_y}} \leq \tilde{c}_{y,j}^{(T_{\text{out}}-1,T_{\text{in}})} \leq \frac{u_{\text{out}}}{n_{\mathcal{T}_y}},$$

and consequently

$$\ell_y \leq \tilde{c}_{y,\mu}^{(T_{\text{out}}-1,T_{\text{in}})} \leq u_{\text{out}}, \quad \frac{\ell_y}{\sqrt{n_{\mathcal{T}_y}}} \leq \|\tilde{\mathbf{c}}_y^{(T_{\text{out}}-1,T_{\text{in}})}\|_2 \leq \frac{u_{\text{out}}}{\sqrt{n_{\mathcal{T}_y}}}.$$

Finally, set  $\Delta_y := u_{\text{out}} - \ell_y$  and  $h_{y,j} := \tilde{c}_{y,j}^{(T_{\text{out}}-1,T_{\text{in}})} - \ell_y/n_{\mathcal{T}_y}$ . Then  $\Delta_y \geq 0$  and

$$h_{y,j} \geq 0, \quad \sum_{j \in \mathcal{I}(\mathcal{T}_y)} h_{y,j} \leq \Delta_y.$$

(iv) Moreover, assume  $G_{\text{max}} < \lambda$ . Then

$$q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t_{\text{in}})}) = 1, \quad i \in \mathcal{I}(\mathcal{T}), \quad t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}, \quad t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\},$$

and

$$q_{\mathcal{S},y}^{(t_{\text{out}},t_{\text{in}})} = 1, \quad t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}, \quad t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}.$$

Consequently,

$$\tilde{u}_y^{(T_{\text{out}}-1,T_{\text{in}})} = \rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}}, \quad \tilde{c}_{y,j}^{(T_{\text{out}}-1,T_{\text{in}})} = \frac{1 - \rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}}}{n_{\mathcal{T}_y}}, \quad j \in \mathcal{I}(\mathcal{T}_y).$$

*Proof.* For (i), Corollary D.4(i) provides  $T_{\text{lin}} \in \{0, 1, \dots, T_{\text{in}}\}$  such that  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}},t)}) = 1$  for all  $i \in \mathcal{I}(\mathcal{T})$  and all  $t \in \{0, 1, \dots, T_{\text{lin}}\}$ , uniformly over  $t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}$ . Lemma D.7(ii) provides  $T_{\mathcal{S},\text{lin}}(t_{\text{out}}) \in \{0, 1, \dots, T_{\text{in}}\}$  such that  $q_{\mathcal{S},y}^{(t_{\text{out}},t)} = 1$  for all  $t \in \{0, 1, \dots, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}$ .

For (ii), Lemma D.7(iii) yields, for every  $(t_{\text{out}}, t_{\text{in}})$ ,

$$\mathbf{z}_y^{\mathcal{S},(t_{\text{out}},t_{\text{in}})} = \tilde{u}_y^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(t_{\text{out}},t_{\text{in}})} \mathbf{z}_j^{\mathcal{T}},$$

with  $\tilde{u}_y^{(t_{\text{out}},t_{\text{in}})} \geq 0$  and  $\tilde{c}_{y,j}^{(t_{\text{out}},t_{\text{in}})} \geq 0$ . Setting  $(t_{\text{out}}, t_{\text{in}}) = (T_{\text{out}} - 1, T_{\text{in}})$  gives the asserted terminal expansion and nonnegativity.

For (iii), Lemma D.7(iii) gives, for every  $(t_{\text{out}}, t_{\text{in}})$ ,

$$\rho_{\mathcal{S}}^{t_{\text{out}}T_{\text{in}}+t_{\text{in}}} \leq \tilde{u}_y^{(t_{\text{out}},t_{\text{in}})} \leq \rho_{\mathcal{S}}^{\Omega_{t_{\text{out}}} + \min\{t_{\text{in}}, T_{\mathcal{S},\text{lin}}(t_{\text{out}})\}}.$$

With  $(t_{\text{out}}, t_{\text{in}}) = (T_{\text{out}} - 1, T_{\text{in}})$  and  $\min\{T_{\text{in}}, T_{\mathcal{S},\text{lin}}(T_{\text{out}} - 1)\} = T_{\mathcal{S},\text{lin}}(T_{\text{out}} - 1)$ , this yields

$$\rho_{\mathcal{S}}^{T_{\text{out}}T_{\text{in}}} \leq \tilde{u}_y^{(T_{\text{out}}-1, T_{\text{in}})} \leq \rho_{\mathcal{S}}^{\Omega_{T_{\text{out}}-1} + T_{\mathcal{S},\text{lin}}(T_{\text{out}}-1)}.$$

Since  $\Omega_{T_{\text{out}}-1} = \sum_{k=0}^{T_{\text{out}}-2} T_{\mathcal{S},\text{lin}}(k)$  by equation 59, we have  $\Omega_{T_{\text{out}}-1} + T_{\mathcal{S},\text{lin}}(T_{\text{out}} - 1) = \sum_{k=0}^{T_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k)$ , proving the displayed bounds for  $\tilde{u}_y^{(T_{\text{out}}-1, T_{\text{in}})}$ .

For the coordinatewise upper bound, apply the upper bound in equation 71 with  $(t_{\text{out}}, t_{\text{in}}) = (T_{\text{out}} - 1, T_{\text{in}})$  to obtain, for each  $j \in \mathcal{I}(\mathcal{T}_y)$ ,

$$\tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})} \leq \frac{\rho_{\mathcal{S}}^{T_{\mathcal{S},\text{lin}}(T_{\text{out}}-1)} (1 - \rho_{\mathcal{S}}^{T_{\text{in}}}) \Gamma_{T_{\text{out}}-1} + (1 - \rho_{\mathcal{S}}^{T_{\text{in}}})}{n_{\mathcal{T}_y}}.$$

By equation 59 with  $t_{\text{out}} = T_{\text{out}} - 1$ ,

$$\Gamma_{T_{\text{out}}-1} = \sum_{r=0}^{T_{\text{out}}-2} \rho_{\mathcal{S}}^{\sum_{k=r+1}^{T_{\text{out}}-2} T_{\mathcal{S},\text{lin}}(k)}, \quad \rho_{\mathcal{S}}^{T_{\mathcal{S},\text{lin}}(T_{\text{out}}-1)} \Gamma_{T_{\text{out}}-1} = \sum_{r=0}^{T_{\text{out}}-2} \rho_{\mathcal{S}}^{\sum_{k=r+1}^{T_{\text{out}}-1} T_{\mathcal{S},\text{lin}}(k)}.$$

Substituting this identity and factoring out  $(1 - \rho_{\mathcal{S}}^{T_{\text{in}}})$  yields

$$\tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})} \leq \frac{u_{\text{out}}}{n_{\mathcal{T}_y}}.$$

For the coordinatewise lower bound, apply the lower bound in equation 71 with  $(t_{\text{out}}, t_{\text{in}}) = (T_{\text{out}} - 1, T_{\text{in}})$  and  $T_{y,\text{lin}}(r) \leq T_{\text{in}}$  to obtain, for each  $j \in \mathcal{I}(\mathcal{T}_y)$ ,

$$\tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})} \geq \frac{\rho_{\mathcal{S}}^{T_{\text{in}}} \underline{\Gamma}_{T_{\text{out}}-1,y} + \rho_{\mathcal{S}}^{T_{\text{in}}-T_{y,\text{lin}}(T_{\text{out}}-1)} (1 - \rho_{\mathcal{S}}^{T_{y,\text{lin}}(T_{\text{out}}-1)})}{n_{\mathcal{T}_y}}.$$

By equation 70 with  $t_{\text{out}} = T_{\text{out}} - 1$  and  $T_{y,\text{lin}}(r) \leq T_{\text{in}}$ ,

$$\underline{\Gamma}_{T_{\text{out}}-1,y} = \sum_{r=0}^{T_{\text{out}}-2} \rho_{\mathcal{S}}^{(T_{\text{out}}-2-r)T_{\text{in}}} \rho_{\mathcal{S}}^{T_{\text{in}}-T_{y,\text{lin}}(r)} (1 - \rho_{\mathcal{S}}^{T_{y,\text{lin}}(r)}),$$

and hence

$$\rho_{\mathcal{S}}^{T_{\text{in}}} \underline{\Gamma}_{T_{\text{out}}-1,y} = \sum_{r=0}^{T_{\text{out}}-2} \rho_{\mathcal{S}}^{(T_{\text{out}}-r)T_{\text{in}}-T_{y,\text{lin}}(r)} (1 - \rho_{\mathcal{S}}^{T_{y,\text{lin}}(r)}).$$

The remaining term equals the  $r = T_{\text{out}} - 1$  summand, so combining the last two displays yields

$$\tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})} \geq \frac{\ell_y}{n_{\mathcal{T}_y}}.$$

Summing the coordinatewise bounds over  $j \in \mathcal{I}(\mathcal{T}_y)$  gives  $\ell_y \leq \tilde{c}_{y,\mu}^{(T_{\text{out}}-1, T_{\text{in}})} \leq u_{\text{out}}$ .

For the  $\ell_2$  bounds, fix a nonnegative vector  $\mathbf{v} = (v_j)_{j=1}^m$  and note that

$$\|\mathbf{v}\|_2^2 = \sum_{j=1}^m v_j^2 \leq m \left( \max_{1 \leq j \leq m} v_j \right)^2, \quad \|\mathbf{v}\|_2^2 = \sum_{j=1}^m v_j^2 \geq m \left( \min_{1 \leq j \leq m} v_j \right)^2.$$

Applying these inequalities with  $\mathbf{v} = \tilde{\mathbf{c}}_y^{(T_{\text{out}}-1, T_{\text{in}})}$  and  $m = n_{\mathcal{T}_y}$  gives

$$\sqrt{n_{\mathcal{T}_y}} \min_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})} \leq \|\tilde{\mathbf{c}}_y^{(T_{\text{out}}-1, T_{\text{in}})}\|_2 \leq \sqrt{n_{\mathcal{T}_y}} \max_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}^{(T_{\text{out}}-1, T_{\text{in}})}.$$

The coordinatewise bounds yield the displayed upper and lower bounds on  $\|\tilde{\mathbf{c}}_y^{(T_{\text{out}}-1, T_{\text{in}})}\|_2$ . Finally,  $\Delta_y \geq 0$  follows from  $\ell_y \leq u_{\text{out}}$ . The definition of  $h_{y,j}$  and the coordinatewise lower bound give  $h_{y,j} \geq 0$ , while the total-mass upper bound gives

$$\sum_{j \in \mathcal{I}(\mathcal{T}_y)} h_{y,j} = \tilde{c}_{y,\mu}^{(T_{\text{out}}-1, T_{\text{in}})} - \ell_y \leq u_{\text{out}} - \ell_y = \Delta_y.$$

For (iv), Corollary D.4(ii) implies  $q_{\mathcal{T},i}(\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}) = 1$  for all  $i \in \mathcal{I}(\mathcal{T})$  and all  $(t_{\text{out}}, t_{\text{in}})$ , so equation 40 gives

$$\mathbf{w}^{(t_{\text{out}}, t_{\text{in}})} = \alpha_{t_{\text{in}}} \mathbf{w}^{(t_{\text{out}}, 0)} + (1 - \alpha_{t_{\text{in}}}) \frac{1}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}, \quad \alpha_{t_{\text{in}}} := \rho^{t_{\text{in}}} \in [0, 1].$$

Under full training activity, the condensed update in class  $y$  satisfies

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}}+1)} = (1 - 2\eta_S q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})}) \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} + 2\eta_S q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \bar{\mathbf{z}}_{\mathcal{T}_y}, \quad \bar{\mathbf{z}}_{\mathcal{T}_y} := \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \mathbf{z}_j^{\mathcal{T}}.$$

Since  $q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \in \{0, 1\}$  and  $0 \leq 2\eta_S \leq 1$ , we have  $1 - 2\eta_S q_{\mathcal{S},y}^{(t_{\text{out}}, t_{\text{in}})} \in [0, 1]$ , and therefore  $\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}}+1)} \in \text{conv}\{\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})}, \bar{\mathbf{z}}_{\mathcal{T}_y}\}$ . Iterating in  $t_{\text{in}}$  and using the inheritance  $\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}+1, 0)} = \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, T_{\text{in}})}$  yields

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \in \text{conv}\{\mathbf{z}_y^{\mathcal{S}, (0, 0)}, \bar{\mathbf{z}}_{\mathcal{T}_y}\}, \quad t_{\text{out}} \in \{0, 1, \dots, T_{\text{out}} - 1\}, \quad t_{\text{in}} \in \{0, 1, \dots, T_{\text{in}}\}.$$

Fix  $(t_{\text{out}}, t_{\text{in}})$ , and write

$$\mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} = \beta_{t_{\text{out}}, t_{\text{in}}} \mathbf{z}_y^{\mathcal{S}, (0, 0)} + (1 - \beta_{t_{\text{out}}, t_{\text{in}}}) \bar{\mathbf{z}}_{\mathcal{T}_y}$$

for some  $\beta_{t_{\text{out}}, t_{\text{in}}} \in [0, 1]$ . Bilinearity gives

$$\begin{aligned} \langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{\mathcal{S}, (t_{\text{out}}, t_{\text{in}})} \rangle &= \alpha_{t_{\text{in}}} \beta_{t_{\text{out}}, t_{\text{in}}} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle + \alpha_{t_{\text{in}}} (1 - \beta_{t_{\text{out}}, t_{\text{in}}}) \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \bar{\mathbf{z}}_{\mathcal{T}_y} \rangle \\ &\quad + (1 - \alpha_{t_{\text{in}}}) \beta_{t_{\text{out}}, t_{\text{in}}} \left\langle \frac{1}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \right\rangle + (1 - \alpha_{t_{\text{in}}}) (1 - \beta_{t_{\text{out}}, t_{\text{in}}}) \left\langle \frac{1}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}, \bar{\mathbf{z}}_{\mathcal{T}_y} \right\rangle. \end{aligned}$$

Since the coefficients in the last display are nonnegative and sum to 1, this inner product is bounded above by the maximum of the four corner values. Under the high-probability event of Lemma C.4, we have

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle| \leq C_\kappa \sigma_S \sigma_{\mathbf{w}} d, \quad |\langle \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle| \leq C_\kappa \sigma_S \left( \|\boldsymbol{\mu}\|_2 \sqrt{d} + \sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d} \right).$$

Assumptions D.1A1 and D.1A4 imply that the right-hand sides are bounded by 1 and  $\lambda$ , respectively, and therefore

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \rangle| \leq 1, \quad \left| \left\langle \frac{1}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}, \mathbf{z}_y^{\mathcal{S}, (0, 0)} \right\rangle \right| \leq 1.$$

Moreover, Lemma D.3(i) gives, for every  $t_{\text{out}}$ ,

$$|\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \bar{\mathbf{z}}_{\mathcal{T}_y} \rangle| = \left| \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_j^{\mathcal{T}} \rangle \right| \leq \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} |\langle \mathbf{w}^{(t_{\text{out}}, 0)}, \mathbf{z}_j^{\mathcal{T}} \rangle| \leq B_{\text{init}} < 1.$$

Finally, Lemma D.3(ii) implies  $\langle \bar{\mathbf{z}}_{\mathcal{T}}, \bar{\mathbf{z}}_{\mathcal{T}_y} \rangle \leq G_{\max}$ , hence

$$\left\langle \frac{1}{\lambda} \bar{\mathbf{z}}_{\mathcal{T}}, \bar{\mathbf{z}}_{\mathcal{T}_y} \right\rangle \leq \frac{G_{\max}}{\lambda} < 1.$$

Combining these bounds yields  $\langle \mathbf{w}^{(t_{\text{out}}, t_{\text{in}})}, \mathbf{z}_y^{S, (t_{\text{out}}, t_{\text{in}})} \rangle < 1$  for every  $(t_{\text{out}}, t_{\text{in}})$ , and therefore  $q_{S,y}^{(t_{\text{out}}, t_{\text{in}})} = 1$  throughout. Under full activity, Lemma D.7(ii) gives

$$\mathbf{z}_y^{S, (t_{\text{out}}, T_{\text{in}})} = \rho_S^{T_{\text{in}}} \mathbf{z}_y^{S, (t_{\text{out}}, 0)} + (1 - \rho_S^{T_{\text{in}}}) \bar{\mathbf{z}}_{\mathcal{T}_y},$$

and iterating the inheritance  $\mathbf{z}_y^{S, (t_{\text{out}}+1, 0)} = \mathbf{z}_y^{S, (t_{\text{out}}, T_{\text{in}})}$  yields

$$\mathbf{z}_y^{S, (T_{\text{out}}-1, T_{\text{in}})} = \rho_S^{T_{\text{out}} T_{\text{in}}} \mathbf{z}_y^{S, (0, 0)} + (1 - \rho_S^{T_{\text{out}} T_{\text{in}}}) \bar{\mathbf{z}}_{\mathcal{T}_y}.$$

Expanding  $\bar{\mathbf{z}}_{\mathcal{T}_y} = \frac{1}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \mathbf{z}_j^T$  gives the coefficient identities in (iv).  $\square$

## D.5 Proof of Lemma 3.7

The proof of Lemma 3.7 is organized into one geometric risk bound and two signal-to-residual corollaries. The first result reduces the population risk of a directional evaluator to a one-dimensional sub-Gaussian tail bound and then lower bounds the effective margin using only a signal-residual decomposition of the two-sample difference vector.

**Lemma D.9.** *Assume the test distribution  $\mathcal{D}$  follows the additive model*

$$\mathbf{x} = y\boldsymbol{\mu} + \boldsymbol{\xi}_{\text{test}}, \quad y \in \{\pm 1\},$$

where  $\boldsymbol{\xi}_{\text{test}}$  is mean-zero sub-Gaussian with proxy covariance  $\boldsymbol{\Sigma}_{\xi}$  in the sense of Definition C.1. Let  $\mathcal{S} = \{(\mathbf{x}_+, +1), (\mathbf{x}_-, -1)\}$  be any two-point dataset and define the two-sample difference vector

$$\mathbf{s}(\mathcal{S}) := \mathbf{x}_+ - \mathbf{x}_-, \quad \widehat{\mathbf{w}}_{\mathcal{S}} := \frac{\mathbf{s}(\mathcal{S})}{\|\mathbf{s}(\mathcal{S})\|_2}.$$

Define the effective margin

$$m(\mathcal{S}) := \langle \widehat{\mathbf{w}}_{\mathcal{S}}, \boldsymbol{\mu} \rangle.$$

Then the population classification risk satisfies

$$\mathcal{R}(\mathcal{S}) := \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}}(y \langle \widehat{\mathbf{w}}_{\mathcal{S}}, \mathbf{x} \rangle < 0) \leq \exp\left(-\frac{m(\mathcal{S})_+^2}{2\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}}\right). \quad (72)$$

Moreover, assume that  $\mathbf{s}(\mathcal{S})$  admits the decomposition

$$\mathbf{s}(\mathcal{S}) = A(\mathcal{S})\boldsymbol{\mu} + \mathbf{r}(\mathcal{S}) \quad \text{with } A(\mathcal{S}) > 0. \quad (73)$$

Define the relative signal-to-residual ratio

$$\widehat{\text{SR}}(\mathcal{S}) := \frac{A(\mathcal{S})\|\boldsymbol{\mu}\|_2}{\|\mathbf{r}(\mathcal{S})\|_2} \in (0, +\infty], \quad (74)$$

where  $\widehat{\text{SR}}(\mathcal{S}) = +\infty$  if  $\mathbf{r}(\mathcal{S}) = \mathbf{0}$ . Then

$$m(\mathcal{S}) \geq \|\boldsymbol{\mu}\|_2 \left( \frac{\widehat{\text{SR}}(\mathcal{S}) - 1}{\widehat{\text{SR}}(\mathcal{S}) + 1} \right), \quad (75)$$

and therefore

$$\mathcal{R}(\mathcal{S}) \leq \exp\left(-\frac{\|\boldsymbol{\mu}\|_2^2}{2\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}} \left( \frac{\widehat{\text{SR}}(\mathcal{S}) - 1}{\widehat{\text{SR}}(\mathcal{S}) + 1} \right)_+^2\right). \quad (76)$$

*Proof.* Fix  $\mathcal{S}$  and abbreviate  $\widehat{\mathbf{w}} := \widehat{\mathbf{w}}_{\mathcal{S}}$ . For  $(\mathbf{x}, y) \sim \mathcal{D}$ , the error event is

$$y\langle \widehat{\mathbf{w}}, \mathbf{x} \rangle < 0 \iff \langle \widehat{\mathbf{w}}, \boldsymbol{\mu} \rangle + \langle \widehat{\mathbf{w}}, y\boldsymbol{\xi}_{\text{test}} \rangle < 0.$$

Since  $\boldsymbol{\xi}_{\text{test}}$  is sub-Gaussian with proxy covariance  $\boldsymbol{\Sigma}_{\xi}$ , the scalar  $\langle \widehat{\mathbf{w}}, y\boldsymbol{\xi}_{\text{test}} \rangle$  is mean-zero sub-Gaussian with proxy variance  $\widehat{\mathbf{w}}^{\top} \boldsymbol{\Sigma}_{\xi} \widehat{\mathbf{w}} \leq \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}$ . Hence, for any  $t > 0$ ,

$$\mathbb{P}(\langle \widehat{\mathbf{w}}, y\boldsymbol{\xi}_{\text{test}} \rangle \leq -t) \leq \exp\left(-\frac{t^2}{2\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}}\right).$$

Taking  $t = m(\mathcal{S})_+ = \langle \widehat{\mathbf{w}}, \boldsymbol{\mu} \rangle_+$  yields equation 72.

Assume now that equation 73 holds. Using  $\widehat{\mathbf{w}} = \mathbf{s}(\mathcal{S})/\|\mathbf{s}(\mathcal{S})\|_2$ ,

$$m(\mathcal{S}) = \frac{\langle \mathbf{s}(\mathcal{S}), \boldsymbol{\mu} \rangle}{\|\mathbf{s}(\mathcal{S})\|_2}.$$

For the numerator, Cauchy–Schwarz yields

$$\langle \mathbf{s}(\mathcal{S}), \boldsymbol{\mu} \rangle = A(\mathcal{S})\|\boldsymbol{\mu}\|_2^2 + \langle \mathbf{r}(\mathcal{S}), \boldsymbol{\mu} \rangle \geq A(\mathcal{S})\|\boldsymbol{\mu}\|_2^2 - \|\mathbf{r}(\mathcal{S})\|_2 \|\boldsymbol{\mu}\|_2.$$

For the denominator, the triangle inequality yields

$$\|\mathbf{s}(\mathcal{S})\|_2 = \|A(\mathcal{S})\boldsymbol{\mu} + \mathbf{r}(\mathcal{S})\|_2 \leq A(\mathcal{S})\|\boldsymbol{\mu}\|_2 + \|\mathbf{r}(\mathcal{S})\|_2.$$

Dividing the last two displays gives

$$m(\mathcal{S}) \geq \|\boldsymbol{\mu}\|_2 \frac{A(\mathcal{S})\|\boldsymbol{\mu}\|_2 - \|\mathbf{r}(\mathcal{S})\|_2}{A(\mathcal{S})\|\boldsymbol{\mu}\|_2 + \|\mathbf{r}(\mathcal{S})\|_2}.$$

Substituting the definition equation 74 yields equation 75. Finally, substituting equation 75 into equation 72 yields equation 76.  $\square$

We next instantiate Lemma D.9 for a coresnet formed by selecting one training sample per class.

**Corollary D.10.** *Assume Assumption D.1 and the additive model  $\mathbf{x}_i^{\top} = y_i^{\top} \boldsymbol{\mu} + \boldsymbol{\xi}_i$  with mean-zero sub-Gaussian noise parameter  $\boldsymbol{\Sigma}_{\xi}$ . Let  $C_{\kappa}$  denote the constant appearing in Lemma C.2. Pick any indices  $i_+ \in \mathcal{I}(\mathcal{T}_+)$  and  $i_- \in \mathcal{I}(\mathcal{T}_-)$  and form  $\mathcal{S}_{\text{core}} = \{(\mathbf{x}_{i_+}^{\top}, +1), (\mathbf{x}_{i_-}^{\top}, -1)\}$ . Then*

$$\mathbf{s}(\mathcal{S}_{\text{core}}) = \mathbf{x}_{i_+}^{\top} - \mathbf{x}_{i_-}^{\top} = 2\boldsymbol{\mu} + \mathbf{r}_{\text{core}}, \quad \mathbf{r}_{\text{core}} := \boldsymbol{\xi}_{i_+} - \boldsymbol{\xi}_{i_-},$$

and

$$\|\mathbf{r}_{\text{core}}\|_2^2 \leq 2(1 + 2C_{\kappa}) \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d. \tag{77}$$

Define

$$\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) := \frac{2\|\boldsymbol{\mu}\|_2}{\|\mathbf{r}_{\text{core}}\|_2}.$$

Then equation 77 implies the interpretable lower bound

$$\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) \geq \frac{2\|\boldsymbol{\mu}\|_2}{\sqrt{2(1 + 2C_{\kappa}) \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d}}. \tag{78}$$

Consequently,

$$\mathcal{R}(\mathcal{S}_{\text{core}}) \leq \exp\left(-\frac{\|\boldsymbol{\mu}\|_2^2}{2\|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}}} \left(\frac{\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) - 1}{\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) + 1}\right)_+^2\right),$$

and the same bound remains valid after replacing  $\widehat{\text{SR}}(\mathcal{S}_{\text{core}})$  by the lower bound in equation 78.

*Proof.* The identity  $\mathbf{s}(\mathcal{S}_{\text{core}}) = \mathbf{x}_{i_+}^T - \mathbf{x}_{i_-}^T = 2\boldsymbol{\mu} + (\boldsymbol{\xi}_{i_+} - \boldsymbol{\xi}_{i_-})$  follows from the additive model.

By Lemma C.2 equation 25 and  $\text{tr}(\boldsymbol{\Sigma}_\xi) \leq d \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}$ , for each  $i \in \mathcal{I}(\mathcal{T})$ ,

$$\|\boldsymbol{\xi}_i\|_2^2 \leq \text{tr}(\boldsymbol{\Sigma}_\xi) + C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d \leq (1 + C_\kappa) \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d.$$

By Lemma C.2 equation 24, for  $i \neq i'$ ,

$$|\langle \boldsymbol{\xi}_i, \boldsymbol{\xi}_{i'} \rangle| \leq C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d.$$

Therefore,

$$\begin{aligned} \|\mathbf{r}_{\text{core}}\|_2^2 &= \|\boldsymbol{\xi}_{i_+} - \boldsymbol{\xi}_{i_-}\|_2^2 \\ &= \|\boldsymbol{\xi}_{i_+}\|_2^2 + \|\boldsymbol{\xi}_{i_-}\|_2^2 - 2\langle \boldsymbol{\xi}_{i_+}, \boldsymbol{\xi}_{i_-} \rangle \\ &\leq (1 + C_\kappa) \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d + (1 + C_\kappa) \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d + 2C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d \\ &= 2(1 + 2C_\kappa) \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d, \end{aligned}$$

which is equation 77. The lower bound equation 78 follows from  $\widehat{\text{SR}}(\mathcal{S}_{\text{core}}) = 2\|\boldsymbol{\mu}\|_2/\|\mathbf{r}_{\text{core}}\|_2$ . The risk bound is Lemma D.9 applied with  $A(\mathcal{S}_{\text{core}}) = 2$  and  $\mathbf{r}(\mathcal{S}_{\text{core}}) = \mathbf{r}_{\text{core}}$ .  $\square$

We finally apply Lemma D.9 to the terminal condensed samples. The proof uses the terminal expansion and the coefficient bounds established in Lemma D.8. The strong-regularization regime  $G_{\text{max}} < \lambda$  is treated as a special case at the end.

**Corollary D.11.** *Assume the hypotheses of Lemma D.8. Let  $C_\kappa$  denote a constant that can be chosen to dominate the constants appearing in Lemma C.2, Lemma C.4, and Lemma C.5. For  $y \in \{\pm 1\}$ , set  $\tilde{u}_y := \tilde{u}_y^{(T_{\text{out}}-1, T_{\text{in}})}$ ,  $\tilde{\mathbf{c}}_y := \tilde{\mathbf{c}}_y^{(T_{\text{out}}-1, T_{\text{in}})}$ , and  $\tilde{c}_{y,\mu} := \tilde{c}_{y,\mu}^{(T_{\text{out}}-1, T_{\text{in}})}$ . Let  $\mathbf{z}_y^S := \mathbf{z}_y^{S, (T_{\text{out}}-1, T_{\text{in}})}$  be the terminal signed condensed samples and define  $\mathcal{S}_{\text{dc}} = \{(\mathbf{z}_+^S, +1), (-\mathbf{z}_-^S, -1)\}$ . Define the two-sample difference vector and its signal-residual decomposition by*

$$\mathbf{s}(\mathcal{S}_{\text{dc}}) := \mathbf{z}_+^S - (-\mathbf{z}_-^S) = \mathbf{z}_+^S + \mathbf{z}_-^S = A_{\text{dc}}\boldsymbol{\mu} + \mathbf{r}_{\text{dc}}, \quad A_{\text{dc}} := \tilde{c}_{+, \mu} + \tilde{c}_{-, \mu}.$$

Define the relative signal-to-residual ratio

$$\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) := \frac{A_{\text{dc}}\|\boldsymbol{\mu}\|_2}{\|\mathbf{r}_{\text{dc}}\|_2}.$$

For  $y \in \{\pm 1\}$  and  $r \in \{0, \dots, T_{\text{out}} - 1\}$ , set  $T_{y, \text{lin}}(r) := \min\{T_{\text{lin}}, T_{\mathcal{S}, \text{lin}}(r)\}$ , and define

$$B_{\text{out}} := 1 + \sum_{r=0}^{T_{\text{out}}-2} \rho_S^{\sum_{k=r+1}^{T_{\text{out}}-1} T_{\mathcal{S}, \text{lin}}(k)}.$$

Define  $u_{\text{out}} := (1 - \rho_S^{T_{\text{in}}})B_{\text{out}}$ . For  $y \in \{\pm 1\}$ , define

$$\ell_y := \sum_{r=0}^{T_{\text{out}}-1} \rho_S^{(T_{\text{out}}-r)T_{\text{in}} - T_{y, \text{lin}}(r)} (1 - \rho_S^{T_{y, \text{lin}}(r)}).$$

Set  $\Delta_y := u_{\text{out}} - \ell_y$ . Set

$$L_{\text{dc}} := \frac{\ell_+ + \ell_-}{B_{\text{out}}\beta_p + \sum_{y \in \{\pm 1\}} \Delta_y (\sqrt{n_T} - p_y^{-1/2})}.$$

Then the population risk satisfies

$$\mathcal{R}(\mathcal{S}_{\text{dc}}) \leq \exp\left(-\frac{\|\boldsymbol{\mu}\|_2^2}{2\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \left(\frac{\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) - 1}{\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) + 1}\right)_+^2\right).$$

Moreover,

$$\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) \geq \frac{\sqrt{n_{\mathcal{T}}} L_{\text{dc}} \|\boldsymbol{\mu}\|_2}{\sqrt{C_{\kappa} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d}},$$

and the same risk bound remains valid after replacing  $\widehat{\text{SR}}(\mathcal{S}_{\text{dc}})$  by the lower bound above.

Moreover, if  $G_{\text{max}} < \lambda$ , then we have

$$\tilde{u}_y = \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}, \quad \tilde{c}_{y,\mu} = 1 - \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}, \quad \|\tilde{\mathbf{c}}_y\|_2 = \frac{1 - \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}}{\sqrt{n_{\mathcal{T}_y}}}, \quad y \in \{\pm 1\},$$

and hence

$$L_{\text{dc}} = \frac{2(1 - \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}})}{\beta_p \sum_{q=0}^{T_{\text{out}}-1} \rho_{\mathcal{S}}^{q T_{\text{in}}}}, \quad \Delta_y = 0, \quad y \in \{\pm 1\},$$

so the stated lower bound becomes

$$\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) \geq \frac{2(1 - \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}) \|\boldsymbol{\mu}\|_2}{\sqrt{C_{\kappa} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d} \left( \sum_{q=0}^{T_{\text{out}}-1} \rho_{\mathcal{S}}^{q T_{\text{in}}} \right) \frac{\beta_p}{\sqrt{n_{\mathcal{T}}}}}$$

*Proof.* By Lemma D.8(ii), for each  $y \in \{\pm 1\}$ ,

$$\mathbf{z}_y^{\mathcal{S}} = \tilde{u}_y \mathbf{z}_y^{\mathcal{S},(0,0)} + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j} \mathbf{z}_j^{\mathcal{T}}.$$

Since  $\mathbf{z}_j^{\mathcal{T}} := y_j^{\mathcal{T}} \mathbf{x}_j^{\mathcal{T}}$  and  $\mathbf{x}_j^{\mathcal{T}} = y_j^{\mathcal{T}} \boldsymbol{\mu} + \boldsymbol{\xi}_j$ , we have  $\mathbf{z}_j^{\mathcal{T}} = \boldsymbol{\mu} + y_j^{\mathcal{T}} \boldsymbol{\xi}_j$ , and for  $j \in \mathcal{I}(\mathcal{T}_y)$ ,  $y_j^{\mathcal{T}} = y$ . Therefore,

$$\sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j} \mathbf{z}_j^{\mathcal{T}} = \tilde{c}_{y,\mu} \boldsymbol{\mu} + y \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j} \boldsymbol{\xi}_j, \quad \tilde{c}_{y,\mu} := \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j}.$$

Summing  $y = +1$  and  $y = -1$  gives

$$\mathbf{s}(\mathcal{S}_{\text{dc}}) = \mathbf{z}_+^{\mathcal{S}} + \mathbf{z}_-^{\mathcal{S}} = (\tilde{c}_{+,\mu} + \tilde{c}_{-,\mu}) \boldsymbol{\mu} + \mathbf{r}_{\text{init}} + \mathbf{r}_{\text{noise}},$$

where

$$\mathbf{r}_{\text{init}} := \tilde{u}_+ \mathbf{x}_+^{\mathcal{S},(0,0)} + \tilde{u}_- \mathbf{x}_-^{\mathcal{S},(0,0)}, \quad \mathbf{r}_{\text{noise}} := \sum_{j \in \mathcal{I}(\mathcal{T}_+)} \tilde{c}_{+,j} \boldsymbol{\xi}_j - \sum_{j \in \mathcal{I}(\mathcal{T}_-)} \tilde{c}_{-,j} \boldsymbol{\xi}_j.$$

Thus  $\mathbf{r}_{\text{dc}} = \mathbf{r}_{\text{init}} + \mathbf{r}_{\text{noise}}$  and  $A_{\text{dc}} = \tilde{c}_{+,\mu} + \tilde{c}_{-,\mu}$ . The risk bound is Lemma D.9 applied with  $A(\mathcal{S}_{\text{dc}}) = A_{\text{dc}}$  and  $\mathbf{r}(\mathcal{S}_{\text{dc}}) = \mathbf{r}_{\text{dc}}$ .

By Lemma C.4 equation 32,  $\|\mathbf{z}_y^{\mathcal{S},(0,0)}\|_2^2 \leq C_{\kappa} \sigma_{\mathcal{S}}^2 d$ , hence

$$\|\mathbf{r}_{\text{init}}\|_2 \leq \sqrt{C_{\kappa} d} \sigma_{\mathcal{S}} (\tilde{u}_+ + \tilde{u}_-).$$

For each  $y \in \{\pm 1\}$ , Lemma D.8(iii) gives the decomposition  $\tilde{c}_{y,j} = \ell_y / n_{\mathcal{T}_y} + h_{y,j}$  with  $h_{y,j} \geq 0$  and

$$\sum_{j \in \mathcal{I}(\mathcal{T}_y)} h_{y,j} \leq \Delta_y.$$

Therefore

$$\sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j} \boldsymbol{\xi}_j = \frac{\ell_y}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \boldsymbol{\xi}_j + \sum_{j \in \mathcal{I}(\mathcal{T}_y)} h_{y,j} \boldsymbol{\xi}_j.$$

For the first term, Lemma C.5 equation 35, applied with the failure probability split between the two fixed classwise-uniform vectors, yields simultaneously for  $y \in \{\pm 1\}$ ,

$$\left\| \frac{\ell_y}{n_{\mathcal{T}_y}} \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \boldsymbol{\xi}_j \right\|_2 \leq \sqrt{C_{\kappa} \|\boldsymbol{\Sigma}_{\xi}\|_{\text{op}} d} \frac{\ell_y}{\sqrt{n_{\mathcal{T}_y}}}.$$

For the second term, Lemma C.2 equation 25 and  $\text{tr}(\mathbf{\Sigma}_\xi) \leq \|\mathbf{\Sigma}_\xi\|_{\text{op}}d$  imply  $\max_j \|\xi_j\|_2 \leq \sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d}$ , and hence

$$\left\| \sum_{j \in \mathcal{I}(\mathcal{T}_y)} h_{y,j} \xi_j \right\|_2 \leq \sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d} \Delta_y.$$

Combining the last three displays gives

$$\left\| \sum_{j \in \mathcal{I}(\mathcal{T}_y)} \tilde{c}_{y,j} \xi_j \right\|_2 \leq \sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d} \left( \frac{\ell_y}{\sqrt{n_{\mathcal{T}_y}}} + \Delta_y \right).$$

Using the triangle inequality on  $\mathbf{r}_{\text{noise}}$  yields

$$\|\mathbf{r}_{\text{noise}}\|_2 \leq \sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d} \sum_{y \in \{\pm 1\}} \left( \frac{\ell_y}{\sqrt{n_{\mathcal{T}_y}}} + \Delta_y \right).$$

Lemma D.8(iii) also gives  $\tilde{u}_y \leq 1$  and  $\tilde{c}_{y,\mu} \geq \ell_y$ , hence

$$\tilde{u}_+ + \tilde{u}_- \leq 2, \quad A_{\text{dc}} \geq \ell_+ + \ell_-.$$

By Assumption D.1A4, we have

$$2\sigma_S \leq \sqrt{\|\mathbf{\Sigma}_\xi\|_{\text{op}}} \rho_S^{T_{\text{in}}} \frac{\beta_p}{\sqrt{n_{\mathcal{T}}}},$$

and since  $B_{\text{out}} \geq 1$ , this implies

$$\|\mathbf{r}_{\text{init}}\|_2 \leq 2\sqrt{C_\kappa d} \sigma_S \leq \sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d} \rho_S^{T_{\text{in}}} B_{\text{out}} \frac{\beta_p}{\sqrt{n_{\mathcal{T}}}}.$$

Combining the bounds on  $\mathbf{r}_{\text{init}}$  and  $\mathbf{r}_{\text{noise}}$ , and using  $n_{\mathcal{T}_y} = p_y n_{\mathcal{T}}$ , gives

$$\|\mathbf{r}_{\text{dc}}\|_2 \leq \frac{\sqrt{C_\kappa \|\mathbf{\Sigma}_\xi\|_{\text{op}}d}}{\sqrt{n_{\mathcal{T}}}} \left[ \rho_S^{T_{\text{in}}} B_{\text{out}} \beta_p + \sum_{y \in \{\pm 1\}} \ell_y p_y^{-1/2} + \sqrt{n_{\mathcal{T}}} \sum_{y \in \{\pm 1\}} \Delta_y \right].$$

Since  $\ell_y = u_{\text{out}} - \Delta_y$  and  $u_{\text{out}} = (1 - \rho_S^{T_{\text{in}}})B_{\text{out}}$ , the bracket equals

$$B_{\text{out}} \beta_p + \sum_{y \in \{\pm 1\}} \Delta_y (\sqrt{n_{\mathcal{T}}} - p_y^{-1/2}).$$

Substituting  $A_{\text{dc}} \geq \ell_+ + \ell_-$  and the preceding bound on  $\|\mathbf{r}_{\text{dc}}\|_2$  into the definition  $\widehat{\text{SR}}(\mathcal{S}_{\text{dc}}) = A_{\text{dc}} \|\boldsymbol{\mu}\|_2 / \|\mathbf{r}_{\text{dc}}\|_2$  gives the stated lower bound.

Finally, if  $G_{\text{max}} < \lambda$ , then Lemma D.8(iv) gives  $\tilde{c}_{y,\mu} = 1 - \rho_S^{T_{\text{out}}T_{\text{in}}}$  and  $\|\tilde{\mathbf{c}}_y\|_2 = (1 - \rho_S^{T_{\text{out}}T_{\text{in}}})/\sqrt{n_{\mathcal{T}_y}}$ . In this regime,  $T_{\mathcal{S},\text{lin}}(k) = T_{\text{in}}$  for all  $k$ , hence  $B_{\text{out}} = \sum_{q=0}^{T_{\text{out}}-1} \rho_S^{qT_{\text{in}}}$ ,  $\ell_y = u_{\text{out}} = 1 - \rho_S^{T_{\text{out}}T_{\text{in}}}$ , and  $\Delta_y = 0$  from the displayed identities.  $\square$

The lower bound in Corollary D.11 isolates the effect of activity switching through the deficits  $\Delta_y = u_{\text{out}} - \ell_y$ . The next deterministic comparison controls these deficits using the lengths of the certified active windows. Conditioned on  $T_{\text{in}}$  and  $T_{\mathcal{S},\text{lin}}(r)$ , the comparison follows algebraically from the definitions of  $B_{\text{out}}$ ,  $\ell_y$ , and  $u_{\text{out}}$ . It identifies regimes in which the factor  $L_{\text{dc}}$  preserves the class-averaging scale in Corollary D.11.

**Lemma D.12.** *Assume the hypotheses of Corollary D.11, and let  $B_{\text{out}}$ ,  $u_{\text{out}}$ ,  $\ell_y$ ,  $\Delta_y$ , and  $L_{\text{dc}}$  be the quantities defined there. Set  $T := T_{\text{in}}$  and  $\rho := \rho_S$ . For  $r \in \{0, \dots, T_{\text{out}} - 1\}$ , define*

$$D_r := \sum_{k=r+1}^{T_{\text{out}}-1} (T - T_{\mathcal{S},\text{lin}}(k)), \quad M_S := \max_{0 \leq r \leq T_{\text{out}}-1} D_r,$$

and

$$\mathcal{W}_{\text{win}} := (1 - \rho^T)(1 - \rho^{Ms}) + \max_{\substack{y \in \{\pm 1\} \\ 0 \leq r \leq T_{\text{out}} - 1}} (1 - \rho^{T - T_{y, \text{lin}}(r)}).$$

Then, for each  $y \in \{\pm 1\}$ ,

$$0 \leq \Delta_y = u_{\text{out}} - \ell_y \leq B_{\text{out}} \mathcal{W}_{\text{win}}, \quad \ell_y \geq B_{\text{out}}((1 - \rho^T) - \mathcal{W}_{\text{win}}).$$

Moreover,

$$\sum_{y \in \{\pm 1\}} \Delta_y \leq 2B_{\text{out}} \mathcal{W}_{\text{win}}.$$

The following consequences for  $L_{\text{dc}}$  hold along any sequence of problem instances for which  $\beta_p > 0$ .

(i) If  $G_{\text{max}} < \lambda$ , then  $\mathcal{W}_{\text{win}} = 0$  and

$$L_{\text{dc}} = \frac{2(1 - \rho^T)}{\beta_p}.$$

In particular, if  $p_+ = p_- = 1/2$ , then  $L_{\text{dc}} = (1 - \rho^T)/\sqrt{2}$ .

(ii) If

$$\mathcal{W}_{\text{win}} = o(1 - \rho^T), \quad \mathcal{W}_{\text{win}} = o\left(\frac{\beta_p}{\sqrt{n_{\mathcal{T}}}}\right),$$

then

$$L_{\text{dc}} = (1 + o(1)) \frac{2(1 - \rho^T)}{\beta_p}.$$

(iii) If

$$\mathcal{W}_{\text{win}} = o(1 - \rho^T), \quad \mathcal{W}_{\text{win}} = O\left(\frac{\beta_p}{\sqrt{n_{\mathcal{T}}}}\right),$$

then

$$L_{\text{dc}} = \Omega\left(\frac{1 - \rho^T}{\beta_p}\right).$$

(iv) Let  $n_{\text{min}} := \min_{y \in \{\pm 1\}} n_{\mathcal{T}_y}$ . If  $n_{\text{min}} \rightarrow \infty$  and

$$\sum_{y \in \{\pm 1\}} \Delta_y = \Omega(B_{\text{out}}(1 - \rho^T)),$$

then

$$L_{\text{dc}} = O(n_{\mathcal{T}}^{-1/2}).$$

*Proof.* For each  $r$ , set

$$F_r := \sum_{k=r+1}^{T_{\text{out}}-1} T_{\mathcal{S}, \text{lin}}(k), \quad s_{y,r} := T - T_{y, \text{lin}}(r).$$

Then  $F_r = (T_{\text{out}} - 1 - r)T - D_r$ , and the  $r$ th contribution to  $u_{\text{out}}$  is  $a_r := (1 - \rho^T)\rho^{F_r}$ , while the  $r$ th contribution to  $\ell_y$  is

$$b_{y,r} := \rho^{(T_{\text{out}}-r)T - T_{y, \text{lin}}(r)} (1 - \rho^{T_{y, \text{lin}}(r)}) = \rho^{F_r + D_r + s_{y,r}} (1 - \rho^{T - s_{y,r}}).$$

Hence

$$a_r - b_{y,r} = \rho^{F_r} [(1 - \rho^T) - \rho^{D_r + s_{y,r}} (1 - \rho^{T - s_{y,r}})].$$

The bracket admits the exact decomposition

$$(1 - \rho^T)(1 - \rho^{D_r}) + \rho^{D_r} (1 - \rho^{s_{y,r}}),$$

which is nonnegative and bounded above by  $\mathcal{W}_{\text{win}}$ . Summing over  $r$  gives

$$0 \leq u_{\text{out}} - \ell_y = \sum_{r=0}^{T_{\text{out}}-1} (a_r - b_{y,r}) \leq \mathcal{W}_{\text{win}} \sum_{r=0}^{T_{\text{out}}-1} \rho^{Fr} = B_{\text{out}} \mathcal{W}_{\text{win}}.$$

The lower bound on  $\ell_y$  follows from  $\ell_y = u_{\text{out}} - \Delta_y$ ,  $u_{\text{out}} = B_{\text{out}}(1 - \rho^T)$ , and the preceding display. Summing the bounds on  $\Delta_y$  over  $y \in \{\pm 1\}$  gives  $\sum_y \Delta_y \leq 2B_{\text{out}} \mathcal{W}_{\text{win}}$ .

If  $G_{\text{max}} < \lambda$ , Lemma D.8(iv) gives full activity, so  $T_{\mathcal{S},\text{lin}}(r) = T$  and  $T_{y,\text{lin}}(r) = T$  for all  $y, r$ . Thus  $M_{\mathcal{S}} = 0$ ,  $\mathcal{W}_{\text{win}} = 0$ ,  $B_{\text{out}} = \sum_{q=0}^{T_{\text{out}}-1} \rho^{qT}$ , and  $\ell_y = u_{\text{out}} = 1 - \rho^{T_{\text{out}}T}$ . Substituting these identities into the definition of  $L_{\text{dc}}$  gives

$$L_{\text{dc}} = \frac{2(1 - \rho^{T_{\text{out}}T})}{B_{\text{out}}\beta_p} = \frac{2(1 - \rho^T)}{\beta_p}.$$

For (ii) and (iii), write the numerator of  $L_{\text{dc}}$  as  $N := \ell_+ + \ell_-$  and its denominator as

$$D := B_{\text{out}}\beta_p + \sum_{y \in \{\pm 1\}} \Delta_y (\sqrt{n_{\mathcal{T}}} - p_y^{-1/2}).$$

The bounds above imply

$$2B_{\text{out}}((1 - \rho^T) - \mathcal{W}_{\text{win}}) \leq N \leq 2B_{\text{out}}(1 - \rho^T),$$

and

$$B_{\text{out}}\beta_p \leq D \leq B_{\text{out}}\beta_p + 2B_{\text{out}}\sqrt{n_{\mathcal{T}}} \mathcal{W}_{\text{win}}.$$

If  $\mathcal{W}_{\text{win}} = o(1 - \rho^T)$  and  $\mathcal{W}_{\text{win}} = o(\beta_p/\sqrt{n_{\mathcal{T}}})$ , then  $N = (1 + o(1))2B_{\text{out}}(1 - \rho^T)$  and  $D = (1 + o(1))B_{\text{out}}\beta_p$ , proving (ii). If  $\mathcal{W}_{\text{win}} = o(1 - \rho^T)$  and  $\mathcal{W}_{\text{win}} = O(\beta_p/\sqrt{n_{\mathcal{T}}})$ , then  $N = \Omega(B_{\text{out}}(1 - \rho^T))$  and  $D = O(B_{\text{out}}\beta_p)$ , proving (iii).

For (iv), the numerator in  $L_{\text{dc}}$  is at most  $2B_{\text{out}}(1 - \rho^T)$ , while its denominator is at least

$$(\sqrt{n_{\mathcal{T}}} - \max_y p_y^{-1/2}) \sum_{y \in \{\pm 1\}} \Delta_y.$$

Since  $\max_y p_y^{-1/2} = \sqrt{n_{\mathcal{T}}/n_{\text{min}}} = o(\sqrt{n_{\mathcal{T}}})$  when  $n_{\text{min}} \rightarrow \infty$ , the denominator is  $\Omega(\sqrt{n_{\mathcal{T}}} B_{\text{out}}(1 - \rho^T))$  under the stated lower bound on  $\sum_y \Delta_y$ . Hence  $L_{\text{dc}} = O(n_{\mathcal{T}}^{-1/2})$ .  $\square$

## D.6 Proofs for the Multiclass OvR Results

This subsection proves the auxiliary results used for Theorems B.3 and B.4, Corollary B.5, and Remark B.6 in Appendix B.3. The argument follows the binary affine-recursion proof after introducing head-wise OvR signed coordinates; the additional class and head indices enter through the activity and coefficient windows.

The first step is to certify active windows for the head-wise OvR hinge indicators. After passing to signed OvR samples, each fixed head has the same fully-active weight trajectory as in the binary analysis, so the binary margin-envelope argument can be applied head-wise.

**Lemma D.13.** *Assume Assumption B.1,  $0 < \eta_{\mathbf{w}}\lambda \leq 1$ , and  $0 < 2\eta_{\mathcal{S}} \leq 1$ . Set  $\rho := 1 - \eta_{\mathbf{w}}\lambda$  and  $\rho_{\mathcal{S}} := 1 - 2\eta_{\mathcal{S}}$ . For each class  $k$ , define*

$$B_{\mathcal{T},k}^{\text{OvR}} := \max_{\substack{0 \leq r < T_{\text{out}} \\ h \in [K] \\ i \in \mathcal{I}_k(\mathcal{T})}} \langle \mathbf{w}_h^{(r,0)}, \mathbf{z}_{i,h}^{\mathcal{T}} \rangle, \quad G_{\mathcal{T},k}^{\text{OvR}} := \max_{\substack{h \in [K] \\ i \in \mathcal{I}_k(\mathcal{T})}} \langle \bar{\mathbf{z}}_h^{\mathcal{T}}, \mathbf{z}_{i,h}^{\mathcal{T}} \rangle,$$

and let

$$T_{\mathcal{T},k,\text{lin}}^{\text{OvR}} := \max \left\{ m \in \{0, \dots, T_{\text{in}}\} : \rho^t B_{\mathcal{T},k}^{\text{OvR}} + \frac{1 - \rho^t}{\lambda} G_{\mathcal{T},k}^{\text{OvR}} < 1 \text{ for all } t = 0, \dots, m - 1 \right\}.$$

For each outer loop  $r$  and each class  $k$ , set  $\bar{\mathbf{z}}_{k,h}^{\mathcal{T}} := n_{\mathcal{T},k}^{-1} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \mathbf{z}_{i,h}^{\mathcal{T}}$  and define

$$B_{S,k}^{\text{OvR}}(r) := \max_{h \in [K]} \max \left\{ \left\langle \mathbf{w}_h^{(r,0)}, \mathbf{z}_{k,h}^{\mathcal{S},(r,0)} \right\rangle, \left\langle \mathbf{w}_h^{(r,0)}, \bar{\mathbf{z}}_{k,h}^{\mathcal{T}} \right\rangle \right\},$$

$$G_{S,k}^{\text{OvR}}(r) := \max_{h \in [K]} \max \left\{ \left\langle \bar{\mathbf{z}}_h^{\mathcal{T}}, \mathbf{z}_{k,h}^{\mathcal{S},(r,0)} \right\rangle, \left\langle \bar{\mathbf{z}}_h^{\mathcal{T}}, \bar{\mathbf{z}}_{k,h}^{\mathcal{T}} \right\rangle \right\}.$$

Then define

$$T_{S,k,\text{lin}}^{\text{OvR}}(r) := \max \left\{ m \in \{0, \dots, T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}\} : \rho^t B_{S,k}^{\text{OvR}}(r) + \frac{1-\rho^t}{\lambda} G_{S,k}^{\text{OvR}}(r) < 1 \text{ for all } t = 0, \dots, m-1 \right\}.$$

Under classwise OvR gradient matching, for every outer loop  $r$ ,

$$q_{i,h}^{\mathcal{T},(r,t)} = 1 \quad \forall i \in \mathcal{I}_k(\mathcal{T}), h \in [K], t < T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, \quad q_{k,h}^{\mathcal{S},(r,t)} = 1 \quad \forall h \in [K], t < T_{S,k,\text{lin}}^{\text{OvR}}(r).$$

In particular, if the two displayed margin inequalities in the definitions of  $T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$  and  $T_{S,k,\text{lin}}^{\text{OvR}}(r)$  hold for every  $t = 0, \dots, T_{\text{in}} - 1$ , then the corresponding window length is  $T_{\text{in}}$ .

*Proof.* Fix a head  $h$ . In the signed OvR coordinate  $\mathbf{z}_{i,h}^{\mathcal{T}} = \tau_{h,y_i^{\mathcal{T}}} \mathbf{x}_i^{\mathcal{T}}$ , the training activity is  $q_{i,h}^{\mathcal{T},(t)} = \mathbb{1}\{1 - \langle \mathbf{w}_h^{(t)}, \mathbf{z}_{i,h}^{\mathcal{T}} \rangle > 0\}$ , and the weight update can be written as

$$\mathbf{w}_h^{(t^+)} = (1 - \eta_{\mathbf{w}} \lambda) \mathbf{w}_h^{(t)} + \eta_{\mathbf{w}} \frac{1}{n_{\mathcal{T}}} \sum_{i=1}^{n_{\mathcal{T}}} q_{i,h}^{\mathcal{T},(t)} \mathbf{z}_{i,h}^{\mathcal{T}}.$$

Thus, for fixed  $h$ , the training-side dynamics have the same signed-sample form as the binary dynamics used in Corollary D.4. The proof of Corollary D.4, with the signed samples replaced by  $\{\mathbf{z}_{i,h}^{\mathcal{T}}\}_i$  and then with the maximum taken over  $h$  and  $i \in \mathcal{I}_k(\mathcal{T})$ , gives

$$\left\langle \mathbf{w}_h^{(r,t)}, \mathbf{z}_{i,h}^{\mathcal{T}} \right\rangle \leq \rho^t B_{\mathcal{T},k}^{\text{OvR}} + \frac{1-\rho^t}{\lambda} G_{\mathcal{T},k}^{\text{OvR}}.$$

The definition of  $T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$  therefore implies  $q_{i,h}^{\mathcal{T},(r,t)} = 1$  for all  $i \in \mathcal{I}_k(\mathcal{T})$ , all heads  $h$ , and all  $t < T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$ .

We next certify the condensed-sample window. Fix an outer loop  $r$ , a class  $k$ , and a head  $h$ . In a joint activity induction, when the training activities for class  $k$  and the condensed-sample activities for class  $k$  are active for all heads up to the current inner time, the averaged quantities satisfy  $\bar{q}_k^{\mathcal{S}} = 1$  and  $\bar{b}_{k,i} = 1$ . Multiplying the condensed-sample update by  $\tau_{h,k}$  gives

$$\mathbf{z}_{k,h}^{\mathcal{S},(t^+)} = \rho_{\mathcal{S}} \mathbf{z}_{k,h}^{\mathcal{S},(t)} + (1 - \rho_{\mathcal{S}}) \bar{\mathbf{z}}_{k,h}^{\mathcal{T}},$$

which is the full-active signed condensed-sample recursion used in Lemma D.7. Since  $T_{S,k,\text{lin}}^{\text{OvR}}(r) \leq T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$ , the training activity conclusion above gives the full-active head trajectory on this window:

$$\mathbf{w}_h^{(r,t)} = \rho^t \mathbf{w}_h^{(r,0)} + \frac{1-\rho^t}{\lambda} \bar{\mathbf{z}}_h^{\mathcal{T}}.$$

The signed condensed sample is a convex combination  $\mathbf{z}_{k,h}^{\mathcal{S},(r,t)} = \rho_{\mathcal{S}}^t \mathbf{z}_{k,h}^{\mathcal{S},(r,0)} + (1 - \rho_{\mathcal{S}}^t) \bar{\mathbf{z}}_{k,h}^{\mathcal{T}}$  on the same induction window. Therefore

$$\left\langle \mathbf{w}_h^{(r,t)}, \mathbf{z}_{k,h}^{\mathcal{S},(r,t)} \right\rangle \leq \rho^t B_{S,k}^{\text{OvR}}(r) + \frac{1-\rho^t}{\lambda} G_{S,k}^{\text{OvR}}(r).$$

The definition of  $T_{S,k,\text{lin}}^{\text{OvR}}(r)$  implies that the last display is strictly smaller than 1 for all  $t < T_{S,k,\text{lin}}^{\text{OvR}}(r)$ , and hence  $q_{k,h}^{\mathcal{S},(r,t)} = 1$  on that window. The statement about full-length windows follows from the definitions of the window lengths.  $\square$

Once the active windows are fixed, the multiclass condensed-sample update can be unrolled class by class. The next lemma is the OvR counterpart of the binary coefficient expansion: it gives a nonnegative expansion, the corresponding signal–noise decomposition, and the coefficient bounds needed later for the pairwise risk analysis.

**Lemma D.14.** *Assume  $0 < 2\eta_S \leq 1$  and  $n_{\mathcal{T},k} \geq 1$  for every  $k \in [K]$ , and set  $\rho_S := 1 - 2\eta_S$  and  $\mathbf{t}^* := (T_{\text{out}} - 1, T_{\text{in}})$ . Under classwise OvR gradient matching, for every  $k \in [K]$  there exist terminal coefficients*

$$u_k^{(\mathbf{t}^*)} \geq 0, \quad c_{k,i}^{(\mathbf{t}^*)} \geq 0, \quad i \in \mathcal{I}_k(\mathcal{T}),$$

such that

$$\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)} = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{\mathcal{S},(0,0)} + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \mathbf{x}_i^{\mathcal{T}}.$$

Set  $c_{k,\mu}^{(\mathbf{t}^*)} := \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)}$  and  $\mathbf{c}_k^{(\mathbf{t}^*)} := (c_{k,i}^{(\mathbf{t}^*)})_{i \in \mathcal{I}_k(\mathcal{T})}$ . If the  $K$ -class additive model satisfies  $\mathbf{x}_i^{\mathcal{T}} = \boldsymbol{\mu}_k + \boldsymbol{\xi}_i$  for  $i \in \mathcal{I}_k(\mathcal{T})$ , then

$$\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)} = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{\mathcal{S},(0,0)} + c_{k,\mu}^{(\mathbf{t}^*)} \boldsymbol{\mu}_k + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \boldsymbol{\xi}_i.$$

Assume additionally that the OvR training and condensed-sample linear windows for class  $k$  have lengths  $T_{\mathcal{T},k,\text{lin}}^{\text{OvR}} \in \{0, \dots, T_{\text{in}}\}$  and  $T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r) \in \{0, \dots, T_{\text{in}}\}$ ,  $r \in \{0, \dots, T_{\text{out}} - 1\}$ , in the sense that, for every outer loop  $r$ ,

$$q_{i,h}^{\mathcal{T},(r,t)} = 1 \quad \forall i \in \mathcal{I}_k(\mathcal{T}), h \in [K], t < T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, \quad q_{k,h}^{\mathcal{S},(r,t)} = 1 \quad \forall h \in [K], t < T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r).$$

Define

$$T_{k,\text{lin}}^{\text{OvR}}(r) := \min\{T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}, T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r)\},$$

$$B_{k,\text{out}} := 1 + \sum_{r=0}^{T_{\text{out}}-2} \rho_S^{\sum_{s=r+1}^{T_{\text{out}}-1} T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(s)}, \quad u_{k,\text{out}} := (1 - \rho_S^{T_{\text{in}}}) B_{k,\text{out}},$$

and

$$\ell_k := \sum_{r=0}^{T_{\text{out}}-1} \rho_S^{(T_{\text{out}}-r)T_{\text{in}}-T_{k,\text{lin}}^{\text{OvR}}(r)} (1 - \rho_S^{T_{k,\text{lin}}^{\text{OvR}}(r)}).$$

Then

$$\rho_S^{T_{\text{out}}T_{\text{in}}} \leq u_k^{(\mathbf{t}^*)} \leq \rho_S^{\sum_{r=0}^{T_{\text{out}}-1} T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(r)}.$$

Moreover, for every  $i \in \mathcal{I}_k(\mathcal{T})$ ,

$$\frac{\ell_k}{n_{\mathcal{T},k}} \leq c_{k,i}^{(\mathbf{t}^*)} \leq \frac{u_{k,\text{out}}}{n_{\mathcal{T},k}},$$

and consequently

$$\ell_k \leq c_{k,\mu}^{(\mathbf{t}^*)} \leq u_{k,\text{out}}, \quad \frac{\ell_k}{\sqrt{n_{\mathcal{T},k}}} \leq \|\mathbf{c}_k^{(\mathbf{t}^*)}\|_2 \leq \frac{u_{k,\text{out}}}{\sqrt{n_{\mathcal{T},k}}}.$$

With  $\Delta_k := u_{k,\text{out}} - \ell_k$  and  $h_{k,i} := c_{k,i}^{(\mathbf{t}^*)} - \ell_k/n_{\mathcal{T},k}$ , we have  $\Delta_k \geq 0$ ,  $h_{k,i} \geq 0$ , and  $\sum_{i \in \mathcal{I}_k(\mathcal{T})} h_{k,i} \leq \Delta_k$ .

If all OvR training and condensed-sample activities are active for every update time in the bilevel procedure, then

$$\mathbf{p}_k^{\mathcal{S},(\mathbf{t}^*)} = \rho_S^{T_{\text{out}}T_{\text{in}}} \mathbf{p}_k^{\mathcal{S},(0,0)} + (1 - \rho_S^{T_{\text{out}}T_{\text{in}}}) \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}}, \quad \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}} := \frac{1}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \mathbf{x}_i^{\mathcal{T}}.$$

*Proof.* We derive the affine recursion for the classwise update. Fix  $k \in [K]$  and an inner-loop time index  $\mathbf{t} = (t_{\text{out}}, t_{\text{in}})$  with  $t_{\text{in}} < T_{\text{in}}$ , and let  $\mathbf{t}^+ = (t_{\text{out}}, t_{\text{in}} + 1)$ . For a fixed head  $h$ , write

$$\mathbf{r}_{k,h}^{(\mathbf{t})} := q_{k,h}^{\mathcal{S},(\mathbf{t})} \tau_{h,k} \mathbf{p}_k^{\mathcal{S},(\mathbf{t})} - \frac{1}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} q_{i,h}^{\mathcal{T},(\mathbf{t})} \tau_{h,k} \mathbf{x}_i^{\mathcal{T}}.$$

Then  $D_k^{\text{OvR}}(\mathcal{P}^{(\mathbf{t})}; \mathbf{W}^{(\mathbf{t})}) = K^{-1} \sum_{h=1}^K \|\mathbf{r}_{k,h}^{(\mathbf{t})}\|_2^2$ . Since all hinge activities are evaluated at time  $\mathbf{t}$  and treated as fixed with respect to  $\mathbf{p}_k^S$ , we have

$$\nabla_{\mathbf{p}_k^S} \left\| \mathbf{r}_{k,h}^{(\mathbf{t})} \right\|_2^2 = 2q_{k,h}^{S,(\mathbf{t})} \tau_{h,k} \mathbf{r}_{k,h}^{(\mathbf{t})}.$$

Substituting the definition of  $\mathbf{r}_{k,h}^{(\mathbf{t})}$  and using  $\tau_{h,k}^2 = 1$  and  $(q_{k,h}^{S,(\mathbf{t})})^2 = q_{k,h}^{S,(\mathbf{t})}$  gives

$$\nabla_{\mathbf{p}_k^S} \left\| \mathbf{r}_{k,h}^{(\mathbf{t})} \right\|_2^2 = 2q_{k,h}^{S,(\mathbf{t})} \mathbf{p}_k^{S,(\mathbf{t})} - \frac{2}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} q_{k,h}^{S,(\mathbf{t})} q_{i,h}^{\mathcal{T},(\mathbf{t})} \mathbf{x}_i^{\mathcal{T}}.$$

Averaging over  $h$  yields

$$\nabla_{\mathbf{p}_k^S} D_k^{\text{OvR}}(\mathcal{P}^{(\mathbf{t})}; \mathbf{W}^{(\mathbf{t})}) = 2\bar{q}_k^{S,(\mathbf{t})} \mathbf{p}_k^{S,(\mathbf{t})} - \frac{2}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \bar{b}_{k,i}^{(\mathbf{t})} \mathbf{x}_i^{\mathcal{T}}.$$

Inserting this gradient into the classwise condensed-sample step gives

$$\mathbf{p}_k^{S,(\mathbf{t}^+)} = \left(1 - 2\eta_S \bar{q}_k^{S,(\mathbf{t})}\right) \mathbf{p}_k^{S,(\mathbf{t})} + \frac{2\eta_S}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \bar{b}_{k,i}^{(\mathbf{t})} \mathbf{x}_i^{\mathcal{T}}.$$

We next unroll this update. Define the global coefficients by

$$u_k^{(0,0)} := 1, \quad c_{k,i}^{(0,0)} := 0, \quad i \in \mathcal{I}_k(\mathcal{T}).$$

For an inner-loop update, set

$$u_k^{(\mathbf{t}^+)} := \left(1 - 2\eta_S \bar{q}_k^{S,(\mathbf{t})}\right) u_k^{(\mathbf{t})}$$

and, for each  $i \in \mathcal{I}_k(\mathcal{T})$ ,

$$c_{k,i}^{(\mathbf{t}^+)} := \left(1 - 2\eta_S \bar{q}_k^{S,(\mathbf{t})}\right) c_{k,i}^{(\mathbf{t})} + \frac{2\eta_S}{n_{\mathcal{T},k}} \bar{b}_{k,i}^{(\mathbf{t})}.$$

Across the outer-loop boundary, set

$$u_k^{(t_{\text{out}}+1,0)} := u_k^{(t_{\text{out}}, T_{\text{in}})}, \quad c_{k,i}^{(t_{\text{out}}+1,0)} := c_{k,i}^{(t_{\text{out}}, T_{\text{in}})}.$$

These recursions match those in Lemma D.6(i) after replacing the binary activity  $q_{S,y}$  by  $\bar{q}_k^S$ , the product  $q_{S,y} q_{\mathcal{T},j}$  by  $\bar{b}_{k,i}$ , and the signed training sample  $\mathbf{z}_j^{\mathcal{T}}$  by  $\mathbf{x}_i^{\mathcal{T}}$ . The affine-recursion argument in Lemma D.6(i) depends only on the update form, the inequalities  $0 \leq \bar{b}_{k,i}^{(\mathbf{t})} \leq \bar{q}_k^{S,(\mathbf{t})} \leq 1$ , and  $0 < 2\eta_S \leq 1$ . It therefore yields the expansion and coefficient nonnegativity for every condensed state, in particular at  $\mathbf{t}^*$ .

Under the additive model, substituting  $\mathbf{x}_i^{\mathcal{T}} = \boldsymbol{\mu}_k + \boldsymbol{\xi}_i$  into the coefficient expansion gives

$$\mathbf{p}_k^{S,(\mathbf{t})} = u_k^{(\mathbf{t})} \mathbf{p}_k^{S,(0,0)} + \left( \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t})} \right) \boldsymbol{\mu}_k + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t})} \boldsymbol{\xi}_i,$$

and setting  $\mathbf{t} = \mathbf{t}^*$  gives the stated signal–noise expansion.

The terminal coefficient estimates follow from the coefficient-bound calculation in Lemma D.6(ii) and the outer-loop coefficient recursion in Lemma D.7(iii). The substitutions are  $n_{\mathcal{T}_y} \mapsto n_{\mathcal{T},k}$ ,  $T_{S,\text{lin}}(r) \mapsto T_{S,k,\text{lin}}^{\text{OvR}}(r)$ ,  $T_{\text{lin}} \mapsto T_{\mathcal{T},k,\text{lin}}^{\text{OvR}}$ , and  $T_{y,\text{lin}}(r) \mapsto T_{k,\text{lin}}^{\text{OvR}}(r)$ . The coefficient-bound argument remains valid for the averaged activities, since the only inequalities used in the upper and lower bounds are  $0 \leq \bar{b}_{k,i} \leq \bar{q}_k^S \leq 1$ ,  $\rho_S \leq 1 - 2\eta_S \bar{q}_k^S \leq 1$ , and  $\bar{q}_k^S = \bar{b}_{k,i} = 1$  on the joint linear window. Applying this recursion gives the displayed bounds on  $u_k^{(\mathbf{t}^*)}$ ,  $c_{k,i}^{(\mathbf{t}^*)}$ ,  $c_{k,\mu}^{(\mathbf{t}^*)}$ , and  $\|\mathbf{c}_k^{(\mathbf{t}^*)}\|_2$ . The inequalities for  $\Delta_k$  and  $h_{k,i}$  follow, as in Lemma D.8, from the coordinatewise lower bound and the total-mass upper bound.

Finally, in the full-activity case,  $\bar{d}_k^{S,(r,t)} = 1$  and  $\bar{b}_{k,i}^{(r,t)} = 1$  throughout the bilevel procedure. The affine recursion reduces at every update to

$$\mathbf{p}_k^{S,(t^+)} = \rho_S \mathbf{p}_k^{S,(t)} + (1 - \rho_S) \widehat{\boldsymbol{\mu}}_k^{\mathcal{T}}.$$

Iterating this recursion through the  $T_{\text{out}} T_{\text{in}}$  condensed-sample updates, with the usual outer-loop inheritance, yields the displayed contraction toward  $\widehat{\boldsymbol{\mu}}_k^{\mathcal{T}}$ .  $\square$

The preceding lemma concerns only the dynamics. The following lemma is a geometric reduction for the nearest-prototype evaluator: if every condensed sample is close to its class mean, then the multiclass error is controlled by pairwise class separation.

**Lemma D.15.** *Assume the  $K$ -class additive test model  $\mathbf{x} = \boldsymbol{\mu}_y + \boldsymbol{\xi}$ , where  $\boldsymbol{\xi}$  is mean-zero sub-Gaussian with proxy covariance  $\boldsymbol{\Sigma}_\xi$ . Let  $\mathcal{P} = \{(\mathbf{p}_k, k) : k \in [K]\}$  be any multiclass condensed set used with the nearest-prototype evaluator, let  $\pi_k := \mathbb{P}(y = k)$ , and define  $\epsilon_k := \|\mathbf{p}_k - \boldsymbol{\mu}_k\|_2$ ,  $\gamma_{k\ell} := \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_\ell\|_2$ ,  $R_{k\ell} := \epsilon_k + \epsilon_\ell$ , and  $\widehat{\text{SR}}_{k\ell}^{\text{mc}} := \gamma_{k\ell}/R_{k\ell}$ , with value  $+\infty$  when  $R_{k\ell} = 0$ . For  $s \in (0, +\infty)$ , set  $\Psi_{\text{mc}}(s) := ((1 - s^{-1})_+)^4 / (1 + s^{-1})^2$ , with  $\Psi_{\text{mc}}(0) := 0$  and  $\Psi_{\text{mc}}(+\infty) := 1$ . Then*

$$\mathcal{R}_{\text{mc}}(\mathcal{P}) \leq \sum_{k=1}^K \pi_k \sum_{\ell \neq k} \exp\left(-\frac{\gamma_{k\ell}^2}{8\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \Psi_{\text{mc}}\left(\widehat{\text{SR}}_{k\ell}^{\text{mc}}\right)\right).$$

*Proof.* Fix  $k \neq \ell$  and condition on  $y = k$ , so that  $\mathbf{x} = \boldsymbol{\mu}_k + \boldsymbol{\xi}$ . The event  $E_{k \rightarrow \ell}$  is equivalent to

$$\langle \boldsymbol{\xi}, \mathbf{p}_\ell - \mathbf{p}_k \rangle \geq \frac{1}{2} (\|\mathbf{p}_\ell - \boldsymbol{\mu}_k\|_2^2 - \|\mathbf{p}_k - \boldsymbol{\mu}_k\|_2^2).$$

Let

$$A_{k\ell} := (\gamma_{k\ell}^2 - 2\gamma_{k\ell}\epsilon_\ell + \epsilon_\ell^2 - \epsilon_k^2)_+.$$

Since  $\|\mathbf{p}_\ell - \boldsymbol{\mu}_k\|_2 \geq \gamma_{k\ell} - \epsilon_\ell$  and  $\|\mathbf{p}_k - \boldsymbol{\mu}_k\|_2 = \epsilon_k$ , the positive part of the threshold is at least  $A_{k\ell}/2$ . Also  $\|\mathbf{p}_\ell - \mathbf{p}_k\|_2 \leq \gamma_{k\ell} + R_{k\ell}$ . The one-dimensional sub-Gaussian tail bound gives

$$\mathbb{P}(E_{k \rightarrow \ell}) \leq \exp\left(-\frac{A_{k\ell}^2}{8\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}(\gamma_{k\ell} + R_{k\ell})^2}\right).$$

Moreover,  $A_{k\ell} \geq (\gamma_{k\ell} - R_{k\ell})_+^2$ , and hence

$$\mathbb{P}(E_{k \rightarrow \ell}) \leq \exp\left(-\frac{((\gamma_{k\ell} - R_{k\ell})_+)^4}{8\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}(\gamma_{k\ell} + R_{k\ell})^2}\right) = \exp\left(-\frac{\gamma_{k\ell}^2}{8\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \Psi_{\text{mc}}\left(\widehat{\text{SR}}_{k\ell}^{\text{mc}}\right)\right),$$

with the stated conventions when  $R_{k\ell} = 0$ . The multiclass risk is bounded by applying the union bound over  $\ell \neq k$  and then averaging over the class prior  $\pi_k$ .  $\square$

Combining the coefficient expansion with the high-probability noise controls converts the terminal OvR condensed set into pairwise certified multiclass signal-to-residual ratios, which can then be substituted into the preceding nearest-prototype risk bound.

**Corollary D.16.** *Assume the hypotheses of Lemma D.14 and the high-probability event in Assumption B.1. Let  $C_\kappa$  denote the corresponding uniform concentration constant, let  $\mathbf{t}^* = (T_{\text{out}} - 1, T_{\text{in}})$ , and define the terminal multiclass condensed set  $\mathcal{P}_{\text{dc}} := \{(\mathbf{p}_k^{S,(\mathbf{t}^*)}, k) : k \in [K]\}$ . Let  $\ell_k$ ,  $\Delta_k$ , and  $T_{S,k,\text{lin}}^{\text{OvR}}(r)$  be the coefficient-window quantities from Lemma D.14. Set  $p_k := n_{\mathcal{T},k}/n_{\mathcal{T}}$ ,*

$$\alpha_S := \frac{\sigma_S \sqrt{n_{\mathcal{T}}}}{\sqrt{\|\boldsymbol{\Sigma}_\xi\|_{\text{op}}}}, \quad \alpha_{\mu,k} := \frac{\|\boldsymbol{\mu}_k\|_2 \sqrt{n_{\mathcal{T}}}}{\sqrt{C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}},$$

and, for each class  $k$ , let

$$\bar{u}_k := \rho_S \sum_{r=0}^{T_{\text{out}}-1} T_{S,k,\text{lin}}^{\text{OvR}}(r)$$

and, for  $k \neq \ell$ , define the pairwise multiclass averaging factor

$$L_{k\ell}^{\text{mc}} := \left[ \alpha_S(\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k}(1 - \ell_k) + \alpha_{\mu,\ell}(1 - \ell_\ell) + \ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2} + \sqrt{n_{\mathcal{T}}}(\Delta_k + \Delta_\ell) \right]^{-1}.$$

Let

$$\widehat{\text{SR}}_{k\ell}^{\text{mc}}(\mathcal{P}_{\text{dc}}) := \frac{\gamma_{k\ell}}{\left\| \mathbf{p}_k^{S,(\mathbf{t}^*)} - \boldsymbol{\mu}_k \right\|_2 + \left\| \mathbf{p}_\ell^{S,(\mathbf{t}^*)} - \boldsymbol{\mu}_\ell \right\|_2}.$$

If the denominator is zero, set  $\widehat{\text{SR}}_{k\ell}^{\text{mc}}(\mathcal{P}_{\text{dc}}) = +\infty$ . Then

$$\widehat{\text{SR}}_{k\ell}^{\text{mc}}(\mathcal{P}_{\text{dc}}) \geq \underline{\text{SR}}_{k\ell}^{\text{mc}} := \frac{\sqrt{n_{\mathcal{T}}} L_{k\ell}^{\text{mc}} \gamma_{k\ell}}{\sqrt{C_\kappa} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}, \quad k \neq \ell.$$

Consequently,

$$\mathcal{R}_{\text{mc}}(\mathcal{P}_{\text{dc}}) \leq \sum_{k=1}^K \pi_k \sum_{\ell \neq k} \exp \left( -\frac{\gamma_{k\ell}^2}{8 \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \Psi_{\text{mc}}(\underline{\text{SR}}_{k\ell}^{\text{mc}}) \right).$$

In particular, if  $\gamma_{\min} := \min_{k \neq \ell} \gamma_{k\ell}$  and  $\underline{\text{SR}}_{\min}^{\text{mc}} := \min_{k \neq \ell} \underline{\text{SR}}_{k\ell}^{\text{mc}}$ , then

$$\mathcal{R}_{\text{mc}}(\mathcal{P}_{\text{dc}}) \leq (K-1) \exp \left( -\frac{\gamma_{\min}^2}{8 \|\boldsymbol{\Sigma}_\xi\|_{\text{op}}} \Psi_{\text{mc}}(\underline{\text{SR}}_{\min}^{\text{mc}}) \right).$$

If all OvR activities are active throughout training, then  $L_{k\ell}^{\text{mc}}$  can be replaced by

$$L_{k\ell, \text{full}}^{\text{mc}} := \left[ \rho_S^{T_{\text{out}} T_{\text{in}}} (2\alpha_S + \alpha_{\mu,k} + \alpha_{\mu,\ell}) + (1 - \rho_S^{T_{\text{out}} T_{\text{in}}}) (p_k^{-1/2} + p_\ell^{-1/2}) \right]^{-1}$$

in the definition of  $\underline{\text{SR}}_{k\ell}^{\text{mc}}$ .

*Proof.* The signal–noise expansion in Lemma D.14 gives

$$\mathbf{p}_k^{S,(\mathbf{t}^*)} - \boldsymbol{\mu}_k = u_k^{(\mathbf{t}^*)} \mathbf{p}_k^{S,(0,0)} - (1 - c_{k,\mu}^{(\mathbf{t}^*)}) \boldsymbol{\mu}_k + \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \boldsymbol{\xi}_i.$$

The same coefficient recursion also gives  $u_k^{(\mathbf{t})} + c_{k,\mu}^{(\mathbf{t})} \leq 1$  for every time index: initially the sum is one, and one update changes it by at most

$$u_k^{(\mathbf{t}^+)} + c_{k,\mu}^{(\mathbf{t}^+)} \leq \left( 1 - 2\eta_S \bar{q}_k^{S,(\mathbf{t})} \right) \left( u_k^{(\mathbf{t})} + c_{k,\mu}^{(\mathbf{t})} \right) + 2\eta_S \bar{q}_k^{S,(\mathbf{t})} \leq 1,$$

because  $n_{\mathcal{T},k}^{-1} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \bar{b}_{k,i}^{(\mathbf{t})} \leq \bar{q}_k^{S,(\mathbf{t})}$ . Hence  $0 \leq c_{k,\mu}^{(\mathbf{t}^*)} \leq 1$  and  $1 - c_{k,\mu}^{(\mathbf{t}^*)} \leq 1 - \ell_k$ .

The upper bound on  $u_k^{(\mathbf{t}^*)}$  from Lemma D.14 gives  $u_k^{(\mathbf{t}^*)} \leq \bar{u}_k$ . On the event in Assumption B.1, the initialization and noise bounds give  $\|\mathbf{p}_k^{S,(0,0)}\|_2 \leq \sqrt{C_\kappa d} \sigma_S$ ,

$$\left\| \frac{1}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \boldsymbol{\xi}_i \right\|_2 \leq \sqrt{\frac{C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}{n_{\mathcal{T},k}}}, \quad \max_{i \in \mathcal{I}_k(\mathcal{T})} \|\boldsymbol{\xi}_i\|_2 \leq \sqrt{C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}.$$

Using the decomposition  $c_{k,i}^{(\mathbf{t}^*)} = \ell_k/n_{\mathcal{T},k} + h_{k,i}$  from Lemma D.14, with  $h_{k,i} \geq 0$  and  $\sum_i h_{k,i} \leq \Delta_k$ , we obtain

$$\left\| \sum_{i \in \mathcal{I}_k(\mathcal{T})} c_{k,i}^{(\mathbf{t}^*)} \boldsymbol{\xi}_i \right\|_2 \leq \ell_k \left\| \frac{1}{n_{\mathcal{T},k}} \sum_{i \in \mathcal{I}_k(\mathcal{T})} \boldsymbol{\xi}_i \right\|_2 + \sum_{i \in \mathcal{I}_k(\mathcal{T})} h_{k,i} \|\boldsymbol{\xi}_i\|_2 \leq \sqrt{C_\kappa \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d} \left( \frac{\ell_k}{\sqrt{n_{\mathcal{T},k}}} + \Delta_k \right).$$

Combining these three estimates and using  $n_{\mathcal{T},k} = p_k n_{\mathcal{T}}$  yields

$$\left\| \mathbf{p}_k^{S,(\mathbf{t}^*)} - \boldsymbol{\mu}_k \right\|_2 \leq \frac{\sqrt{C_\kappa} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}{\sqrt{n_{\mathcal{T}}}} \left[ \alpha_S \bar{u}_k + \alpha_{\mu,k}(1 - \ell_k) + \ell_k p_k^{-1/2} + \sqrt{n_{\mathcal{T}}} \Delta_k \right].$$

Therefore, for every pair  $k \neq \ell$ ,

$$\left\| \mathbf{p}_k^{\mathcal{S},(t^*)} - \boldsymbol{\mu}_k \right\|_2 + \left\| \mathbf{p}_\ell^{\mathcal{S},(t^*)} - \boldsymbol{\mu}_\ell \right\|_2 \leq \frac{\sqrt{C_\kappa} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}{\sqrt{n\mathcal{T}}} \frac{1}{L_{k\ell}^{\text{mc}}}.$$

By the definition of the pairwise multiclass ratio, this gives  $\widehat{\text{SR}}_{k\ell}^{\text{mc}}(\mathcal{P}_{\text{dc}}) \geq \underline{\text{SR}}_{k\ell}^{\text{mc}}$ .

Applying Lemma D.15 and using that  $\Psi_{\text{mc}}$  is nondecreasing on  $[0, +\infty]$  gives the pairwise risk bound. The compact bound with  $\gamma_{\min}$  and  $\underline{\text{SR}}_{\min}^{\text{mc}}$  follows from  $\gamma_{k\ell} \geq \gamma_{\min}$ ,  $\underline{\text{SR}}_{k\ell}^{\text{mc}} \geq \underline{\text{SR}}_{\min}^{\text{mc}}$ , and  $\sum_{k=1}^K \pi_k = 1$ . The full-activity version follows from the exact contraction in Lemma D.14: in that case  $\bar{u}_k = \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}$ ,  $1 - \ell_k = \rho_{\mathcal{S}}^{T_{\text{out}} T_{\text{in}}}$ , and  $\Delta_k = 0$  for every class  $k$ .  $\square$

The final deterministic comparison interprets the factor  $L_{k\ell}^{\text{mc}}$  in the certified ratio. It isolates the effect of active-window gaps through  $\mathcal{W}_k^{\text{OvR}}$  and shows when the pairwise bound retains the class-averaging scale.

**Lemma D.17.** *Assume the hypotheses of Corollary D.16. Set  $T := T_{\text{in}}$  and  $\rho := \rho_{\mathcal{S}}$ . For each class  $k$  and outer-loop index  $r$ , define*

$$D_{k,r} := \sum_{s=r+1}^{T_{\text{out}}-1} (T - T_{\mathcal{S},k,\text{lin}}^{\text{OvR}}(s)), \quad M_{\mathcal{S},k} := \max_{0 \leq r \leq T_{\text{out}}-1} D_{k,r},$$

and

$$\mathcal{W}_k^{\text{OvR}} := (1 - \rho^T)(1 - \rho^{M_{\mathcal{S},k}}) + \max_{0 \leq r \leq T_{\text{out}}-1} (1 - \rho^{T - T_{k,\text{lin}}^{\text{OvR}}(r)}).$$

Then, for every  $k \in [K]$ ,

$$0 \leq \Delta_k \leq B_{k,\text{out}} \mathcal{W}_k^{\text{OvR}}, \quad \ell_k \geq B_{k,\text{out}} ((1 - \rho^T) - \mathcal{W}_k^{\text{OvR}}).$$

For  $k \neq \ell$ ,

$$(L_{k\ell}^{\text{mc}})^{-1} = \alpha_{\mathcal{S}}(\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k}(1 - \ell_k) + \alpha_{\mu,\ell}(1 - \ell_\ell) + \sqrt{n\mathcal{T}}(\Delta_k + \Delta_\ell) + \ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}.$$

Moreover, the window lengths control the nonuniform coefficient contribution through

$$\Delta_k + \Delta_\ell \leq B_{k,\text{out}} \mathcal{W}_k^{\text{OvR}} + B_{\ell,\text{out}} \mathcal{W}_\ell^{\text{OvR}}.$$

The following consequences hold along any sequence of problem instances for which  $\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2} > 0$ .

(i) If

$$\alpha_{\mathcal{S}}(\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k}(1 - \ell_k) + \alpha_{\mu,\ell}(1 - \ell_\ell) = o\left(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}\right),$$

and

$$\Delta_k + \Delta_\ell = o\left(\frac{\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}}{\sqrt{n\mathcal{T}}}\right),$$

then

$$L_{k\ell}^{\text{mc}} = (1 + o(1)) \left(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}\right)^{-1},$$

and

$$\underline{\text{SR}}_{k\ell}^{\text{mc}} = (1 + o(1)) \frac{\sqrt{n\mathcal{T}} \gamma_{k\ell}}{(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}) \sqrt{C_\kappa} \|\boldsymbol{\Sigma}_\xi\|_{\text{op}} d}.$$

(ii) If

$$\alpha_{\mathcal{S}}(\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k}(1 - \ell_k) + \alpha_{\mu,\ell}(1 - \ell_\ell) = O\left(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}\right),$$

and

$$\Delta_k + \Delta_\ell = O\left(\frac{\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}}{\sqrt{n\mathcal{T}}}\right),$$

then

$$L_{k\ell}^{\text{mc}} = \Omega\left(\left(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}\right)^{-1}\right), \quad \underline{\text{SR}}_{k\ell}^{\text{mc}} = \Omega\left(\frac{\sqrt{n\mathcal{T}} \gamma_{k\ell}}{(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}) \sqrt{C_\kappa \|\Sigma_\xi\|_{\text{op}} d}}\right).$$

(iii) If  $\Delta_k + \Delta_\ell = \Omega(1)$ , then

$$L_{k\ell}^{\text{mc}} = O(n\mathcal{T}^{-1/2}), \quad \underline{\text{SR}}_{k\ell}^{\text{mc}} = O\left(\frac{\gamma_{k\ell}}{\sqrt{C_\kappa \|\Sigma_\xi\|_{\text{op}} d}}\right).$$

(iv) If the OvR training and condensed-sample windows for classes  $k$  and  $\ell$  have full length  $T$  at every outer loop, then

$$(L_{k\ell}^{\text{mc}})^{-1} = \rho^{T_{\text{out}}T} (2\alpha_S + \alpha_{\mu,k} + \alpha_{\mu,\ell}) + (1 - \rho^{T_{\text{out}}T}) (p_k^{-1/2} + p_\ell^{-1/2}).$$

Thus the full-window transition is governed by the comparison between the two terms in this denominator. In particular, if

$$\rho^{T_{\text{out}}T} (2\alpha_S + \alpha_{\mu,k} + \alpha_{\mu,\ell}) = O\left(\left(1 - \rho^{T_{\text{out}}T}\right) (p_k^{-1/2} + p_\ell^{-1/2})\right),$$

then the conclusion in (ii) holds with  $\ell_j = 1 - \rho^{T_{\text{out}}T}$  for  $j \in \{k, \ell\}$ . A sufficient scale for this regime, when  $0 < \rho < 1$ , is

$$T_{\text{out}}T = \Omega\left(\frac{\left[\log\left(\frac{2\alpha_S + \alpha_{\mu,k} + \alpha_{\mu,\ell}}{p_k^{-1/2} + p_\ell^{-1/2}}\right)\right]_+}{-\log \rho}\right).$$

*Proof.* For each  $k$  and  $r$ , set

$$F_{k,r} := \sum_{s=r+1}^{T_{\text{out}}-1} T_{S,k,\text{lin}}^{\text{OvR}}(s), \quad s_{k,r} := T - T_{k,\text{lin}}^{\text{OvR}}(r).$$

Then  $F_{k,r} = (T_{\text{out}} - 1 - r)T - D_{k,r}$ . The  $r$ th contribution to  $u_{k,\text{out}}$  is  $a_{k,r} := (1 - \rho^T) \rho^{F_{k,r}}$ , while the  $r$ th contribution to  $\ell_k$  is

$$b_{k,r} := \rho^{(T_{\text{out}}-r)T - T_{k,\text{lin}}^{\text{OvR}}(r)} (1 - \rho^{T_{k,\text{lin}}^{\text{OvR}}(r)}) = \rho^{F_{k,r} + D_{k,r} + s_{k,r}} (1 - \rho^{T - s_{k,r}}).$$

Therefore

$$a_{k,r} - b_{k,r} = \rho^{F_{k,r}} \left[ (1 - \rho^T) - \rho^{D_{k,r} + s_{k,r}} (1 - \rho^{T - s_{k,r}}) \right].$$

The bracket has the exact decomposition

$$(1 - \rho^T)(1 - \rho^{D_{k,r}}) + \rho^{D_{k,r}} (1 - \rho^{s_{k,r}}),$$

which is nonnegative and bounded above by  $\mathcal{W}_k^{\text{OvR}}$ . Summing over  $r$  yields

$$0 \leq u_{k,\text{out}} - \ell_k \leq \mathcal{W}_k^{\text{OvR}} \sum_{r=0}^{T_{\text{out}}-1} \rho^{F_{k,r}} = B_{k,\text{out}} \mathcal{W}_k^{\text{OvR}}.$$

Since  $\Delta_k = u_{k,\text{out}} - \ell_k$  and  $u_{k,\text{out}} = B_{k,\text{out}}(1 - \rho^T)$ , the first display follows.

By definition,

$$\begin{aligned} (L_{k\ell}^{\text{mc}})^{-1} &= \alpha_{\mathcal{S}}(\bar{u}_k + \bar{u}_\ell) + \alpha_{\mu,k}(1 - \ell_k) + \alpha_{\mu,\ell}(1 - \ell_\ell) \\ &\quad + \ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2} + \sqrt{n\mathcal{T}}(\Delta_k + \Delta_\ell). \end{aligned}$$

This proves the displayed identity for  $(L_{k\ell}^{\text{mc}})^{-1}$ . Summing the two bounds  $\Delta_j \leq B_{j,\text{out}} \mathcal{W}_j^{\text{OvR}}$  over  $j \in \{k, \ell\}$  gives the displayed control of the nonuniform coefficient contribution.

The two assumptions in (i) make the initialization, signal-recovery, and nonuniform coefficient terms jointly  $o(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2})$ , and the two assumptions in (ii) make them jointly  $O(\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2})$ . The first two asymptotic conclusions therefore follow by comparing the denominator with the class-averaging term  $\ell_k p_k^{-1/2} + \ell_\ell p_\ell^{-1/2}$  and then using the definition of  $\text{SR}_{k\ell}^{\text{mc}}$ . If  $\Delta_k + \Delta_\ell = \Omega(1)$ , then the denominator of  $L_{k\ell}^{\text{mc}}$  is at least  $\sqrt{n\mathcal{T}}(\Delta_k + \Delta_\ell) = \Omega(\sqrt{n\mathcal{T}})$ , which proves (iii). If the windows for  $k$  and  $\ell$  are full, then  $T_{\mathcal{S},j,\text{lin}}^{\text{OvR}}(r) = T_{j,\text{lin}}^{\text{OvR}}(r) = T$  for  $j \in \{k, \ell\}$  and all  $r$ . Hence  $\bar{u}_j = \rho^{T_{\text{out}}T}$ ,  $\ell_j = 1 - \rho^{T_{\text{out}}T}$ , and  $\Delta_j = 0$ , which gives the displayed full-window formula. The stated sufficient scale is the logarithmic form of the contraction requirement when  $0 < \rho < 1$ .  $\square$

Table 3: Default hyperparameters for synthetic experiments.

Item	Value
Dataset size	$n_{\mathcal{T}} = 1000$ (balanced), train/test split 60%/40%
Dimension	$d = 500$
Signal and noise	$\ \boldsymbol{\mu}\ _2 = 1.0$ , $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ with $\sigma = 1.0$
Regularization	$\lambda = 1.0$
Learning rates	$\eta_{\mathbf{w}} = 0.05$ , $\eta_{\mathcal{S}} = 0.25$
Bilevel steps	$T_{\text{out}} = 50$ , $T_{\text{in}} = 10$
Initialization scales	$\sigma_{\mathbf{w}} = 0.01$ , $\sigma_{\mathcal{S}} = 0.01$
Seeds	$n_{\text{seeds}} = 5$ with seeds $\{0, 1, 2, 3, 4\}$

Table 4: Default hyperparameters for the  $K$ -class OvR synthetic experiment.

Item	Value
Dataset size	$K = 5$ , 200 generated samples per class, train/test split 60%/40%
Dimension	$d = 500$
Class means and noise	Near-orthogonal class means with $\ \boldsymbol{\mu}_k\ _2 = 2.0$ , $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ with $\sigma = 1.0$
Regularization	$\lambda = 1.0$
Learning rates	$\eta_{\mathbf{w}} = 0.05$ , $\eta_{\mathcal{S}} = 0.25$
Bilevel steps	$T_{\text{out}} = 50$ , $T_{\text{in}} = 10$
Initialization scales	$\sigma_{\mathbf{w}} = 0.01$ , $\sigma_{\mathcal{S}} = 0.01$
Seeds	$n_{\text{seeds}} = 5$ with seeds $\{0, 1, 2, 3, 4\}$

## E Additional Simulation Details

This section provides the reproducibility details, numerical tables, and auxiliary diagnostics underlying the simulation results in Section 4.

### E.1 Experimental setup and reproducibility

Unless otherwise stated, each experiment compares the learned condensed set with a random one-shot coreset that selects one real training sample per class. We report mean  $\pm$  one standard deviation over the specified random seeds.

#### E.1.1 Parameter settings

Table 3, Table 4, and Table 5 summarize the default hyperparameters used in the binary synthetic, multiclass synthetic, and KMNIST experiments, respectively.

#### E.1.2 Environment and random seeds

All experiments are implemented in Python with PyTorch. Each script iterates `seed` in `range(n_seeds)` with default seeds  $\{0, 1, 2, 3, 4\}$ . At the beginning of each trial, we set both `torch.manual_seed(seed)` and `np.random.seed(seed)`; the same seed is propagated to synthetic data generation, the train/test split, random one-shot coreset sampling, and model initialization. This convention makes the reported Monte Carlo summaries reproducible from the experiment logs.

Table 5: Default hyperparameters for KMNIST experiments (two-class subset).

Item	Value
Dataset	KMNIST (class 0 vs class 1), official train/test split
Input	$d = 784$ (flattened), zero-mean unit-variance normalization
Regularization	$\lambda = 0.01$
Learning rates	$\eta_{\mathbf{w}} = 0.005, \eta_{\mathcal{S}} = 0.1$
Bilevel steps	$T_{\text{out}} = 10, T_{\text{in}} \in \{5, 10, 20, 30\}$
Initialization scales	$\sigma_{\mathbf{w}} = 0.01, \sigma_{\mathcal{S}} = 0.01$
Seeds	$n_{\text{seeds}} = 5$ with seeds $\{0, 1, 2, 3, 4\}$

Table 6: Synthetic binary aggregation results (mean  $\pm$  std over seeds). Alignment is  $\cos(\mathbf{s}(\mathcal{S}^*), \boldsymbol{\mu})$ , and SNR is the signal-to-residual proxy used in the synthetic diagnostic.

Method	Alignment	SNR proxy	SVM acc.	Nearest-centroid acc.
Condensed (GM)	$0.725 \pm 0.024$	$0.526 \pm 0.036$	$0.772 \pm 0.013$	$0.773 \pm 0.013$
Random one-shot	$0.070 \pm 0.028$	$0.006 \pm 0.004$	$0.517 \pm 0.029$	$0.506 \pm 0.010$
Full data	–	–	$0.768 \pm 0.021$	–

### E.1.3 Implementation conventions

The implementation follows the deterministic hinge indicator convention used in the theory:  $q(\mathbf{w}) = \mathbb{1}\{u(\mathbf{w}) < 1\}$ , with a strict inequality, and the indicator is not differentiated through. In the binary single-sample-per-class setting  $n_{\mathcal{S}} = 2$ , the condensed update at state  $\mathbf{w}$  uses the explicit formula  $\mathbf{z}_y \leftarrow \mathbf{z}_y + \eta_{\mathcal{S}} q_y^{\mathcal{S}}(\mathbf{w}) y (g_{\mathcal{S}}(\mathbf{w}) - g_{\mathcal{T}}(\mathbf{w}))$ , followed by the inner-loop model update  $\mathbf{w} \leftarrow \mathbf{w} - \eta_{\mathbf{w}} g_{\mathcal{T}}(\mathbf{w})$ . The random one-shot baseline constructs  $\mathcal{S}_{\text{rand}}$  by sampling one training point per class from  $\mathcal{T}$  and evaluates it under the same downstream training procedure as  $\mathcal{S}^*$ .

## E.2 Raw tables for synthetic condensed-sample aggregation

Table 6 reports the geometric and downstream quantities underlying Fig. 1b. The condensed set has much stronger alignment with the signal direction than the random one-shot baseline (0.725 versus 0.070), and its SVM test accuracy is comparable to full-data training (0.772 versus 0.768) while clearly exceeding random selection (0.518).

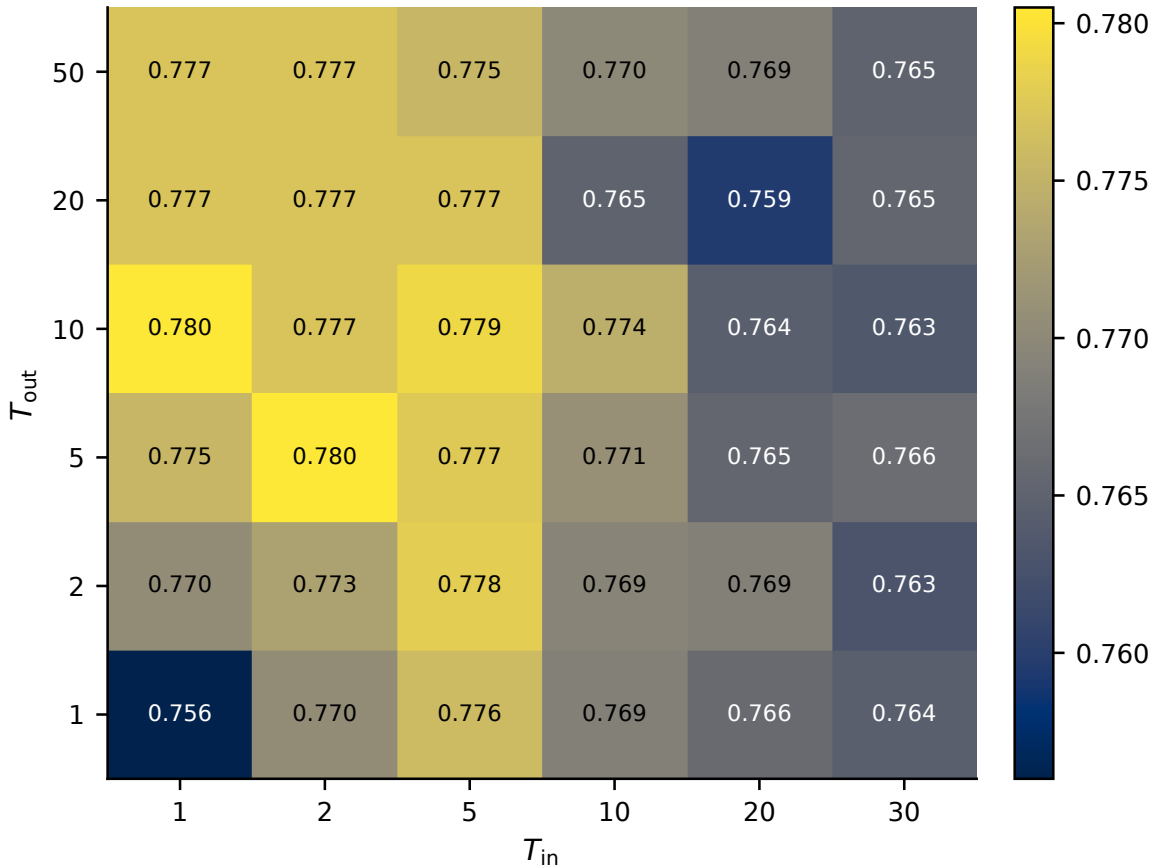
## E.3 Accuracy heatmap and grid values for the inner/outer sweep

Figure 3 complements Fig. 2a in the main text by visualizing downstream test accuracy over the same  $(T_{\text{in}}, T_{\text{out}})$  grid, and Table 7 lists the corresponding values. The peak accuracy 0.780 is attained for small  $T_{\text{in}} \in \{1, 2\}$  with moderate  $T_{\text{out}} \in \{5, 10\}$ ; larger inner-loop lengths lead to a modest but consistent reduction, in line with the weaker alignment values in Fig. 2a.

## E.4 Dynamics diagnostics

Figure 4 reports a representative stepwise diagnostic that tracks the training-set activity ratio  $\text{act}_{\mathcal{T}}$ , the condensed-set activity ratio  $\text{act}_{\mathcal{S}}$ , gradient norms, and the gradient-matching discrepancy as functions of the inner-loop step. In this default configuration,  $\text{act}_{\mathcal{S}}$  remains equal to one throughout the displayed trajectory, while  $\text{act}_{\mathcal{T}}$  and the training-gradient norm decrease along the inner loop. Thus the diagnostic illustrates a training-side loss of active gradient signal, rather than a condensed-sample activity switch.

The displayed diagnostic corresponds to the representative grid point  $(T_{\text{out}}, T_{\text{in}}) = (10, 10)$ , averaged over the five random seeds. Within each outer-loop restart, the training-set activity begins near one and then

Figure 3: Test accuracy over the  $(T_{in}, T_{out})$  grid on synthetic data, complementing Fig. 2a in the main text.Table 7: Grid values for Fig. 3: test accuracy over  $(T_{in}, T_{out})$  on synthetic data.

$T_{in} \backslash T_{out}$	1	2	5	10	20	50
1	0.756	0.770	0.775	0.780	0.777	0.777
2	0.770	0.773	0.780	0.777	0.777	0.777
5	0.776	0.778	0.777	0.779	0.777	0.775
10	0.769	0.769	0.771	0.774	0.765	0.770
20	0.766	0.769	0.765	0.764	0.759	0.769
30	0.764	0.763	0.766	0.763	0.765	0.765

decreases along the inner loop, reaching roughly 0.79 by the final inner step. The norm of the training gradient follows the same pattern, falling from about 1.39 at the beginning of an inner loop to about 0.53 near its end. By contrast, the condensed-set activity remains equal to one, so the condensed samples stay in the hinge-active region while the real-data gradient becomes progressively sparser. The gradient-matching discrepancy decreases rapidly during the early active part of the trajectory, but it stops decreasing monotonically once the training gradient has weakened. This behavior is consistent with the active-window interpretation in the theory: the early inner steps provide the main classwise averaging signal, whereas later steps operate with fewer active real samples and therefore contribute less stable signal accumulation.

### E.5 Multiclass OvR experiment

Table 8 and Fig. 5 report the  $K$ -class classwise OvR experiment described in Appendix B. Since the multiclass theory evaluates the learned set through nearest-prototype classification, the table uses the same evaluator

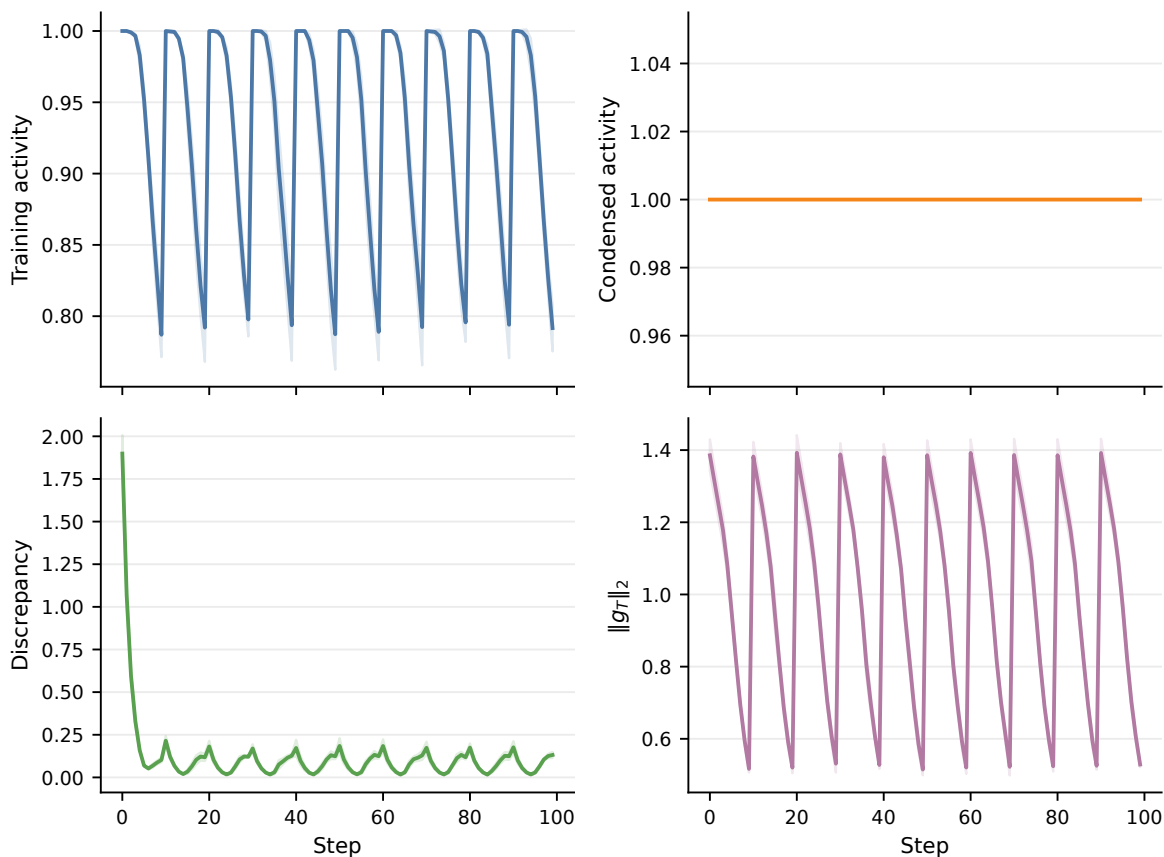


Figure 4: Representative dynamics diagnostic on synthetic data at  $(T_{\text{out}}, T_{\text{in}}) = (10, 10)$ . The condensed-set activity ratio remains active in this default setting, while the training-set activity and gradient norms decrease along the inner loop.

Table 8:  $K$ -class OvR condensation on synthetic additive data with one condensed sample per class, evaluated by nearest-prototype classification (mean  $\pm$  std over seeds).

Method	Test accuracy
GM-OvR	$0.624 \pm 0.008$
Random one-shot	$0.231 \pm 0.030$
Full-data nearest-prototype	$0.624 \pm 0.018$

for GM-OvR, the random one-shot baseline, and the full-data nearest-prototype reference whose prototypes are the empirical class means. The learned multiclass condensed set matches this full-data nearest-prototype reference within Monte Carlo variability and substantially outperforms random one-shot selection.

## E.6 Cross-model evaluation tables

Table 9 reports the mean $\pm$ std test accuracy underlying Fig. 2b, comparing the condensed set (GM) and the random one-shot baseline across multiple evaluators on KMnist. Table 10 provides the analogous cross-model comparison on synthetic data.

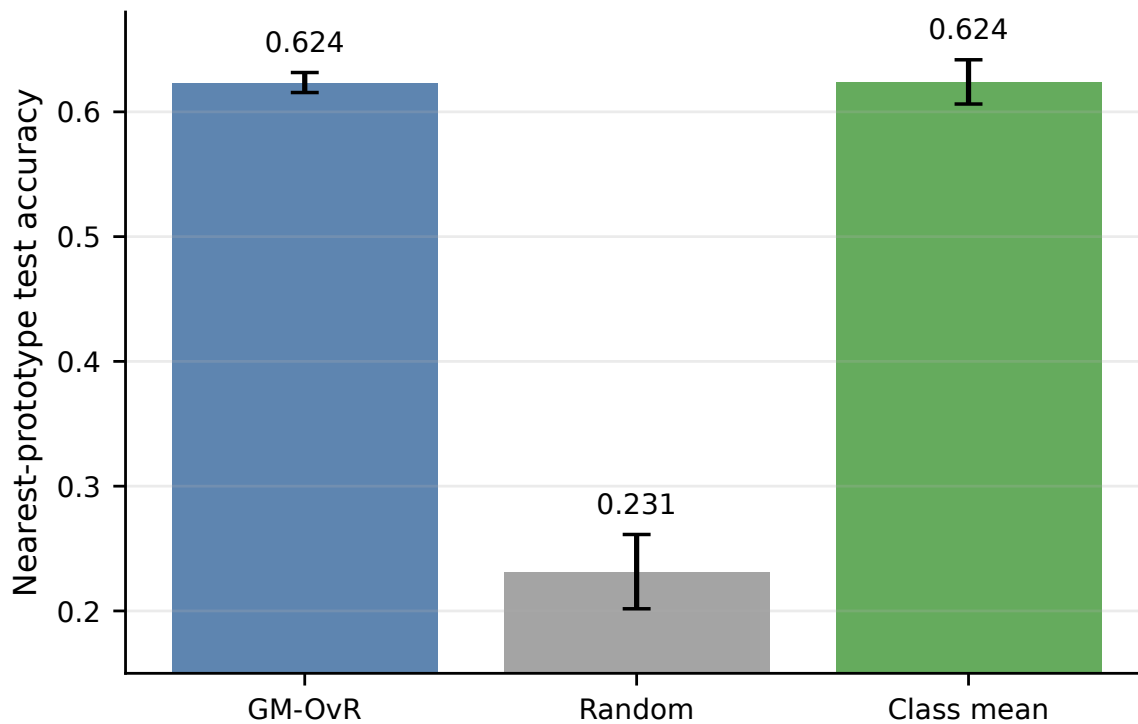


Figure 5: Multiclass OvR condensation with one condensed sample per class. The figure reports  $K$ -way test accuracy under the nearest-prototype evaluator, comparing the learned multiclass condensed set against a random one-shot coreset and the full-data nearest-prototype reference based on empirical class means.

Table 9: KMNIST cross-model evaluation (mean  $\pm$  std over seeds) on the same condensed set  $\mathcal{S}^*$  (GM) compared with a random one-shot coreset (Random).

Evaluator	GM	Random
Linear SVM	$0.960 \pm 0.006$	$0.828 \pm 0.062$
Logistic Regression	$0.960 \pm 0.005$	$0.838 \pm 0.053$
Nearest Centroid	$0.960 \pm 0.006$	$0.669 \pm 0.153$
Two-layer ReLU MLP	$0.953 \pm 0.003$	$0.811 \pm 0.059$

Table 10: Synthetic cross-model evaluation (mean  $\pm$  std over seeds) on the same condensed set  $\mathcal{S}^*$  (GM) compared with a random one-shot coreset (Random).

Evaluator	GM	Random
Linear SVM	$0.991 \pm 0.007$	$0.648 \pm 0.102$
Logistic Regression	$0.992 \pm 0.005$	$0.640 \pm 0.052$
Nearest Centroid	$0.991 \pm 0.007$	$0.598 \pm 0.038$
Two-layer ReLU MLP	$0.564 \pm 0.012$	$0.556 \pm 0.023$