# RISE: REGRESSION IMBALANCE HANDLING USING SWITCHING EXPERTS

## Anonymous authors

Paper under double-blind review

## **ABSTRACT**

Deep Imbalanced Regression (DIR) is challenging due to skewed label distributions and the need to preserve target continuity. Existing DIR methods rely on a single, monolithic model, yet empirical analysis shows that standard benchmarks exhibit strong distributional heterogeneity, exposing a core limitation of such approaches. We theoretically prove that this property creates an irreducible bias for any single model, leading to poor performance in data-scarce regions. This creates a core challenge for algorithmic fairness, as these regions often correspond to marginalized demographic groups. To address this, we propose RISE—Regression Imbalance handling via Switching Experts—a modular Mixture-of-Experts—inspired framework, theoretically motivated by our analysis. RISE employs a novel imbalance-aware algorithm to identify underperforming regions via validation loss and trains dedicated experts with targeted upsampling. As a complementary framework, RISE achieves new state-of-the-art performance while improving fairness, highlighting a principled new direction for imbalanced regression.

## 1 Introduction

Imbalanced data distributions—common in real-world settings—create severe challenges for regression models, producing high variance on minority labels and bias toward majority ones Wang et al. (2020); Gong et al. (2022). Unlike classification, where imbalance has been extensively studied, Deep Imbalanced Regression (DIR) is more complex due to its continuous and unbounded label space. This limitation has critical fairness implications: in healthcare, underestimating rapid disease progression delays care for underrepresented patients Cross et al. (2024), while in environmental policy, smoothing over pollution spikes overlooks harms concentrated in marginalized communities Su et al. (2024)—highlighting DIR as both a technical challenge and a fairness imperative in high-stakes domains.

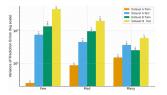
In Fig.1, we compare state-of-the-art (SOTA) methods for DIR, including LDS-FDSYang et al. (2021) and SRL Dong et al. (2025), on Dataset A Moschoglou et al. (2017). While these approaches reduce training error in tail (few-label) regions, their gains vanish at test time, revealing overfitting and poor generalization on underrepresented labels. Standard remedies such as frequency-based oversampling Steininger et al. (2021) partially close this gap in the tail but consistently degrade performance on head (many-label) regions, exposing a persistent head—tail trade-off Xu et al. (2021). A key observation is that performance across label bands is highly sensitive to the specific sampling realization of the training data, suggesting that the observed dataset is but one draw from a richer underlying distribution, and oversampling schemes represent alternative draws.

We hypothesize that the persistent head–tail discrepancy in DIR arises from two factors: (a) different label regions exhibit distinct, and sometimes conflicting, conditional distributions P(y|x); and (b) a single monolithic model lacks the capacity to jointly capture these heterogeneous mappings Sattler et al. (2020). We empirically validate distributional heterogeneity in standard DIR benchmarks, providing the first direct evidence in this setting. First, independent linear predictors trained on frozen ResNet-50 features for the many, medium-, and few-label bands of Dataset A and Dataset B Rothe et al. (2018) yield nearly orthogonal weight vectors, with cosine

**Table 1:** Cosine similarities

Dataset A	$w_{\mathrm{few}}$	$w_{\mathrm{med}}$	$w_{many}$
$w_{ m few} \ w_{ m med} \ w_{ m many}$	1.00	0.04	0.03
	0.04	1.00	0.09
	0.03	0.09	1.00
Dataset B	$w_{\mathrm{few}}$	$w_{\mathrm{med}}$	$w_{many}$
$w_{ m few} \ w_{ m med} \ w_{ m many}$	1.00	0.02	0.03
	0.02	1.00	0.18
	0.03	0.18	1.00

similarities as low as 0.03 (Table 1), indicating fundamentally different predictive functions across



**Figure 1:** Dataset A: SOTA DIR methods cut tail error but worsen head, exposing a persistent head–tail trade-off.

**Figure 2:** Heteroscedasticity in Model Error-SRL Dong et al. (2025)

regions. Second, to demonstrate the consequences of this heterogeneity, we analyze the error profile of a single global model SRL. We find pronounced heteroscedasticity in its prediction errors: on test data, variance in the few-label band is up to  $7\times$  higher than in the many-label band, while the opposite trend holds on training data—classic overfitting to scarce samples (Fig. 2). This instability arises precisely because a monolithic model cannot simultaneously capture distinct conditional distributions P(y|x) across regions. Together, these findings show that the core challenge in DIR is not merely label imbalance but distributional heterogeneity, motivating architectures that explicitly specialize across label regions.

This necessitates an architectural shift towards a multi-expert paradigm. We therefore propose RISE (Regression Imbalance handling using Switching Experts), a framework that directly confronts this challenge by learning specialized representations for different data regions. Crucially, RISE is not a generic Mixture of Experts (MoE) Mu & Lin (2025). Its novelty lies in its imbalance-aware algorithm that operationalizes the MoE specifically for DIR. Unlike generic MoEs that partition data by feature similarity, RISE identifies expert domains by analyzing the failure modes of a global model revealed through its validation loss. Each expert is then trained with targeted upsampling, ensuring it focuses on the underrepresented data that challenges a single, monolithic network. This end-to-end approach transforms the MoE from a general tool for heterogeneity into a targeted, principled solution for DIR.

Below we summarize our key contributions:

- 1. To the best of our knowledge, we are the first to identify and empirically validate that standard DIR benchmarks exhibit distributional heterogeneity, reframing the core challenge from mere label imbalance to representational imbalance.
- 2. We prove that any monolithic model in DIR suffers from an irreducible heterogeneity bias amplified by imbalance (Theorem 1), and show that targeted expert specialization trades bias reduction against estimation variance (Theorem 2).
- 3. Building on this, we propose **RISE**, a modular and model-agnostic framework that complements existing SOTA methods by explicitly addressing distributional heterogeneity, overcoming the persistent head–tail trade-off, and improving performance across all regions (as shown in Fig. 1).
- 4. RISE sets new SOTA on multiple DIR benchmarks Moschoglou et al. (2017), Rothe et al. (2018), outperforming all baselines, highlighting its effectiveness and establishing a new direction for DIR.

## 2 Imbalanced Regression Problem Formulation

In DIR, we are given a dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  with inputs  $x_i \in \mathcal{X} \subset \mathbb{R}^p$  and continuous labels  $y_i \in \mathcal{Y} \subset \mathbb{R}$ . The label marginal p(y) is highly non-uniform (long-tailed), producing majority and scarce (tail) regions where conventional models systematically fail. Motivated by empirical evidence (Sec. 1), we argue that the core difficulty is not merely imbalance in p(y), but a deeper distributional heterogeneity in the conditional  $P(y \mid x)$ . We posit a latent partition of the problem space into K regions, with region k comprising fraction  $\rho_k = n_k/n$  of the data and governed by a distinct conditional distribution  $P_k(y \mid x)$ . Because the fractions  $\{\rho_k\}$  are highly non-uniform, a single monolithic predictor trained on the pooled data is dominated by majority regions and induces a persistent bias in scarce ones, a limitation we formalize in Theorem 1. This heterogeneity makes a MoE architecture the natural modeling choice. We therefore model the global conditional distribution

as a mixture of these latent, region-specific distributions:

$$P(y|x) = \sum_{k=1}^{K} \pi_k(x) P_k(y|x),$$
(1)

where each component  $P_k(y|x)$  is modeled by an expert network  $E_k$  and the mixing coefficients  $\pi_k(x)$  are determined by a gating network  $g_{\phi}$ . The final prediction is the expectation under this mixture:  $\hat{y} = \sum_{k=1}^K g_{\phi}(x)_k \cdot E_k(x)$ . The learning task is thus transformed from fitting a single complex function into discovering this latent partition (the gate) and mastering each sub-problem (the experts), even when data is sparse—the core challenge our RISE framework is designed to solve.

#### 2.1 RELATED WORK

Deep Imbalanced Regression: DIR is challenging as it must preserve label continuity under skewed distributions. Prior methods modify loss functions or label densities: LDS-FDS Yang et al. (2021) and Balanced-MSE Ren et al. (2022) address global imbalance but ignore local heterogeneity; RankSim Gong et al. (2022), ConR Keramati et al. (2024), and SRL Dong et al. (2025) add feature-space regularization (ranking, contrastive, or latent uniformity) yet assume homogeneous features. Regression-via-classification approaches Pintea et al. (2023); Pu et al. (2025); Xiong & Yao (2024) discretize the target and leverage classification training, aligning with MoE formulations, but lack mechanisms to detect minority regions or train dedicated experts, limiting their ability to capture distributional heterogeneity.

Ensembling and Mixture of Experts: A common approach to imbalance is partitioning data by class sizes and training separate experts. Ensemble-based methods Xiang et al. (2020); Cui et al. (2023); Cai et al. (2021) follow this strategy in classification but do not extend naturally to regression, where targets are continuous and lack softmax-style aggregation. In long-tailed recognition, multi-expert models such as BBN Zhou et al. (2020) (two-branch fusion for head/tail) and RIDE Wang et al. (2020) (diversity-regularized experts) reduce bias, yet their applicability to DIR—where label continuity and regional heterogeneity are central—remains unexplored.

#### 3 THEORETICAL INSIGHTS: WHY MONOLITHIC MODELS FAIL ON DIR

We formalize the core difficulty we empirically observe in DIR: when data comes from a mixture of region-specific mechanisms, a single global predictor suffers cross-region interference, amplified by label imbalance. To study this, we adopt a simplified linear regression setting, a standard tool for analyzing generalization in complex models Belkin et al. (2018); Lin et al. (2023).

**Setup.** We consider heterogeneous linear regression with K latent regions, each occurring with probability  $\rho_k = n_k/n$  (Sec. 2). For a sample (x,y) from region k, such that  $x \sim \mathcal{N}(0,\Sigma)$ , and  $y = w_k^{*\top}x + \varepsilon$ , where  $w_k^* \in \mathbb{R}^p$  is the region-specific parameter,  $\varepsilon \sim \mathcal{N}(0,\sigma_k^2)$  is independent noise, and  $\Sigma \succ 0$  is the common feature covariance matrix  $^1$ . Heterogeneity is captured entirely by  $\{w_k^*\}$ , which define distinct  $P_k(y \mid x)$ . Stacking all  $n = \sum_{k=1}^K n_k$  samples gives the design matrix  $X \in \mathbb{R}^{n \times p}$  and the label vector  $Y \in \mathbb{R}^n$ . The pooled(or global) Ordinary Least Squares (OLS) estimator is  $\widehat{w} = (X^\top X)^{-1} X^\top Y$ , trained on all n samples. We evaluate performance by the region-weighted generalization error:  $\mathcal{G}_{\rho}(\widehat{w}) = \sum_{k=1}^K \rho_k \|\widehat{w} - w_k^*\|^2$ .

**Theorem 1** (Generalization error under imbalance and heterogeneity). Let  $w_{\text{avg}} = \sum_{k=1}^{K} \rho_k w_k^*$  and  $\bar{\sigma}^2 = \max_k \sigma_k^2$ . Under Gaussian design with n > p+1, the expected region-weighted error of the pooled OLS estimator decomposes as

$$E[\mathcal{G}_{\rho}(\widehat{w})] = \underbrace{\frac{\bar{\sigma}^{2} \operatorname{tr}(\Sigma^{-1})}{n-p-1}}_{Estimation \ Variance \ (shrinks \ with \ n)} + \underbrace{\sum_{k=1}^{K} \rho_{k} \|w_{k}^{*} - w_{\operatorname{avg}}\|^{2}}_{Heterogeneity \ Bias \ (persists)}, \tag{2}$$

<sup>&</sup>lt;sup>1</sup>In Appendix A we relax this assumption to region-dependent covariances  $\Sigma_k$  and noise  $\sigma_k^2$  and show the same qualitative conclusions hold.

**Proof sketch.** The decomposition follows from  $\mathcal{G}_{\rho}(\widehat{w}) = \|\widehat{w} - w_{\mathrm{avg}}\|^2 + \sum_k \rho_k \|w_k^* - w_{\mathrm{avg}}\|^2$ , since  $\sum_k \rho_k (w_k^* - w_{\mathrm{avg}}) = 0$ . The first term is bounded using inverse-Wishart moments for Gaussian design, yielding the variance term. The second term is deterministic and captures irreducible heterogeneity. Full derivations, and generalizations to  $\Sigma_k$ ,  $\sigma_k^2$  are provided in Appendix A.

**Implications.** Theorem 1 shows that imbalance amplifies heterogeneity:  $w_{\text{avg}}$  is dominated by head regions, yielding persistent error on tails when  $w_t^*$  lies far away. Even with infinite data, a monolithic model converges to this biased average. Since the Heterogeneity Bias cannot be reduced by more data or reweighting, a natural remedy is architectural: partition the space and assign specialized predictors, so each operates in a more homogeneous region and achieves better generalization.

## 4 PROPOSED METHOD: RISE

Our proposed method, **RISE**, as illustrated in Fig. 3, operates as an orthogonal meta-framework designed to systematically enhance any pre-trained DIR baseline. Its core architectural choice—replacing a single monolithic model with a system of specialized experts—is a direct response to the distributional heterogeneity we identified in Sec. 1. First, RISE-Identify takes the trained baseline model ( $f_{\theta}$ ) and analyzes its performance on a held-out validation set to discover its specific failure modes. By using held-out data, we identify regions of true generalization error, not artifacts of training set memorization. Second, RISE-Train creates a set of dedicated experts, each one targeting a specific failure region identified in the first stage. These experts are trained on the train-dataset with targeted upsampling, a strategy that encourages specialization while regularizing against overfitting. Finally, RISE-Inference learns a gating mechanism, also on the held-out set, that dynamically routes new inputs to the most appropriate expert at test time. Complete implementation details and pseudo code are provided in Appendix D.1.

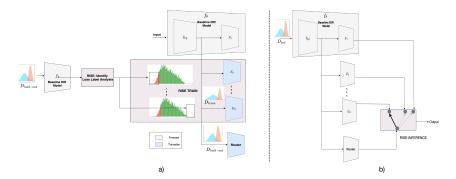


Figure 3: Overview of RISE framework.

#### 4.1 RISE-IDENTIFY: DISCOVERING LATENT FAILURE REGIONS

The first stage of RISE is to identify the latent regions where a baseline model fails, corresponding to the distinct components of the heterogeneous data distribution we posited in our problem formulation. The overall dataset  $\mathcal D$  is first split into a training set  $\mathcal D_{\text{train}}$  and a held-out validation set  $\mathcal D_{\text{val}}$ . A naive approach, implicitly used by frequency-based methods Cui et al. (2023); Yang et al. (2021), is to partition data using label-density bins from  $\mathcal D_{\text{train}}$ . Specifically, the continuous label space is first discretized into bins, and the frequency of labels in each bin is computed Yang et al. (2021). A K'-component Gaussian Mixture Model (GMM) is then fitted to these frequencies:

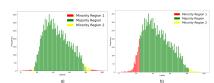
**Table 2:** Frequency-Loss Relationship Analysis for Dataset A

Label Band	Freq	Held-out Loss
0-20	231	8.86
20-40	4,913	6.30
40-60	4,609	7.44
60-80	2,244	7.79
80-100	208	9.34

 $p(\nu) = \sum_{j=1}^{K'} \pi'_j \mathcal{N}(\nu | \mu_j, \sigma_j^2)$  where  $\nu$  denotes the bin frequency. Component with the largest weight  $\pi'_j$  corresponds to the majority region, while the remaining components capture minority regions. However, this frequency-based approach is a flawed proxy for two key reasons. First, as our analysis shows in Table 2, error (or held-out loss) and frequency are not perfectly correlated; a region can have moderate data density yet still exhibit high generalization error. The 40-60 label band shows

higher loss (7.44) than the 20–40 band (6.30), despite similar sample sizes, indicating that frequency alone does not explain model error—performance is not strictly inversely proportional to frequency, highlighting the role of distributional heterogeneity. Second, frequency-based partitioning often creates non-contiguous regions in the label space as shown in Fig. 4a, which is problematic for regression tasks where nearby labels are highly correlated and should be modeled coherently Yang et al. (2021); Gong et al. (2022).

RISE adopts a more direct and principled strategy: we identify regions based on the model's generalization error, a direct signal of where the single, monolithic model is failing. First, we take a pre-trained DIR baseline,  $f_{\theta}$ , trained on  $\mathcal{D}_{\text{train}}$ . We then use this model to make predictions on the disjoint  $\mathcal{D}_{\text{val}}$ . For each sample  $(x_i, y_i) \in \mathcal{D}_{\text{val}}$ , we compute its pointwise prediction error,  $e_i = \mathcal{L}(f_{\theta}(x_i), y_i)$ , where  $\mathcal{L}$  is a loss function such as the absolute error (L1) or squared error (L2). To identify contiguous regions of high error, we model the joint distribution of these errors and their corresponding labels. This joint modeling ensures that identified regions are contiguous



**Figure 4:** Comparison of minority region identification approaches on Dataset A. a) Frequency-based analysis leads to disconnected minority regions (red, yellow) separated by majority regions (green). b) Loss-Label Distribution analysis produces contiguous minority regions

in label space—nearby labels with similar error patterns are grouped together—which is crucial for regression tasks where adjacent target values should be handled by similar predictive functions. Following Yang et al. (2021) we partition the continuous label range of  $\mathcal{D}_{\text{val}}$  into B disjoint, uniform-width bins,  $\{B_1,\ldots,B_B\}$ . For each bin b, we define the set of sample indices it contains as  $\mathcal{I}_b = \{i \mid y_i \in B_b\}$  and compute its average generalization error:  $\ell_b = \frac{1}{|\mathcal{I}_b|} \sum_{i \in \mathcal{I}_b} e_i$ . We model the resulting distribution of (average error, label bin center) pairs,  $\{(\ell_b, y_b)\}_{b=1}^B$ , using a K'-component GMM:

$$p(\ell, y) = \sum_{j=1}^{K'} \pi'_j \mathcal{N}((\ell, y) | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j).$$
(3)

The GMM naturally clusters the label bins into K' distinct and contiguous performance regions  $\mathcal{R}'_j$ . The component with the lowest mean error (along the error dimension of  $\mu_j$ ) is designated the well-performing "majority" region, while the remaining components correspond to distinct failure modes requiring specialized experts. K' is a hyperparameter that needs to be tuned. Using a held-out set ensures the identified regions correspond to true generalization failures, not training set memorization—a phenomenon we explicitly observe (Fig. 1). As illustrated in Fig. 4b, this approach produces continuous minority regions, aligning with the principle of region similarity and enabling more homogeneous expert training. By defining regions based on error, we directly target the heterogeneity bias identified as the key limitation in Theorem 1.

#### 4.2 RISE-TRAIN: EXPERT TRAINING

Having identified the baseline model's failure regions, the next stage is to train a dedicated expert for each one. To maintain computational efficiency and leverage the powerful representations learned by the baseline  $f_{\theta}$ , we adopt a parameter-efficient fine-tuning approach Kirichenko et al. (2023). As shown in Fig. 3 each expert  $E_i$  shares the frozen backbone of the pre-trained model; only its final layers are trained for specialization. For RISE-Train, we evaluate two strategies for expert specialization. A naive strategy (T1: Subgroup-Specific Training) trains each expert exclusively on its assigned data partition. This hard partitioning forces experts to learn from a severely restricted support of the data distribution, inducing high estimation variance and overfitting. It also prevents learning smooth functions across the label space in regression, leading to poor generalization and discontinuities at region boundaries. We therefore propose a more robust and principled strategy: T2: Cross-Group Training with Upsampling. For each expert  $E_i$ , we train on the full dataset  $\mathcal{D}_{\text{train}}$ with sample weights: target region samples  $\mathcal{R}'_{i}$  are upsampled by  $\alpha_{i} > 1$ , others by 1. This weighted empirical risk minimization both regularizes and specializes: exposure to the full dataset prevents high variance and discontinuities by constraining experts to remain well-behaved across the manifold, while  $\alpha_i > 1$  amplifies gradients from  $\mathcal{R}'_i$ , biasing the expert toward its designated failure regions. Our ablations (Sec. 6.3) empirically confirm that T2 is a far superior strategy, and we adopt it for all experiments.

## 4.3 RISE-INFERENCE: EXPERT SWITCHING STRATEGY

The final stage of RISE is to dynamically route each new input to the most suitable expert at test time. We compare three strategies: (I1) Expert Averaging, a simple ensemble baseline which aggregates predictions from all experts via weighted averaging; (I2) Train-based Router, a gating network  $g_{\phi}$  trained on the training set  $\mathcal{D}_{\text{train}}$ ; and our proposed (I3) Held-out-based Router, a gating network  $g_{\phi}$  trained on the held-out validation set  $\mathcal{D}_{\text{val}}$ . Our ablations (Sec. 6.4) show that the held-out router (I3) is decisively superior. A router trained on  $\mathcal{D}_{\text{train}}$  tends to select experts that best fit training artifacts, whereas training on  $\mathcal{D}_{\text{val}}$  turns routing into a meta-learning task: it learns to pick the expert that generalizes best. We therefore adopt I3 as our standard strategy. The router is implemented as a small multi-layer perceptron (MLP) that takes the shared features from the baseline's backbone as input and outputs a probability distribution over the K' experts. The final prediction  $\hat{y}$  is the output of the single expert selected by the router:  $\hat{y} = E_{\hat{y}^*}(x)$ , where  $\hat{y}^* = \arg \max_{\hat{y}} g_{\phi}(x)_{\hat{y}}$ .

## 5 THEORETICAL JUSTIFICATION FOR RISE

Having established in Theorem 1 that a global model suffers an irreducible heterogeneity bias, the natural question is: under what conditions can a region-specialized architecture overcome this limitation? We provide a formal result showing when RISE strictly outperforms the pooled model.

**Theorem 2** (Generalization advantage of RISE). Building on the heterogeneous regression setup of Theorem 1, let K' experts be trained with per-region upsampling factors  $\alpha_j$  and routing probabilities  $q_k(j)$  (the probability that a sample from region k is assigned to expert j). The effective sample size for expert j is  $n_{\text{eff}}^{(j)} = (\alpha_j - 1)n_j + n$ . From Theorem 1 we know the pooled/global model incurs region-weighted risk (or generalisation error)  $\mathcal{G}_{\text{pooled}} = V_{\text{glob}} + \Delta_{\text{glob}}$ , where  $\Delta_{\text{glob}} = \sum_{k=1}^K \rho_k \|w_k^* - w_{\text{avg}}\|$ , and  $V_{\text{glob}} = O(p/n)$  is the estimation variance of the global model, while generalisation error of RISE satisfies

$$\mathcal{G}_{RISE} = B_{det}(\alpha, q) + V_{est}(\alpha, q) + R_{cross}(\alpha, q),$$

where  $B_{\rm det}$  is deterministic bias from imperfect specialization (including possible  $K' \neq K$  or overlapping experts),  $V_{\rm est}(\alpha,q) = O(p/n_{\rm eff}^{(j)})$  is expert estimation variance, and  $R_{\rm cross}(\alpha,q) = O(\sqrt{p/n_{\rm eff}^{(j)}})$  are vanishing cross-terms. RISE outperforms the pooled model whenever

$$\Delta_{\text{glob}} - B_{\text{det}}(\alpha, q) > V_{\text{est}}(\alpha, q) - V_{\text{glob}} + R_{\text{cross}}(\alpha, q).$$

**Proof Sketch and Implications.** The pooled model converges to the data-weighted average  $w_{\rm avg}$ , incurring a persistent heterogeneity bias  $\Delta_{\rm glob}$ . RISE reduces this bias by upsampling scarce regions and routing them to specialized experts, so their effective targets move closer to  $w_k^*$ . Any mismatch between the number of experts and true regions  $(K' \neq K \text{ or overlaps})$  is absorbed into the deterministic bias term  $B_{\rm det}(\alpha,q)$ . The trade-off is increased finite-sample variance  $V_{\rm est}(\alpha,q) = O(p/n_{\rm eff}^{(j)})$  and negligible cross-terms  $R_{\rm cross}(\alpha,q) = O(\sqrt{p/n_{\rm eff}^{(j)}})$ , both of which decay with sample size. Thus, whenever the bias reduction dominates these penalties, RISE achieves strictly better generalization than the pooled model. Detailed proofs are in Appendix B. In practice, imbalance-aware upsampling (RISE-Train T2) increases  $n_{\rm eff}$  in scarce regions and the learned router (RISE-Inference I3) keeps the maximum routing error  $\epsilon = \max_k (1 - q_k(k))$  small, directly satisfying the theorem's condition  $\mathcal{G}_{\rm RISE} < \mathcal{G}_{\rm pooled}$ . We provide empirical validation of this effect in Sec. 6.5.

## 6 EXPERIMENTS AND RESULTS

We evaluate the utility of RISE through the following research questions-

- **RQ1:** How effective is RISE compared to SOTA baselines across different datasets?
- **RQ2:** How do expert training strategies and hyperparameters affect RISE performance?
- RQ3: How do different RISE-INFERENCE strategies affect overall performance?
- **RQ4:** How Practically Achievable are the Theoretical Conditions (Theorem 2) for RISE's Success?
- RQ5: Do RISE's performance gains stem from its specialized architecture or model capacity?

#### 6.1 EXPERIMENTAL SETUP

**Algorithms:** We compare RISE with four SOTA DIR methods and a *Vanilla* ResNet-50 backbone He et al. (2016). Since RISE is a modular framework that complements existing approaches, we evaluate it in combination with *Vanilla*, *LDS+FDS* Yang et al. (2021), *RankSIM* Gong et al. (2022), *BalancedMSE* Ren et al. (2022), and *SRL* Dong et al. (2025). For baselines we use released weights or official implementations. All RISE experts are trained with MSE loss. We tune the number of experts K', upsampling ratio  $\alpha$ , selecting the best configuration by validation performance. Further details are in Appendix D.2.

**Datasets:** We evaluate RISE on four DIR benchmarks across modalities: Dataset A Moschoglou et al. (2017)(images, target values in range 0–101), Dataset B Rothe et al. (2018) (images, range 0–186), STS-B Cer et al. (2017) from GLUE Wang et al. (2018) (text, similarity 0–5), and UCI-Abalone Nash et al. (1994) (tabular, range 1–29). Following confidentiality requirements, we anonymize Dataset A and Dataset B by omitting their names. Full details are in Appendix C.

Metrics: Following Yang et al. (2021); Gong et al. (2022); Dong et al. (2025), we report performance overall and across Many (>100 samples), Medium (20–100), and Few (<20) label bands. For Dataset A and B, we use Mean Absolute Error (MAE)↓, Mean Squared Error (MSE)↓, and Geometric Mean Error (GMEAN)↓. For STS-B, we additionally report Pearson↑ and Spearman↑ correlation. To assess fairness — defined as minimizing performance disparities across these bands—we also report balanced-MAE (bMAE)↓ Ren et al. (2022), which averages MAE over uniformly partitioned label bins to capture regional performance gaps (see Appendix Section E.2).

#### 6.2 RO1: Performance of RISE on public benchmark datasets

Table 4 shows that RISE consistently improves strong baselines (LDS+FDS, RankSIM, SRL) on Dataset A across all label bands (similar results for other datasets are provided in Appendix Sec. E.1). The largest relative gains occur in the Few and Medium regions, where monolithic models suffer most. For example, SRL+RISE reduces Few-MAE by 15% while simultaneously lowering Many-MAE by 10%, thereby overcoming the common head-tail performance trade-off. The performance gains from RISE scale directly

**Table 3:** Balanced-MAE (bMAE) ↓ on Dataset A

		bMa	AE↓	
Method	All	Many	Med	Few
SRL SRL + RISE	8.32 <b>7.39</b>	6.64 <b>6.00</b>	8.34 <b>7.25</b>	11.74 <b>10.33</b>

with the quality of the learned router. Weak backbones (e.g., *Vanilla*, with a router accuracy of  $\approx 0.44$ ) lead to unstable tail performance. In contrast, strong backbones (e.g., *SRL*, with a router accuracy of  $\approx 0.87$ ) enable RISE to fully realize the theoretical advantage of specialization (Theorem. 2). This confirms that the benefit from reducing heterogeneity bias dominates once the routing error is sufficiently low, while the variance cost remains controlled. Additional results (Appendix E.4) show that using an optimal router trained on the best feature representation yields significantly better performance than the baseline router, due to higher routing accuracy.

We assess fairness via bMAE in Table 3 (full results in Appendix Sec. E.2). By significantly improving Few and Medium-band performance while preserving Many-band accuracy, *SRL+RISE* directly mitigates the bias towards head regions exhibited by the baseline. This reduces performance disparities across label bands and demonstrably more equitable performance across all label bands.

## 6.3 RQ2: ABLATION ON EXPERT TRAINING AND HYPERPARAMETERS

We ablate RISE's core design choices on Dataset A with SRL as backbone in Tables 5 and 6. Results on Dataset B is in Appendix Sec. E.3. Our adopted expert training strategy, T2 (full-dataset training with region specific upsampling), consistently outperforms T1 (region specific training). T1's hard partitioning causes severe overfitting, whereas T2's full-dataset exposure acts as a powerful regularizer that promotes smooth generalization while upsampling encourages specialization. Our analysis of the number of experts (K') and upsampling ratio ( $\alpha$ ) reveals a clear U-shaped performance curve. This empirically validates our theory's cost-benefit trade-off (Theorem 2) and directly operationalizes it: the upsampling factor  $\alpha$  is a key lever to control the expert's estimation variance ( $V_{\rm est}$ ) while still achieving the primary goal of reducing heterogeneity bias ( $B_{\rm det}$ ). Performance peaks at moderate values (e.g., K'=3,  $\alpha=3$ ) before degrading as the costs of data fragmentation and overfitting outweigh the benefits of heterogeneity reduction.

**Table 4:** Results on Dataset AMoschoglou et al. (2017). For each baseline/RISE pair, the better score is in **bold**; the best overall is underlined. Router accuracy of RISE is shown in parentheses.

	L1 (MAE) ↓				GME	EAN↓		$MSE\downarrow$				
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few
VANILLA +RISE (0.44)	11.05 <b>10.43</b>	9.96 <b>9.40</b>	12.79 <b>11.62</b>	<b>16.53</b> 16.93	7.06 <b>6.55</b>	6.27 <b>5.85</b>	8.37 <b>7.47</b>	13.48 <b>13.16</b>	202.09 <b>181.61</b>	165.09 <b>148.38</b>	270.75 <b>221.57</b>	<b>361.74</b> 384.95
BalancedMSE +RISE (0.47)	8.70 <b>7.71</b>	8.44 <b>7.23</b>	8.99 <b>8.16</b>	10.26 <b>10.02</b>	5.58 <b>4.83</b>	5.44 <b>4.52</b>	5.87 <b>5.10</b>	<b>6.17</b> 6.87	127.05 <b>103.39</b>	118.69 <b>91.14</b>	133.94 <b>114.84</b>	<b>187.01</b> 187.41
LDS+FDS +RISE (0.53)	7.47 <b>7.28</b>	6.91 <b>6.79</b>	8.27 <b>8.07</b>	10.58 <b>9.72</b>	4.77 <b>4.49</b>	4.44 <b>4.25</b>	5.33 <b>4.88</b>	6.87 <b>6.04</b>	95.32 <b>92.79</b>	79.71 <b>78.88</b>	118.52 <b>116.49</b>	178.58 <b>158.63</b>
RankSIM +RISE (0.54)	7.02 <b>6.94</b>	6.58 <b>6.50</b>	7.86 <b>7.38</b>	9.72 <b>9.10</b>	4.55 <b>4.35</b>	4.14 <b>4.08</b>	5.39 <b>4.80</b>	6.97 <b>6.04</b>	83.55 <b>82.70</b>	74.34 <b>71.96</b>	99.30 <b>91.20</b>	149.51 <b>138.15</b>
SRL +RISE (0.87)	7.23 <b>6.57</b>	6.64 <u><b>6.16</b></u>	8.28 <b>7.36</b>	9.85 <b>8.30</b>	4.53 <u>3.61</u>	4.17 <b>3.40</b>	5.32 <b>4.14</b>	6.35 <b>4.33</b>	91.79 <b>82.01</b>	77.20 <b>70.88</b>	115.83 <b>91.20</b>	163.15 <b>134.93</b>

**Table 5:** Ablation on Dataset A: MAE for varying upsampling (with fixed K' = 3, left) and varying experts (with fixed  $\alpha = 3$ , right). Best RISE configuration beating baseline SRL is in **bold**.

	L1 (MAE) ↓					
Config	All	Many	Med	Few		
SRL	7.23	6.64	8.28	9.85		
SRL+RISE						
$\alpha=2$	6.72	6.23	7.69	8.66		
$\alpha=3$	6.57	6.16	7.36	8.30		
$\alpha=4$	6.73	6.32	7.51	8.43		
$\alpha$ =5	6.89	6.52	7.49	8.68		

	L1 (MAE) ↓					
Config	All	Many	Med	Few		
SRL	7.23	6.64	8.28	9.85		
SRL+RISE						
K'=2	6.88	6.41	7.70	9.06		
K'=3	6.57	6.16	7.36	8.30		
K'=4	6.89	6.48	7.38	9.29		
K'=5	7.29	6.93	7.67	9.58		

#### 6.4 RQ3: ABLATION ON RISE-INFERENCE STRATEGIES

We compare three routing strategies as mentioned in Section 4.3 on Dataset A in Table 7 (full results in Appendix Sec. E.3): **11** (expert averaging), **12** (router trained on the training set), and our proposed **13** (router trained on a held-out validation set). We observe that I3 is significantly superior. The reason is fundamental—routers trained on the training set (I2) overfit to features already captured by experts, whereas I3 learns which expert generalizes best, providing a robust signal for routing. To confirm the gain is not simply from more data, we trained another router-variant on the full dataset (train+held-out); while this improved over the baseline SRL, it was still outperformed by RISE with the I3 router. This confirms that RISE's advantage stems from its effective use of held-out data for what is essentially a meta-learning task—learning to select the best generalizing expert—not from access to additional training samples.

**Table 6:** Ablation of RISE-Expert training. Best results in **bold**. Full results in Table 14

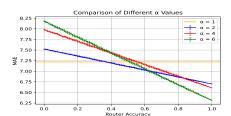
		L1 (M	AE)↓	
Method	All	Many	Med	Few
RISE (T1) RISE (T2)	7.23 <b>6.57</b>	6.77 <b>6.16</b>	7.95 <b>7.36</b>	9.61 <b>8.30</b>

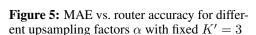
**Table 7:** Ablation of RISE inference strategies. Best results in **bold**. Full results in Table 15.

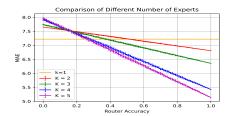
	L1 (MAE) ↓						
Method	All	Many	Med	Few			
Baseline SRL	7.23	6.64	8.28	9.85			
Expert average (I1)	7.23	6.72	8.13	9.54			
Train-based Router (I2)	7.26	6.61	8.34	10.33			
Held-out-based Router (I3)	6.57	6.16	7.36	8.30			
Train+Held-out Router	7.18	6.62	8.15	9.84			

#### 6.5 RQ4: EMPIRICAL VALIDATION OF THEORETICAL TRADE-OFFS IN PRACTICE

Theorem 2 predicts that RISE outperforms a pooled (or monolithic) model whenever the bias reduction from specialization outweighs the added estimation variance and routing cost. To empirically validate this, we conduct controlled experiments on Dataset A with RISE using SRL as the backbone. We simulate router behavior with accuracy  $p \in \{0.01, \ldots, 1.0\}$ , where the correct expert is chosen with probability p. We systematically vary (i) the upsampling factor  $\alpha$  (Fig. 5) and (ii) the number of experts K' (Fig. 6), averaging over 20 trials. Keeping fixed K' = 3, Fig. 5 shows that higher upsampling reduces error under accurate routing but increases sensitivity to poor routing, consistent







**Figure 6:** MAE vs. router accuracy for varying numbers of experts with fixed  $\alpha = 3$ 

with  $\alpha$  reducing bias while amplifying variance. Keeping  $\alpha=3$  fixed in Fig. 6 shows that larger K' improves accuracy when routing is reliable but offers diminishing returns and greater instability when routing is noisy. In both cases,  $\alpha=1$  or K'=1 reduces RISE to the pooled baseline (or standalone SRL model). Overall, the empirical gain,  $\mathcal{G}_{\mathrm{pooled}}-\mathcal{G}_{\mathrm{RISE}}$ , becomes positive once router accuracy exceeds  $\sim\!60\%$  (with moderate  $\alpha,K'$ ), confirming that RISE successfully operationalizes the theoretical trade-off and remains robust to realistic routing imperfections.

## 6.6 RQ5: ABLATION: RISE vs. HIGH-CAPACITY ENSEMBLES.

A critical question is whether RISE's gains stem from its principled architecture or simply from an increased parameter count. To isolate this, we compare RISE against strong, high-capacity ensembles (Table 8). We train ensembles of 3 and 5 SRL models resulting in significantly additional model size than a RISE-augmented model, where each member is trained on a random data subset to induce diversity. We observe that

**Table 8:** Comparison of RISE vs. traditional ensembles on Dataset A. Best results in **bold** 

		L1 (MAE) ↓				
Experiment	Additional Parameters	All	Many	Median	Few	
SRL SRL+ RISE (K'=3) SRL: 3 ensemble SRL: 5 ensemble	0 +2,100,224 +3,150,336 +5,250,560	7.23 <b>6.57</b> 7.22 7.22	6.64 <b>6.16</b> 6.63 6.62	8.28 <b>7.36</b> 8.28 8.30	9.85 <b>8.30</b> 9.86 9.90	

RISE consistently and significantly outperforms these ensembles, even with their much higher capacity. This highlights a fundamental architectural difference. Standard ensembles create diversity through *unstructured*, random data sampling. In contrast, RISE employs a *principled*, *structured specialization*: it uses validation loss to deterministically identify the model's specific failure modes and trains experts to explicitly target those weaknesses. This confirms that RISE's performance gains are not a product of raw model capacity but are a direct result of its intelligent, data-driven approach to resolving distributional heterogeneity.

## 7 CONCLUSION, BROADER IMPACT, AND LIMITATIONS

We presented RISE (Regression Imbalance handling via Switching Experts), a novel framework that addresses the fundamental challenge of distributional heterogeneity in Deep Imbalanced Regression (DIR). RISE employs a three-stage approach: identifying failure regions via validation loss analysis rather than frequency-based heuristics, training experts with cross-group upsampling to encourage specialization while maintaining smoothness, and learning a gating mechanism, that dynamically routes new inputs to the most appropriate expert at test time. This approach consistently outperforms existing methods, improving both predictive accuracy and fairness, especially for underrepresented regions of the target distribution. RISE is broadly applicable to any regression problem with imbalance issues, advancing the development of more reliable and fair AI systems for critical decision-making.

**Limitations:** RISE introduces additional computational overhead due to training multiple experts and a router network; however, this is partially offset by training experts on last-layer features only. The framework also depends on a high-quality, representative validation set for effective minority subgroup identification and router training. The method's performance and fairness gains can degrade if the validation set is noisy or biased, potentially reinforcing existing biases through expert specialization. Future work could explore adaptive validation strategies and more efficient training schemes to further mitigate these limitations.

## REFERENCES

- Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand kernel learning. In *International Conference on Machine Learning*, pp. 541–549. PMLR, 2018.
- Jiarui Cai, Yizhou Wang, and Jenq-Neng Hwang. Ace: Ally complementary experts for solving long-tailed recognition in one-shot. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 112–121, 2021.
- Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. Semeval-2017 task 1: Semantic textual similarity multilingual and crosslingual focused evaluation. In *Proceedings of the 11th International Workshop on Semantic Evaluation*, 2017.
- James L Cross, Michael A Choma, and John A Onofrey. Bias in medical ai: Implications for clinical decision-making. PLOS Digital Health, 3(11):e0000651, 2024.
- Jiequan Cui, Shu Liu, Zhuotao Tian, Zhisheng Zhong, and Jiaya Jia. Reslt: Residual learning for long-tailed recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3): 3695–3706, 2023. doi: 10.1109/TPAMI.2022.3174892.
- Zijian Dong, Yilei Wu, Chongyao Chen, Yingtian Zou, Yichi Zhang, and Juan Helen Zhou. Improve representation for imbalanced regression through geometric constraints. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2025.
- Yu Gong, Greg Mori, and Fred Tung. Ranksim: Ranking similarity regularization for deep imbalanced regression. In *International Conference on Machine Learning*, pp. 7634–7649. PMLR, 2022.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- Mahsa Keramati, Lili Meng, and R. David Evans. Conr: Contrastive regularizer for deep imbalanced regression. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=RIuevDSK5V.
- Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=Zb6c8A-Fghk.
- Sen Lin, Peizhong Ju, Yingbin Liang, and Ness Shroff. Theory on forgetting and generalization of continual learning. In *International Conference on Machine Learning*, pp. 21078–21100. PMLR, 2023.
- Stylianos Moschoglou, Athanasios Papaioannou, Christos Sagonas, Jiankang Deng, Irene Kotsia, and Stefanos Zafeiriou. Agedb: The first manually collected, in-the-wild age database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshop*, volume 2, pp. 5, 2017.
- Siyuan Mu and Sen Lin. A comprehensive survey of mixture-of-experts: Algorithms, theory, and applications. *arXiv preprint arXiv:2503.07137*, 2025.
- Warwick Nash, Tracy Sellers, Simon Talbot, Andrew Cawthorn, and Wes Ford. Abalone. UCI Machine Learning Repository, 1994. DOI: https://doi.org/10.24432/C55C7W.
- Silvia L Pintea, Yancong Lin, Jouke Dijkstra, and Jan C van Gemert. A step towards understanding why classification helps regression. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 19972–19981, 2023.
- Ruizhi Pu, Gezheng Xu, Ruiyi Fang, Bing-Kun Bao, Charles Ling, and Boyu Wang. Leveraging group classification with descending soft labeling for deep imbalanced regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 19978–19985, 2025.
- Jiawei Ren, Mingyuan Zhang, Cunjun Yu, and Ziwei Liu. Balanced mse for imbalanced visual regression. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

- Rasmus Rothe, Radu Timofte, and Luc Van Gool. Deep expectation of real and apparent age from a single image without facial landmarks. *International Journal of Computer Vision*, 126(2-4): 144–157, 2018.
  - Felix Sattler, Klaus-Robert Müller, and Wojciech Samek. Clustered federated learning: Model-agnostic distributed multitask optimization under privacy constraints. *IEEE transactions on neural networks and learning systems*, 32(8):3710–3722, 2020.
  - Michael Steininger, Konstantin Kobs, Padraig Davidson, Anna Krause, and Andreas Hotho. Density-based weighting for imbalanced regression. *Machine Learning*, 110:2187–2211, 2021.
  - Gilbert W Stewart and Ji-guang Sun. Matrix perturbation theory. (No Title), 1990.
  - Jason G Su, Shadi Aslebagh, Vy Vuong, Eahsan Shahriary, Emma Yakutis, Emma Sage, Rebecca Haile, John Balmes, Michael Jerrett, and Meredith Barrett. Examining air pollution exposure dynamics in disadvantaged communities through high-resolution mapping. *Science Advances*, 10 (32):eadm9986, 2024.
  - Joel A. Tropp. An introduction to matrix concentration inequalities, 2015. URL https://arxiv.org/abs/1501.01571.
  - Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, Cambridge, UK, 2018.
  - Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *EMNLP*, 2018.
  - Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu, and Stella X Yu. Long-tailed recognition by routing diverse distribution-aware experts. *arXiv preprint arXiv:2010.01809*, 2020.
  - Liuyu Xiang, Guiguang Ding, and Jungong Han. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *European Conference on Computer Vision*, pp. 247–263. Springer, 2020.
  - Haipeng Xiong and Angela Yao. Deep imbalanced regression via hierarchical classification adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23721–23730, 2024.
  - Da Xu, Yuting Ye, and Chuanwei Ruan. Understanding the role of importance weighting for deep learning. *arXiv preprint arXiv:2103.15209*, 2021.
  - Yuzhe Yang, Kaiwen Zha, Yingcong Chen, Hao Wang, and Dina Katabi. Delving into deep imbalanced regression. In *International Conference on Machine Learning*, pp. 11842–11851. PMLR, 2021.
  - Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9719–9728, 2020.

## **APPENDIX**

## A Proof of Theorem 1

We restate Theorem 1 from the main paper, together with its assumptions, before presenting a complete proof and a refinement using matrix concentration to obtain a tighter bound. In addition, we extend the analysis to *region-dependent feature covariances*, where feature distributions may differ across regions, to make the theory more realistic. This extension leads to the same qualitative conclusion as in the main paper.

## A.1 ASSUMPTIONS AND NOTATION

We work in the classical fixed-p regime. Let  $p, n \in \text{with } n > p+1$ . The condition n > p+1 ensures that the expectation of the inverse-Wishart distribution exists, which is needed to evaluate the estimation variance. For vectors and matrices we use the Euclidean norm  $\|\cdot\|$  and the spectral norm  $\|\cdot\|_{\text{OD}}$ ;  $(\cdot)$  denotes the trace.

**Assumption 1** (Gaussian design). Fix a positive definite covariance matrix  $\Sigma \in \mathbb{R}^{p \times p}$  with eigenvalues  $0 < \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) < \infty$ . Let  $x_1, \ldots, x_n \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \Sigma)$  be the design rows, stacked into  $X \in \mathbb{R}^{n \times p}$ .

The label space is partitioned into K regions indexed by k = 1, ..., K. Each observation i has a region label  $z_i \in \{1, ..., K\}$ , drawn independently of X, with

$$P(z_i = k) = \rho_k, \qquad \rho_k > 0, \quad \sum_{k=1}^K \rho_k = 1.$$

Let  $n_k = \sum_{i=1}^n \mathbf{1}\{z_i = k\}$  be the (random) region counts, with  $\mathbf{E}[n_k] = n\rho_k$ .

**Remark 1** (On independence of z and x). The assumption  $z_i \perp x_i$  is restrictive but crucial for tractability. In practice (e.g., econometrics, biostatistics), features are often predictive of group membership, in which case off-diagonal terms would appear in conditional covariances and the analysis would require more advanced tools.

**Assumption 2** (Linear region-specific models). For each region k there exists a parameter vector  $w_k^* \in \mathbb{R}^p$ . Observations in region k follow

$$y_i = x_i^{\top} w_{z_i}^* + \varepsilon_i, \qquad \varepsilon_i \stackrel{ind}{\sim} \mathcal{N}(0, \sigma_{z_i}^2),$$

with  $\varepsilon_i$  independent of  $x_i$  and other noise variables. Label vector  $y \in \mathbb{R}^n$ .

Define the population-weighted average parameter

$$w_{\text{avg}} \coloneqq \sum_{k=1}^{K} \rho_k w_k^*,$$

and the centered deviations

$$v_k := w_k^* - w_{\text{avg}}, \qquad \sum_{k=1}^K \rho_k v_k = 0.$$

[Pooled OLS and risk] The pooled ordinary least squares estimator is

$$\widehat{w} = (X^{\top} X)^{-1} X^{\top} y,$$

which is well-defined almost surely for n > p.

We measure performance by the  $\rho$ -weighted mean squared error

$$\mathcal{G}_{\rho}(\widehat{w}) = \sum_{k=1}^{K} \rho_k \|\widehat{w} - w_k^*\|^2.$$

#### A.2 EXACT DECOMPOSITION OF EXPECTED ERROR

**Theorem 1.** Under Assumptions 1-2, for n > p + 1,

$$\mathbf{E}\left[\mathcal{G}_{\rho}(\widehat{w})\right] = \underbrace{\mathbf{E}\left[\|(X^{\top}X)^{-1}X^{\top}\varepsilon\|^{2}\right]}_{estimation\ variance} + \underbrace{\mathbf{E}\left[\|(X^{\top}X)^{-1}X^{\top}\delta\|^{2}\right]}_{mismatch\ term} + \underbrace{\Delta}_{irreducible\ heterogeneity}, \tag{4}$$

where  $\delta \in \mathbb{R}^n$  has entries  $\delta_i = x_i^\top v_{z_i}$ , and

$$\Delta := \sum_{k=1}^{K} \rho_k ||w_k^* - w_{\text{avg}}||^2 = \sum_{k=1}^{K} \rho_k ||v_k||^2.$$

Moreover,

$$\mathbf{E}[\|(X^{\top}X)^{-1}X^{\top}\varepsilon\|^{2}] = \frac{(\Sigma^{-1})}{n-p-1} \left(\sum_{k=1}^{K} \rho_{k} \sigma_{k}^{2}\right).$$
 (5)

*Proof.* Expanding  $\sum_k \rho_k \|\widehat{w} - w_k^*\|^2$  yields

$$\mathcal{G}_{\rho}(\widehat{w}) = \|\widehat{w} - w_{\text{avg}}\|^2 + \Delta.$$

The response can be written as  $y = Xw_{\text{avg}} + \delta + \varepsilon$ , where  $\delta_i = x_i^{\top}v_{z_i}$ . Therefore

$$\widehat{w} - w_{\text{avg}} = (X^{\top} X)^{-1} X^{\top} (\delta + \varepsilon).$$

Squaring gives

$$\|\widehat{w} - w_{\text{avg}}\|^2 = \|(X^{\top}X)^{-1}X^{\top}\varepsilon\|^2 + \|(X^{\top}X)^{-1}X^{\top}\delta\|^2 + 2\langle(X^{\top}X)^{-1}X^{\top}\varepsilon,(X^{\top}X)^{-1}X^{\top}\delta\rangle.$$

Taking expectation: the cross-term vanishes because conditional on (X, z),  $\delta$  is fixed and  $\mathbf{E}[\varepsilon | X, z] = 0$ . This proves (4).

For (5), let  $A = (X^{T}X)^{-1}X^{T}$ . Then

$$\mathbf{E} \|A\varepsilon\|^2 = \mathbf{E} \left( A \, \mathbf{E} [\varepsilon \varepsilon^\top | z] \, A^\top \right)$$
$$= \mathbf{E}_{X,z} \left( (X^\top X)^{-1} X^\top (\sigma_{z_1}^2, \dots, \sigma_{z_n}^2) X (X^\top X)^{-1} \right).$$

Independence of z and X implies

$$\mathbf{E}_z[(\sigma_{z_1}^2, \dots, \sigma_{z_n}^2)] = \left(\sum_{k=1}^K \rho_k \sigma_k^2\right) I_n.$$

Thus

$$\mathbf{E} \|A\varepsilon\|^2 = \left(\sum_k \rho_k \sigma_k^2\right) \mathbf{E}((X^\top X)^{-1}).$$

Since  $X^{\top}X \sim \mathcal{W}_p(\Sigma, n)$ ,

$$\mathbf{E}[(X^{\top}X)^{-1}] = \frac{\Sigma^{-1}}{n-p-1}, \quad n > p+1,$$

hence the trace formula (5).

## A.3 REMARKS

 The decomposition (4) provides a transparent separation of error sources: (i) variance due to noise, (ii) a design-dependent mismatch term induced by parameter heterogeneity, and (iii) the irreducible population heterogeneity Δ.

• The estimation variance admits an exact closed form (5), scaling as O(1/n) for fixed p.

- The mismatch term is always nonnegative. Its precise asymptotics depend on higher-order Wishart moment identities; deriving tight general rates is delicate and left for future work.
- As  $n \to \infty$  with p fixed, the total expected error approaches  $\Delta$ , which is the asymptotic bias from heterogeneity.
- Ill-conditioning of  $\Sigma$  (large  $(\Sigma^{-1})$ ) inflates the variance term and slows convergence to  $\Delta$ .
- These conclusions hold in the fixed-p, large-n regime. In high-dimensional settings with  $p/n \not\to 0$ , ridge regularization and random matrix theory tools are needed.

**Assumption 3** (Sub-Gaussian heterogeneous design). For each region k, the covariates  $x_{k,i}$  are independent mean-zero  $K_{\psi}$ -sub-Gaussian vectors with covariance  $\Sigma_k \succ 0$ , i.e. for every unit vector  $u \in \mathbb{R}^p$  and  $t \in \mathbb{R}$ ,

$$E \exp(t u^{\top} x_{k,i}) \le \exp(K_{\psi}^2 t^2 / 2).$$

Define the mixture covariance

$$\Sigma_{\text{mix}} \coloneqq \sum_{k=1}^{K} \rho_k \Sigma_k,$$

and assume  $\lambda_{\min}(\Sigma_{\min}) > 0$ .

**Proposition 1** (Sample-covariance concentration). Under Assumption 3, there exist constants  $c_0, C_0 > 0$  depending only on  $K_{\psi}$  such that if  $n \geq C_0(p + \log(1/\delta))$  then with probability at least  $1 - \delta$ ,

$$\|\widehat{\Sigma} - \Sigma_{\min}\|_{\text{op}} \le c_0 \|\Sigma_{\min}\|_{\text{op}} \sqrt{\frac{p + \log(1/\delta)}{n}}.$$

Consequently, on this event  $\lambda_{\min}(\widehat{\Sigma}) \geq \frac{1}{2}\lambda_{\min}(\Sigma_{\min})$  and  $\|\widehat{\Sigma}^{-1}\|_{\mathrm{op}} \leq 2/\lambda_{\min}(\Sigma_{\min})$ .

**Theorem 1.1** (Finite-sample generalization under heterogeneous covariances). Suppose Assumptions 2 and 3 hold. Let  $\sigma_{\text{avg}}^2 := \sum_{k=1}^K \rho_k \sigma_k^2$  and define  $\Delta := \sum_{k=1}^K \rho_k \|v_k\|^2$ . There exist constants  $C, C_1, C_2 > 0$  depending only on  $K_{\psi}$  and the spectral condition number  $\kappa(\Sigma_{\text{mix}})$  such that if

$$n \ge C(p + \log(1/\delta)),$$

then with probability at least  $1 - \delta$  the pooled least-squares estimator  $\hat{w} = (X^{\top}X)^{-1}X^{\top}y$  satisfies

$$\mathcal{G}_{\rho}(\widehat{w}) \le C_1 \frac{\sigma_{\text{avg}}^2 p}{n \lambda_{\min}(\Sigma_{\text{mix}})} + C_2 \left\| \Sigma_{\text{mix}}^{-1} \sum_{k=1}^K \rho_k \Sigma_k v_k \right\|^2 + \Delta + \frac{C}{n}.$$
 (6)

*Moreover, in the fixed-p,*  $n \to \infty$  *limit,* 

$$\lim_{n \to \infty} E\left[\mathcal{G}_{\rho}(\widehat{w})\right] = \Delta + \left\| \Sigma_{\min}^{-1} \sum_{k=1}^{K} \rho_k \Sigma_k v_k \right\|^2.$$
 (7)

*Proof (proof sketch and main lemmas).* The proof proceeds in six steps. Below we give the key ideas and cite the concentration results used for brevity and readability.

**Step 1: Decomposition.** Write  $y = Xw_{\text{avg}} + \delta + \varepsilon$  where  $\delta_i = x_i^{\top}v_{z_i}$  and  $\varepsilon = (\varepsilon_i)_{i=1}^n$ . Then

$$\widehat{w} - w_{\text{avg}} = \widehat{\Sigma}^{-1} \left( \frac{1}{n} X^{\top} \delta \right) + \widehat{\Sigma}^{-1} \left( \frac{1}{n} X^{\top} \varepsilon \right).$$

Thus

$$\mathcal{G}_{\rho}(\widehat{w}) = \left\| \widehat{\Sigma}^{-1} \frac{1}{n} X^{\top} \varepsilon \right\|^{2} + \left\| \widehat{\Sigma}^{-1} \frac{1}{n} X^{\top} \delta \right\|^{2} + 2 \langle \cdot, \cdot \rangle + \Delta.$$

The three display terms correspond to estimation variance, mismatch, and a cross-term.

Step 2: Control of  $\widehat{\Sigma}$ . Proposition 1 (matrix concentration for sub-Gaussian samples; see Vershynin (2018); Tropp (2015)) implies that for  $n \gtrsim p + \log(1/\delta)$  the event in which  $\|\widehat{\Sigma} - \Sigma_{\text{mix}}\|_{\text{op}}$  is small holds with probability  $1 - \delta$ . On this event one obtains the deterministic bound  $\|\widehat{\Sigma}^{-1}\|_{\text{op}} \lesssim 1/\lambda_{\min}(\Sigma_{\min})$ .

Step 3: Estimation variance term. Conditioning on X and z,  $X^{\top}\varepsilon$  is a mean-zero vector with componentwise variances  $\sigma_{z_i}^2 \|x_i\|^2$ . Standard conditional-sub-Gaussian tail bounds together with operator-norm control of  $\widehat{\Sigma}^{-1}$  yield the displayed O(p/n) bound in (6). One may make this fully explicit by combining Hanson–Wright and matrix Bernstein inequalities (see Vershynin (2018); Tropp (2015)).

## Step 4: Population limit and asymptotic bias. Note

$$\frac{1}{n}X^{\top}y = \frac{1}{n}\sum_{i=1}^{n} x_i x_i^{\top} w_{z_i}^* + \frac{1}{n}X^{\top}\varepsilon.$$

By the law of large numbers and multinomial concentration of region counts,  $\frac{1}{n} \sum_i x_i x_i^\top w_{z_i}^* \to \sum_k \rho_k \Sigma_k w_k^*$  and  $\widehat{\Sigma} \to \Sigma_{\text{mix}}$ . Hence  $\widehat{w} \to w_\infty$  where  $w_\infty = \Sigma_{\text{mix}}^{-1} \sum_k \rho_k \Sigma_k w_k^*$ . Using  $v_k = w_k^* - w_{\text{avg}}$  yields the asymptotic mismatch bias in (7).

Step 5: Finite-sample mismatch fluctuation. The deviation  $\frac{1}{n}X^{\top}\delta - \sum_k \rho_k \Sigma_k v_k$  is a mean-zero sum of sub-Gaussian terms and therefore has Euclidean norm  $O_p(1/\sqrt{n})$ . Multiplication by  $\widehat{\Sigma}^{-1}$ , which is O(1) in operator norm on the concentration event, yields an  $O_p(1/\sqrt{n})$  deviation of the centered estimator; squaring gives the  $O_p(1/n)$  remainder in (6).

**Step 6: Cross-term.** The cross-term is bounded in absolute value via Cauchy–Schwarz and is of smaller order (absorbed into the displayed C/n remainder) under the same sample-size regime.

Combining the bounds in Steps 3–6 yields (6) and the limit (7).

## Remark 2. (References)

• The matrix-concentration proposition can be proved by applying matrix Bernstein / non-commutative Bernstein inequalities as in Tropp (2015) or via Vershynin's sub-Gaussian covariance concentration (see Vershynin (2018)).

 All big-O and constants can be made explicit by tracking constants in Hanson-Wright and matrix Bernstein inequalities; we omitted explicit numerical constants for readability.

**Remark 3** (Interpretation). Unlike the homogeneous-covariance case, pooled OLS error converges not only to the irreducible heterogeneity  $\Delta$  but also to a persistent asymptotic mismatch bias (cf. Eq. (7)). This bias vanishes only under special conditions such as  $\Sigma_k \equiv \Sigma$  for all k or  $\sum_k \rho_k \Sigma_k v_k = 0$ . Finite-sample fluctuations of the mismatch term decay at rate O(1/n), while the estimation variance scales as O(p/n). Both contributions are magnified when  $\Sigma_{\text{mix}}$  is ill-conditioned.

## B THEORETICAL GUARANTEES FOR RISE WITH UPSAMPLING AND ROUTING

We present a rigorous finite-sample analysis of RISE. We first state assumptions, then supporting lemmas, and finally the main theorem with proof. We also derive the exact pooled decomposition, and conclude with a corollary giving explicit sufficient conditions under which RISE improves over pooled OLS.

#### **B.1** Assumptions

**Assumption 4** (Sub-Gaussian design and bounded covariance). For each region  $k \in [K]$ , covariates  $x_{k,i} \in \mathbb{R}^p$  are i.i.d. mean-zero  $K_{\psi}$ -sub-Gaussian vectors with covariance  $\Sigma_k = \mathbf{E}[x_{k,i}x_{k,i}^{\top}] \succ 0$ . Eigenvalues are uniformly bounded:

$$0 < \underline{\lambda} \le \lambda_{\min}(\Sigma_k) \le \lambda_{\max}(\Sigma_k) \le \overline{\lambda} < \infty.$$

**Assumption 5** (Noise tails). For each region k, labels satisfy  $y = x^\top w_k^* + \varepsilon$  with  $\mathbf{E}[\varepsilon \mid x] = 0$ ,  $(\varepsilon \mid x) = \sigma_k^2$ , and  $\sigma_k^2 \leq \sigma_{\max}^2 < \infty$ . Moreover, the noise satisfies a uniform tail condition: either (i)  $\varepsilon$  is sub-Gaussian, or (ii)  $\varepsilon^2$  is sub-exponential (uniform constants). These tail assumptions are used to obtain operator-norm concentration for heteroskedastic noise matrices; if only finite variance is available, replace sample moments by robust estimators (truncation / median-of-means).

**Assumption 6** (Routing). Each population sample is drawn from region k with probability  $\rho_k$ . Conditional on region k, the sample is routed to expert j with fixed probability  $q_k(j)$ , independent of features x and noise  $\varepsilon$ . Each expert j may upweight its own region by a factor  $\alpha_j \geq 1$  (we explain below how this affects the training mixture and the realized counts). All routing probabilities  $\{q_k(j)\}$  are fixed (non-adaptive).

## B.2 EFFECTIVE DISTRIBUTIONS AND A CLARIFYING REMARK ON UPSAMPLING

We use two distinct population-level quantities; reviewers should not conflate them.

(i) Marginal routing probability (controls realized counts). The marginal probability that a random population sample is routed to expert j (before any upsampling normalization) is

$$p_j^{\text{route}} \coloneqq \sum_{k=1}^K \rho_k \, q_k(j).$$

The realized number  $N_j$  of training samples routed to expert j is multinomial/binomial with mean  $np_j^{\text{route}}$ . Lemma 1 below gives precise concentration for  $N_j$ .

(ii) Unnormalized upweight mass and training mixture (controls bias). To describe how upsampling changes the *training mixture* used to estimate each expert, define unnormalized weights

$$\omega_{k \to j} \coloneqq \begin{cases} \alpha_j \rho_j q_j(j), & k = j, \\ \rho_k q_k(j), & k \neq j, \end{cases} \qquad \Omega_j \coloneqq \sum_{k=1}^K \omega_{k \to j}, \qquad \pi_{k \to j} \coloneqq \frac{\omega_{k \to j}}{\Omega_j}.$$

Here  $\pi_{k \to j}$  defines the *population-level training mixture* for expert j: when estimating expert j we (conceptually) mix regions k with proportions  $\pi_{k \to j}$ . These  $\pi_{k \to j}$  enter the deterministic bias  $B_{\text{det}}$  via

$$\Sigma_{j,\text{train}} \coloneqq \sum_{k=1}^{K} \pi_{k \to j} \Sigma_{k}, \qquad w_{j}^{\text{eff}} \coloneqq \Sigma_{j,\text{train}}^{-1} \Big( \sum_{k=1}^{K} \pi_{k \to j} \Sigma_{k} w_{k}^{*} \Big).$$

**Remark:**  $\omega_{k \to j}$  (and hence  $\pi_{k \to j}$ ) involve  $\alpha_j$  and  $\rho_k$  and are *not* probabilities over experts; they describe the training mixture used to form population-level bias terms. The realized counts  $N_j$  (used for variance bounds) are governed by  $p_j^{\text{route}}$ , which depends only on  $\rho_k, q_k(j)$  and not on  $\alpha_j$ . In practice, upsampling can be implemented either by (A) re-sampling from the modified mixture induced by  $\pi_{k \to j}$  (sampling interpretation), or (B) by attaching per-sample weights in the loss (weighting interpretation). The analysis below treats the bias via  $\pi_{k \to j}$  and controls variance via the realized counts  $N_j$ ; if you implement upsampling by weighting, replace  $N_j$  in variance rates by the appropriate ESS (effective sample size) — see Practical Considerations.

Define the population-level weighted noise and effective-sample-size

$$\sigma_{j,\text{eff}}^2 = \sum_{k=1}^K \pi_{k \to j} \sigma_k^2, \qquad n_{\text{eff}}^{(j)} = n \cdot \Omega_j.$$

#### B.3 PRELIMINARY LEMMAS

**Lemma 1** (Routing counts concentration). Let  $p_j^{\text{route}} = \sum_{k=1}^K \rho_k q_k(j)$ . Then  $(N_1, \ldots, N_J) \sim \text{Multinomial}(n; p_1^{\text{route}}, \ldots, p_J^{\text{route}})$ . Fix  $\delta \in (0,1)$ . There exist constants  $c_1, c_2 > 0$  such that for each j and any t > 0,

$$\Pr\left(|N_j - np_j^{\text{route}}| \ge t\right) \le 2\exp\left(-\frac{t^2}{2np_j^{\text{route}} + (2/3)t}\right).$$

Choosing  $t_j = c_1 \sqrt{np_j^{\text{route}} \log(J/\delta)} + c_2 \log(J/\delta)$  and applying a union bound yields that with probability at least  $1 - \delta$ ,

$$|N_j - np_j^{\text{route}}| \le t_j \quad \text{for all } j \in [J].$$

Consequently, if  $np_j^{\text{route}} \gtrsim C(p + \log(J/\delta))$  for all j, then with probability at least  $1 - \delta$  we have  $N_j \geq \frac{1}{2} np_j^{\text{route}}$  for every j.

**Lemma 2** (Design and noise concentration). Assume rows of  $X_j$  are independent  $K_{\psi}$ -sub-Gaussian vectors with covariance  $\Sigma_{j,\text{train}}$ , and assume the noise satisfies the tail condition in Assumption 5 (sub-Gaussian or sub-exponential so that  $\varepsilon_i^2 x_i x_i^{\top}$  has controlled sub-exponential operator-norm). Fix  $\delta \in (0,1)$ . There exist constants  $C_0, C_1, C_2 > 0$  (depending on  $K_{\psi}$  and the noise-tail constants) such that, provided  $N_j \gtrsim p + \log(J/\delta)$  for all j, the following holds with probability at least  $1 - \delta$  simultaneously over  $j \in [J]$ :

$$\left\| \frac{1}{N_j} X_j^{\top} X_j - \Sigma_{j, \text{train}} \right\| \le C_0 \left( \sqrt{\frac{p + \log(J/\delta)}{N_j}} + \frac{p + \log(J/\delta)}{N_j} \right), \tag{8}$$

$$\left\| \frac{1}{N_j} \sum_{i \in \text{train}_j} \varepsilon_i^2 x_i x_i^\top - \sigma_{j, \text{eff}}^2 \Sigma_{j, \text{train}} \right\| \le C_1 \sigma_{\text{max}}^2 \sqrt{\frac{p + \log(J/\delta)}{N_j}} + C_2 \sigma_{\text{max}}^2 \frac{p + \log(J/\delta)}{N_j}. \tag{9}$$

In the usual regime  $N_j \gtrsim p + \log(J/\delta)$  the square-root term dominates and the simpler form with only the  $\sqrt{\cdot}$  term is valid.

**Remarks on the lemmas.** - Lemma 1 is a standard Bernstein/Hoeffding tail for binomial/multinomial counts. - Lemma 2 follows from applying matrix Bernstein / Vershynin concentration to sub-Gaussian rows, and to the heteroskedastic weighted noise matrices  $\varepsilon_i^2 x_i x_i^{\mathsf{T}}$  using the noise-tail assumption. If the noise has only finite variance, replace empirical moments by robust estimators (truncation, MOM) to retain high-probability control.

#### B.4 MAIN THEOREM FOR RISE

Intuition. The decomposition below separates prediction risk into: irreducible noise  $\sigma_{\text{avg}}^2$ ; deterministic bias  $B_{\text{det}}$  due to training-mixture mismatch; estimation variance  $V_{\text{est}}$  governed by realized counts  $N_j$ ; and a cross-term  $R_{\text{cross}}$  of smaller order.

**Theorem 2** (Generalization error of RISE). Suppose Assumptions 4–6 and the noise-tail condition in Assumption 5 hold. Suppose further that the marginal routing masses satisfy  $np_j^{\text{route}} \gtrsim C(p + \log(J/\delta))$  for all j (so Lemma 1 implies  $N_j \gtrsim p$  w.h.p.). Then, conditioning on the joint high-probability event from Lemmas 1–2, with probability at least  $1 - \delta$ ,

$$\mathcal{G}_{\text{RISE}}(\alpha, q) = \sigma_{\text{avg}}^2 + B_{\text{det}}(\alpha, q) + V_{\text{est}}(\alpha, q) + R_{\text{cross}}(\alpha, q),$$

$$B_{\text{det}}(\alpha, q) = \sum_{k=1}^K \rho_k \sum_{j=1}^J q_k(j) \|w_j^{\text{eff}} - w_k^*\|_{\Sigma_k}^2,$$

$$(10)$$

$$V_{\text{est}}(\alpha, q) \leq C_1 \sum_{k=1}^{K} \sum_{j=1}^{J} \rho_k q_k(j) \frac{\sigma_{j, \text{eff}}^2}{N_j} \left( \Sigma_k \Sigma_{j, \text{train}}^{-1} \right)$$
$$+ C_1' \sum_{k, j} \rho_k q_k(j) \frac{\sigma_{\text{max}}^2 p}{N_j} \sqrt{\frac{p + \log(J/\delta)}{N_j}},$$

$$|R_{\text{cross}}(\alpha, q)| \le C_2 \left( \max_{j, k} \|w_j^{\text{eff}} - w_k^*\|_{\Sigma_k} \right) \sqrt{\lambda_{\text{max}}(\Sigma_k \Sigma_{j, \text{train}}^{-1})} \sqrt{\frac{p + \log(J/\delta)}{N_{\text{min}}}}.$$

Here  $N_{\min} = \min_j N_j$ , and constants  $C_1, C'_1, C_2$  depend only on  $K_{\psi}$  and the noise-tail parameters.

*Proof sketch.* All concentration statements below are applied on the joint high-probability event from Lemmas 1 and 2.

**Step 1 (decomposition).** For a test point  $(x, y) \sim \mathcal{R}_k$  routed to expert j,

$$\mathbf{E}[(x^{\top}\widehat{w}_{j} - y)^{2} \mid x] = ||w_{j}^{\text{eff}} - w_{k}^{*}||_{\Sigma_{k}}^{2} + (\Sigma_{k}(\widehat{w}_{j})) + (w_{j}^{\text{eff}} - w_{k}^{*})^{\top} \Sigma_{k}(\widehat{w}_{j} - w_{j}^{\text{eff}}) + \sigma_{k}^{2}.$$

Averaging over (k, j) with weights  $\rho_k q_k(j)$  yields (10) and the definition of  $\sigma_{\text{avg}}^2$ .

**Step 2 (bias).** The first term is exactly  $B_{\text{det}}$ .

Step 3 (variance). By Lemma 2 the sandwich covariance satisfies

$$(\widehat{w}_j) = \frac{\sigma_{j,\text{eff}}^2}{N_j} \Sigma_{j,\text{train}}^{-1} + E_j, \qquad \|E_j\| \le C \frac{\sigma_{\text{max}}^2}{N_j} \sqrt{\frac{p + \log(J/\delta)}{N_j}}.$$

Taking trace against  $\Sigma_k$  and averaging with  $\rho_k q_k(j)$  yields the bound on  $V_{\rm est}$ .

Step 4 (cross-term). By Cauchy-Schwarz,

$$|R_{k,j}| \le ||w_j^{\text{eff}} - w_k^*||_{\Sigma_k} ||\widehat{w}_j - w_j^{\text{eff}}||_{\Sigma_k}.$$

Using operator-norm change of metric and the concentration bound for  $\|\widehat{w}_j - w_j^{\text{eff}}\|_{\Sigma_{j,\text{train}}}$  (of order  $\sqrt{(p + \log)/N_j}$ ) gives the stated bound on  $R_{\text{cross}}$ .

## B.5 POOLED MODEL AND COMPARISON

For the pooled estimator  $\widehat{w}_{pool} = (X^{\top}X)^{-1}X^{\top}y$ , the same decomposition (conditioning on the same high-probability event) yields

$$\mathcal{G}_{\text{pooled}} = \sigma_{\text{avg}}^2 + B_{\text{pooled}} + V_{\text{pooled}},$$

where

$$B_{\text{pooled}} = \sum_{k=1}^{K} \rho_k \|w_{\text{pool}}^{\text{eff}} - w_k^*\|_{\Sigma_k}^2, \qquad w_{\text{pool}}^{\text{eff}} = \left(\sum_k \rho_k \Sigma_k\right)^{-1} \left(\sum_k \rho_k \Sigma_k w_k^*\right),$$

and  $V_{\rm pooled}$  is the pooled estimation variance (bounded by O(p/n) under our assumptions). Subtracting gives the exact comparison

$$\mathcal{G}_{\text{RISE}} - \mathcal{G}_{\text{pooled}} = (B_{\text{det}} - B_{\text{pooled}}) + (V_{\text{est}} - V_{\text{pooled}}) + R_{\text{cross}},$$

since the common  $\sigma_{\text{avg}}^2$  cancels.

#### B.6 ILLUSTRATIVE COROLLARY: SUFFICIENT CONDITIONS FOR IMPROVEMENT

**Corollary 1** (When RISE improves pooled). *Under the conditions of Theorem 2, suppose further that* 

- (i) (Bias reduction)  $B_{\mathrm{pooled}} B_{\mathrm{det}} \ge c_0 \sum_k \rho_k \|w_k^* w_{\mathrm{avg}}\|_{\Sigma_k}^2$  for some  $c_0 > 0$ ;
- (ii) (Sufficient counts)  $\min_j N_j \gtrsim C(p + \log(J/\delta))$  so that the variance and cross-term remainders are small.

*Then with probability at least*  $1 - \delta$ *,* 

$$\mathcal{G}_{\mathrm{RISE}} < \mathcal{G}_{\mathrm{pooled}}$$
.

*Proof sketch.* Under (ii) the variance and cross-term penalties scale as  $O(p/N_j)$  and  $O(\sqrt{p/N_j})$  and can be made small; under (i) the deterministic bias reduction is order  $\Delta_{\text{glob}}$ . Hence the total difference is negative with high probability.

## PRACTICAL CONSIDERATIONS AND LIMITATIONS

The quantities appearing in Theorem 2 (such as  $w_k^*$ ,  $\Sigma_k$ ,  $\sigma_k^2$ , and the induced effective parameters  $w_j^{\text{eff}}$ ) are population-level objects and unknown in practice. In experiments we approximate them with plug-in estimates from held-out validation data; standard perturbation bounds for covariance estimation (Stewart & Sun, 1990; Vershynin, 2018) imply that population inequalities carry over to plug-in versions with sufficient validation sample size (scaling as  $O(p/\gamma^2)$  for margin  $\gamma$ ).

Important limitations and practical conditions:

- Routing independence assumption. We assume  $q_k(j)$  are fixed and independent of x. If routing depends on features (learned gating that uses x), conditional covariances and bias expressions change; the analysis must be adapted to conditional mixtures.
- Implementation of upsampling. Our statements separate the population-level training-mixture  $\pi_{k \to j}$  (used to define deterministic bias) from the realized counts  $N_j$  (used for variance). In practice upsampling can be implemented either by (A) re-sampling from a modified mixture (sampling) or (B) by attaching weights in the loss (weighting). If weighting is used replace all  $N_j$ -based rates by the appropriate effective sample size (ESS) and analyze weighted-OLS (sandwich) covariance (we provide that variant in the appendix on request).
- **Noise tails / robustness.** We assume sub-Gaussian or sub-exponential noise. If only finite variance is available, robust estimators (truncation or median-of-means) are required to obtain comparable high-probability bounds.
- Minimum routing mass required. The bounds require non-negligible routing mass for each expert:  $np_j^{\text{route}} \gtrsim C(p + \log(J/\delta))$ . If some expert is assigned vanishing mass, concentration and OLS asymptotics break down and regularization or enforced minimum routing mass is necessary.

## C DATASET DETAILS

We evaluate our RISE framework on the benchmark datasets on four diverse regression datasets: two datasets from the computer vision domain (Dataset A (Moschoglou et al. (2017)) and Dataset B (Rothe et al. (2018)), one from the natural language processing domain (STS-B Cer et al. (2017)) and one standard tabular regression dataset- UCI Abalone Nash et al. (1994).

- Dataset A (Moschoglou et al. (2017)): An image regression dataset with 12,208 training samples, 2,140 validation samples, and 2,140 test samples. The target range spans from 0 to 101.
- Dataset B (Rothe et al. (2018)): A large-scale image regression dataset containing 191,509 training samples, 11,022 validation samples, and 11,022 test samples. The target range spans from 0 to 186.
- STS-B: A text similarity dataset containing 5,249 training sentence pairs, 1,000 validation pairs, and 1,000 test pairs, with similarity scores ranging from 0 to 5.
- **UCI Abalone**: A standard tabular benchmark predicting shellfish ring from 9 different physical measurements, the dataset consists of of 3155 training, 511 test and 511 validation samples with the target column shellfish ring ranging from 1 to 29.

We follow the train/val/test split provided in Yang et al. (2021)

## D IMPLEMENTATION DETAILS

#### D.1 NETWORK ARCHITECTURE

Figure 3 illustrates the RISE architecture and its key components. Let the full dataset be denoted by  $D=D_{\text{train}}\cup D_{\text{val}}\cup D_{\text{test}}$ . The RISE framework begins by employing a baseline Deep Imbalanced Regression (DIR) model  $f_{\theta}$  for both feature extraction and minority subgroup identification. Input data—whether image, text, or tabular—is first passed through the feature extractor  $h_{\theta}$ , a component of the baseline model  $f_{\theta}$ . This model is pre-trained on  $D_{\text{train}}$  using existing DIR methods such as LDS-FDS (Yang et al. (2021)), RankSim (Gong et al. (2022)), and SRL (Dong et al. (2025)). The architecture of the baseline can be expressed as  $f_{\theta}(x)=E_1(h_{\theta}(x))$ , where  $h_{\theta}(x)$  denotes the backbone feature extractor, typically instantiated as ResNet-50 for images and BiLSTM for text. RISE is agnostic to the specific DIR method and can integrate any baseline model  $f_{\theta}$  built on these backbone architectures.

**RISE-Identify:** To address underperformance in imbalanced regression, we propose RISE-Identify for identifying minority or poorly modeled regions by analyzing the joint distribution of validation

loss and target labels. Specifically, we fit a Gaussian Mixture Model (GMM) to validation data to uncover latent structure in model error patterns, enabling targeted expert specialization.

In regression tasks with heterogeneous label distributions, performance typically degrades in minority subregions of the label space. A key observation is that these regions often exhibit higher and more variable validation losses. By analyzing the joint distribution of validation loss and target values, we can detect structured error patterns that are not captured by traditional frequency-based binning.

Following (Yang et al. (2021)), we partition the continuous label space into disjoint intervals  $B_i$  and compute the average loss in each bin:

$$\ell_i = \frac{1}{|B_i|} \sum_{j \in B_i} \mathcal{L}(f_\theta(x_j), y_j) \tag{11}$$

Here,  $B_i$  is the set of samples whose continuous labels fall within the boundaries of bin i,  $\mathcal{L}$  is typically Mean Squared Error (MSE) or Mean Absolute Error (MAE),  $f_{\theta}$  denotes the baseline model, and  $|B_i|$  is the number of samples in bin i. Importantly, the model is trained and evaluated end-to-end in continuous space—binning is used only for region-level loss estimation, not for converting regression into classification.

Next, we fit a K'-component Gaussian Mixture Model (GMM) over the joint distribution of loss-label pairs:

$$p(\ell, y) = \sum_{j=1}^{K'} \pi'_j \mathcal{N}((\ell, y) | \mu_j, \Sigma_j)$$
(12)

where  $\mu_j$  and  $\Sigma_j$  denote the mean vector and covariance matrix of the j-th component, respectively. The component with the lowest mean loss (along the loss dimension of  $\mu_j$ ) is treated as the majority group, while the remaining components define minority subgroups requiring dedicated experts.

Unlike frequency-based approaches that often result in non-contiguous minority regions, our loss-label distribution analysis produces continuous minority regions, aligning with the principle of region similarity and enabling more homogeneous expert training. We observe a memorization effect where the baseline model achieves the lowest training loss in few-shot regions despite higher test errors. To address this, we use held-out set loss as a more reliable signal for minority subgroup identification, as it better reflects true generalization behavior and mitigates misleading effects of memorization.

Unlike methods based on label frequency or manual binning, our loss-aware formulation is adaptive and reflects the true generalization profile of the baseline model. The identified regions are continuous, semantically meaningful, and sensitive to the model's inductive biases. By relying on the validation—training loss gap, our method is capable of detecting overfitting and memorization—particularly in underrepresented areas. The resulting expert assignments are thus aligned with true generalization performance, enabling smooth transitions between expert domains. This leads to coherent regional specialization and improved overall generalization, especially in long-tailed or imbalanced regression settings.

The RISE-Identify component leverages a held-out validation set  $(80\% \text{ of } D_{val})$  to conduct this loss-label distribution analysis, with cross-validation on the remaining 20% to determine GMM hyperparameters like the number of components K'. As illustrated in Fig.4, this approach successfully identifies continuous minority regions requiring specialized experts - one towards the lower end of the label distribution and another in the higher range.

**RISE-Train:** RISE-Train trains K'-1 additional expert networks  $E_2, E_3, \ldots, E_{K'}$  for the identified minority regions, while the baseline model  $E_1$  (extracted from  $f_{\theta}$ ) serves as the expert for the majority region. Each expert  $E_j$  operates on shared features produced by the frozen backbone  $h_{\theta}$ , and produces predictions as:

$$\hat{y}_i = E_i(h_\theta(x)) \tag{13}$$

To address data imbalance, we adopt a *Cross-Group Training with Upsampling* strategy. This approach (T2) is particularly effective for regression tasks where adjacent labels exhibit strong

#### 1080 **Algorithm 1** RISE Training **Require:** Dataset $D = \{D_{train}, D_{val}\}$ , model $f_{\theta} = \{h_{\theta}, E_1\}$ , experts K', upsampling $\alpha$ 1082 **Ensure:** Experts $\{E_j\}_{j=1..K'}$ , router R1: // Phase 1: RISE-Identify 1084 2: $F \leftarrow \emptyset$ 3: **for** (x, y) in $D_{val}$ **do** $\hat{y} \leftarrow E_1(h_{\theta}(x))$ {Baseline model prediction} 1087 $\ell \leftarrow \mathcal{L}(f_{\theta}(x), y)$ {Compute validation loss per Eq. 11} 1088 $F \leftarrow F \cup \{(\ell, y)\}$ 1089 8: $gmm \leftarrow \text{FitGaussianMixture}(F, K')$ {Fit GMM using Eq. 12} 1090 $\{R'_i\}_{j=1..K'} \leftarrow \text{GetMinorityRegions}(gmm) \{\text{Identify expert regions}\}$ 1091 10: // Phase 2: RISE-Train 1092 11: Initialize experts $E_2$ through $E_{K'}$ 12: **for** i = 2 to K' **do** 1094 13: for epoch = 1 to T do 1095 14: for $(X_b, Y_b)$ in $D_{train}$ do 15: $F \leftarrow h_{\theta}(X_b)$ {Extract shared features} 16: for j=1 to $|X_b|$ do if $y_j \in R'_i$ then 17: 1099 18: $w_j \leftarrow \alpha$ {Upsample minority region samples} 19: 1100 20: $w_i \leftarrow 1$ {Normal weight for other samples} 1101 end if 21: 1102 22: 1103 $Y \leftarrow E_i(F)$ {Get predictions from Eq. 13} 23: 1104 $L \leftarrow \frac{1}{|X_b|} \sum_{j=1}^{|X_b|} w_j (\hat{Y}_j - Y_j)^2$ {Weighted loss from Eq. 14} Update $E_i$ using gradient $\nabla L$ 24: 1105 25: 1106 26: end for 1107 27: end for 1108 28: end for 1109 29: Initialize router R1110 30: **for** epoch = 1 to T' **do** 1111 31: for $(X_b, Y_b)$ in $D_{val}$ do 32: $F \leftarrow h_{\theta}(X_b)$ 1113 33: for j=1 to $|X_b|$ do 1114 $t_i \leftarrow \text{find } i \text{ such that } y_i \in R'_i \{ \text{Assign ground truth expert labels} \}$ 34: 1115 35: $r \leftarrow R(F)$ {Get router probabilities} 1116 36: $\mathcal{L}_{router} \leftarrow \text{CrossEntropy}(r, T_b) \text{ using Eq. 17}$ 1117 Update R using gradient $\nabla \mathcal{L}_{router}$ 38: 1118 39: end for 1119 40: **end for** 1120 41: **return** $\{E_j\}_{j=1..K'}, R$ 1121

correlations, enabling smooth transitions between expert domains while preserving specialization, as confirmed by our empirical analysis. For each identified region  $R'_j$ , we upsample the samples in  $R'_j$  by assigning a higher weight  $\alpha>1$ , while keeping the sample weights unchanged elsewhere. We train each expert using  $D_{train}$  where loss for each expert  $E_j$  is given by:

$$\mathcal{L}_{\text{expert}}^j = \frac{1}{N} \sum_{i=1}^N w_i (y_i - \hat{y}_i)^2$$
 (14)

with sample weights  $w_i$  defined as:

112211231124

11251126

1128 1129

1130

113111321133

## Algorithm 2 RISE-Inference

**Require:** Sample x, backbone  $h_{\theta}$ , router R, experts  $\{E_j\}_{j=1..K'}$ 

Ensure: Prediction  $\hat{y}$ 

1:  $F \leftarrow h_{\theta}(x)$  {Extract features using frozen backbone}

2:  $r \leftarrow R(F)$  {Get router probabilities}

3:  $j^* \leftarrow$  Select expert using Eq. 16

4:  $\hat{y} \leftarrow E_{j^*}(F)$  {Get final prediction using Eq. 18}

5: return *i* 

 $w_i = \begin{cases} \alpha & \text{if } x_i \in R_j' \\ 1 & \text{otherwise} \end{cases}$  (15)

Here,  $\alpha$  is an upsampling hyperparameter that emphasizes minority-region samples, and N is the total number of samples in  $D_{train}$ . Importantly, only the final layer of each new expert  $E_j$  (for j=2,...,K') is trained, while the shared backbone  $h_{\theta}$  and the baseline expert  $E_1$  remain frozen. This facilitates efficient parameter sharing and reduces computational overhead.

**RISE-Inference:** We train a router network (implementing the gating network  $g_{\phi}$  from Eq. 1) using a held-out validation set (80% of  $D_{\text{val}}$ ) to perform dynamic expert selection, with the remaining 20% used for hyperparameter validation. We motivate the choice of using held-out data in Sec. 6.4. Unlike soft routing strategies that blend predictions from multiple experts, we adopt a hard routing approach, where exactly one expert is selected per input. This decision is motivated by Theorem 1, which demonstrates that mixing predictions from heterogeneous regions can lead to interference and degraded performance due to distributional mismatch.

The router is trained as a classification task to predict which expert should handle each input. For each validation sample (x,y), we first determine the ground truth expert assignment by checking which region  $R'_j$  the label y belongs to. The router then learns to map input features to these expert assignments.

Given input x, the router processes shared features  $h_{\theta}(x)$  and outputs mixing coefficients  $\pi_k(x)$  over the K' experts, implementing the gating mechanism from Eq. (1). A hard assignment is then made as follows:

$$j^* = \arg\max_{j \in \{1, \dots, K'\}} g_{\phi}(h_{\theta}(x))_j$$
 (16)

where  $j^*$  denotes the index of the selected expert, consistent with the final prediction  $\hat{y} = E_{j^*}(x)$  described in Section 4.3. The router is trained using an inverse-frequency weighted cross-entropy loss to mitigate expert imbalance:

$$\mathcal{L}_{\text{router}} = -\sum_{j=1}^{K'} w_j t_j \log(p_j)$$
 (17)

Here,  $p_j$  is the predicted probability for expert j,  $t_j$  is the ground truth expert label from the RISE-Identify stage, and  $w_j = \frac{1}{f_j}$  is the inverse frequency of expert j's assigned region, where  $f_j$  is the fraction of samples assigned to expert j in  $D_{\text{val}}$ .

At inference time, the router selects a single expert  $E_{j^*}$  based on the hard assignment, and the final prediction is:

$$\hat{y} = E_{i^*}(h_{\theta}(x)) \tag{18}$$

This hard routing strategy offers several advantages: it prevents distribution mixing that could degrade expert specialization, reduces computation by evaluating only one expert at inference, provides

interpretable routing decisions, and maintains clear accountability for predictions. The complete RISE framework is summarized in Algorithm 1 for training and Algorithm 2 for inference.

#### D.2 TRAINING DETAILS

Experiments were run on an AWS ml.g6.24xlarge instance equipped with 4 NVIDIA GPUs. For all baseline DIR models, we use official released model weights or reproduce their best configuration using the official implementations. For the model architecture, we froze the backbone network (ResNet-50 for images, pretrained on ImageNet; BiLSTM with GloVe embeddings for text) and implemented expert networks with two fully connected layers (dimensions: 2048,512,1) with ReLU activation and dropout (0.2) for ResNet-50. The router network consists of three linear layers with ReLU activation and a final softmax layer. Expert training was conducted for 50 epochs using the Adam optimizer with a learning rate of 3e-5, utilizing a batch size of 64. For image datasets, we applied standard augmentations including random horizontal flips, crops, rotations, affine transformations, and color jittering, followed by normalization. Text data was processed using SpaCy tokenization with a maximum sequence length of 40.

Hyperparameters were tuned through grid search, exploring different numbers of experts  $(K' \in 2, 5]$ ), upsampling ratios (Upsample  $(\alpha) \in [1, 5]$ ) based primarily on validation's overall MAE. For Dataset A (Moschoglou et al. (2017)) and Dataset B (Rothe et al. (2018)) datasets, we set K' = 3 experts, with one expert assigned to the left tail, one to the right tail, and one for the majority region. The upsampling ratio was set to 3 for Dataset A (Moschoglou et al. (2017)) and 2 for Dataset B (Rothe et al. (2018)). For the STS dataset using the RankSim baseline, we used K' = 2 experts, identifying a one-sided under-performing region with an upsampling ratio of 3, while K' = 3 experts with upsampling ratio of 3 were chosen for LDS+FDS and SRL baselines. The number of experts (K') and their assignments were determined based on the baseline model's loss-label distribution and can vary depending on model performance. This approach ensures we only train additional experts for regions where the baseline model underperforms. Further, identified minority regions for experts may differ across baseline models due to variations in their learned representations and performance characteristics.

## E ADDITIONAL EXPERIMENTAL RESULTS

## E.1 RISE PERFORMANCE ON ADDITIONAL DATASETS

To further demonstrate the effectiveness of RISE, we evaluate our method on additional datasets beyond Dataset A (Moschoglou et al. (2017)). Table 9 presents results on Dataset B Rothe et al. (2018) (evaluated using MAE, GMEAN, and MSE) and STS-B (evaluated using MAE, Pearson Correlation, and Spearman Correlation). Additionally, Table 10 shows the MAE and bMAE metrics the UCI-Abalone dataset.

## E.2 BALANCED METRICS FOR RISE

To address the challenges of evaluating models on imbalanced data distributions, particularly for tail labels, we employ three balanced metrics as defined in Ren et al. (2022). These metrics are designed to provide a more equitable assessment across all data regions by dividing the label space into even sub-regions, enabling a fairer evaluation.

The balanced Mean Squared Error (bMSE) is formulated as:

$$bMSE = -\log p_{train}(y|x;\theta) = -\log p_{bal}(y|x;\theta) \cdot \frac{p_{train}(y)}{\int_{Y} p_{bal}(y'|x;\theta) \cdot p_{train}(y')dy'}$$
(19)

This formulation comprises two components: the standard MSE loss and a balancing term to mitigate distribution mismatch between training and testing. Balanced metrics such as balanced Mean Absolute Error (bMAE) and balanced Geometric Mean Error (bGMEAN) are used to fairly assess performance across regions. bMAE averages errors within each sub-region or bins before computing the overall mean; formally for B bins with  $j^{th}$  bin containing  $N_j$  datapoints with y being the golden label and  $\hat{y}$  being the prediction, eq. 20 describes the formula for bMAE computation.

**Table 9:** Results on Dataset B (Rothe et al. (2018)) and STS-B dataset. The best baseline result for each metric and data subset is in red, best RISE version in blue, and the overall best result is in **bold**.

		L1 (N	1AE) ↓			GME	AN↓			MS	E↓	
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few
Dataset B												
Baseline Methods												
VANILLA	8.04	7.21	15.18	25.89	4.53	4.13	10.77	18.80	137.82	108.62	365.43	954.03
BalancedMSE	8.10	7.57	12.27	22.98	4.68	4.46	7.05	13.17	139.70	117.19	305.12	848.52
LDS+FDS	7.68	7.07	12.78	21.87	4.33	4.07	7.48	12.72	129.18	105.55	313.90	785.49
RankSIM	7.68	7.12	12.30	21.46	4.33	4.12	6.61	12.47	129.12	106.19	304.08	799.94
SRL	7.71	7.10	12.81	21.52	4.32	4.09	7.01	13.58	133.16	107.77	339.95	771.71
RISE Methods												
VANILLA+RISE	8.11	7.24	14.98	25.00	4.73	4.17	11.68	17.67	136.60	110.18	319.45	934.62
BalancedMSE+RISE	8.25	7.56	12.87	22.08	4.90	4.58	7.43	13.03	137.13	111.55	309.90	704.25
LDS+FDS+RISE	7.71	7.09	12.94	21.60	4.35	4.08	7.68	13.31	129.84	105.13	316.23	779.27
RankSIM+RISE	<b>7.67</b>	7.07	12.29	21.46	4.32	4.11	6.63	12.53	129.11	106.23	303.55	799.58
SRL +RISE	7.70	7.18	11.92	20.92	4.34	4.15	6.41	11.74	129.20	107.31	294.51	783.00
		L1 (M	1AE) ↓		Pe	arson Corr	elation (%	5)↑	Sp	earman cor	relation (%	) ↑
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few
STS-B												
Baseline Methods												
LDS+FDS	0.77	0.72	0.98	0.75	76.27	74.08	66.07	76.60	76.27	70.75	54.95	74.88
RankSIM	0.75	0.75	0.77	0.67	77.28	72.15	69.32	86.84	77.39	69.57	48.05	89.34
SRL	0.89	0.85	1.07	0.95	68.83	62.98	63.96	73.65	68.92	59.72	51.07	82.14
RISE Methods												
LDS+FDS+RISE	0.75	0.73	0.86	0.68	76.38	72.05	68.81	80.92	75.26	69.31	54.09	79.68
RankSIM+RISE	0.74	0.73	0.75	0.67	77.50	72.16	72.06	86.91	77.41	69.54	45.70	90.15
SRL+RISE	0.84	0.83	0.91	0.81	70.14	64.33	64.83	74.58	69.87	61.26	47.66	76.61

**Table 10:** Mean Absolute Error (MAE) results on UCI-Abalone dataset. Lower values indicate better performance. The best of the baseline and baseline+RISE pair is in **bold** and the best overall metric is underlined.

		$MAE\downarrow$				
Method	Many	Medium	Few	All		
VANILLA	1.77	5.46	9.98	2.56		
VANILLA + RISE	<b>1.59</b>	<b>5.19</b>	<b>9.75</b>	<b>2.34</b>		
BalancedMSE	2.50	5.41	4.61	3.43		
BalancedMSE + RISE	<b>1.30</b>	<b>2.35</b>	<b>4.53</b>	<b>1.53</b>		
LDS+FDS	2.80	4.44	7.64	3.18		
LDS+FDS + RISE	<b>2.07</b>	<b>2.91</b>	<b>7.16</b>	<b>2.30</b>		

$$bMAE = \frac{1}{B} \sum_{j=1}^{B} \frac{1}{N_j} \sum_{i=1}^{N_j} \|y - \hat{y}\|$$
 (20)

bGMEAN is formulated similarly but uses the geometric mean instead of MAE to highlight disparities across regions. These metrics are especially important for long-tailed distributions, where standard metrics may disproportionately reflect majority class performance. For our purposes, we chose to use bMAE to compare different RISE configurations. Due to space limitations for Dataset A (Moschoglou et al. (2017)), we had only reported the SRL result in the main paper. Therefore, we present the bMAE metric across different baselines in Table 11. Similarly we provide bMAE metrics for Dataset B (Rothe et al. (2018)) and STS-B in 12, and the balanced metrics for UCI-Abalone in 13.

**Table 11:** bMAE Results: Baseline vs RISE Methods on Dataset A (Moschoglou et al. (2017)). The best of the baseline and baseline+RISE pair is in **bold** and the best overall metric is underlined.

		Baseline Methods				eline + R	ISE Meth	nods
Method	All	Many	Med	Few	All	Many	Med	Few
VANILLA BalancedMSE LDS+FDS RankSIM	13.14 <b>8.70</b> 8.79 8.06	9.96 8.44 6.91 <b>6.49</b>	12.85 8.96 8.28 7.85	<b>19.81 11.43</b> 12.94 11.40	12.84 8.98 8.40 7.92	9.40 7.23 6.79 6.58	11.66 8.16 8.09 7.36	20.62 13.06 <b>11.87</b> <b>11.01</b>
SRL	8.32	6.64	8.34	11.74	<u>7.39</u>	<u>6.00</u>	<u>7.25</u>	<u>10.33</u>

**Table 12:** Balanced Mean Absolute Error (bMAE) results on Dataset B (Rothe et al. (2018)) and STS-B dataset. Lower values indicate better performance. The best of the baseline and baseline+RISE pair is in **bold** and the best overall metric is underlined.

		Base	eline		Baseline + RISE					
Method	All	Many	Med	Few	All	Many	Med	Few		
	Data	set B (Ro	the et al.	(2018))						
VANILLA	13.93	7.32	15.92	32.80	13.21	7.38	14.97	30.90		
BalancedMSE (Ren et al. (2022))	12.65	7.64	12.69	28.10	12.54	7.62	12.47	28.10		
LDS+FDS (Yang et al. (2021))	12.53	7.14	13.25	28.65	12.42	7.17	13.21	27.95		
RankSIM (Gong et al. (2022))	12.56	7.19	12.80	28.95	12.56	7.18	12.79	27.97		
SRL (Dong et al. (2025))	12.30	7.18	13.09	27.54	<u>12.28</u>	<u>7.14</u>	12.32	26.2		
		S	ΓS-B							
LDS+FDS (Yang et al. (2021))	0.77	0.73	0.84	0.79	0.73	0.74	0.77	0.70		
RankSIM (Gong et al. (2022))	0.72	0.76	0.72	0.66	0.71	0.74	0.71	0.65		
SRL (Dong et al. (2025))	0.87	0.85	0.88	0.88	$\overline{0.80}$	0.84	0.76	0.66		

**Table 13:** Balanced Mean Absolute Error (bMAE) results on UCI-Abalone dataset. Lower values indicate better performance. The best of the baseline and baseline+RISE pair is in **bold** and the best overall metric is underlined.

		bMAE	$\downarrow$	
Method	Many	Medium	Few	All
VANILLA	1.68	5.42	9.75	4.44
VANILLA + RISE	<b>1.58</b>	<b>5.20</b>	<b>9.74</b>	<b>4.32</b>
BalancedMSE + RISE	1.43	2.26	4.86	2.28
	<b>1.31</b>	2.23	<b>4.86</b>	<b>2.21</b>
LDS+FDS	2.64	4.66	7.64	4.22
LDS+FDS + RISE	<b>2.00</b>	<b>4.18</b>	7.64	<b>3.74</b>

#### E.3 COMPLETE ABLATION RESULTS

For brevity, the main paper only presented the L1 (MAE) metric for various ablations on Dataset A. In this section, we present the results across multiple metrics. Table 14 shows the complete ablation for different RISE-Train Strategies, Table 15 shows the ablation for different RISE-Infer strategies, and lastly, Table 16 provides complete results comparing RISE with ensembles with similar and increased capacity.

To strengthen our findings and validate the optimal RISE strategy beyond the Dataset A (Moschoglou et al. (2017)) dataset, we present comprehensive ablation studies on the Dataset B (Rothe et al. (2018)) dataset. Table 17 demonstrates that RISE (T2) consistently outperforms RISE (T1) across all metrics (MAE, GMEAN, and MSE) and data subsets, confirming the superiority of the T2 training configuration observed on Dataset A (Moschoglou et al. (2017)). Furthermore, Table 18 provides detailed architectural ablation results, showing that the optimal configuration uses K=2 experts with an upsampling ratio of 3, which achieves the best overall performance with an MAE of 7.67. Additionally, Table 19 examines different inference strategies, revealing that the held-out-based router (I3) consistently outperforms both expert averaging (I1) and train-based routing (I2), achieving the

best results across all metrics and data subsets with significant improvements in the Few subset. These results on Dataset B (Rothe et al. (2018)) corroborate our Dataset A (Moschoglou et al. (2017)) findings and demonstrate the robustness of our proposed RISE methodology across different long-tailed regression datasets.

**Table 14:** Complete ablation of RISE-Train on Dataset A (Moschoglou et al. (2017)) with SRL (Dong et al. (2025)) backbone, across multiple metrics. Best results in **bold**.

		L1 (M	AE)↓			GMEAN ↓				MSE ↓				
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few		
RISE (T1) RISE (T2)	7.23 <b>6.57</b>	6.77 <b>6.16</b>	7.95 <b>7.36</b>	9.61 <b>8.30</b>	4.44 <b>3.61</b>	4.15 <b>3.40</b>	4.94 <b>4.14</b>	6.19 <b>4.33</b>	92.54 <b>82.01</b>	80.12 <b>70.88</b>	110.96 <b>100.90</b>	158.86 <b>134.93</b>		

**Table 15:** Complete ablation of RISE inference strategies with SRL backbone on Dataset A, across multiple metrics. Best results in **bold** 

	L1 (MAE) ↓				$GMEAN\downarrow$				$MSE\downarrow$			
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few
Baseline SRL	7.23	6.64	8.28	9.85	4.53	4.17	5.32	6.35	91.79	77.20	115.83	163.15
Expert average (I1) Train-based Router (I2) Held-out-based Router (I3) Train+Held-out Router	7.23 7.26 <b>6.57</b> 7.18	6.72 6.61 <b>6.16</b> 6.62	8.13 8.34 <b>7.36</b> 8.15	9.54 10.33 <b>8.30</b> 9.84	4.51 4.56 <b>3.61</b> 4.42	4.20 4.15 <b>3.40</b> 4.11	5.16 5.48 <b>4.14</b> 5.09	6.16 6.76 <b>4.33</b> 5.92	91.73 92.11 <b>82.01</b> 90.79	78.85 76.48 <b>70.88</b> 76.68	112.50 116.26 <b>100.90</b> 112.53	156.00 173.01 <b>134.93</b> 163.86

**Table 16:** Complete comparison of RISE vs. traditional ensembles on Dataset A, across multiple metrics. Best results in **bold** 

			L1 (MAE) ↓						
Experiment	Additional Parameters	All	Many	Median	Few	All	Many	Median	Few
SRL SRL+ RISE (K=3) SRL: 3 ensemble SRL: 5 ensemble	0 2,100,224 3,150,336 5,250,560	91.79 <b>80.72</b> 91.66 91.56	77.20 <b>69.06</b> 77.04 76.75	115.83 <b>99.88</b> 115.65 115.83	163.15 <b>137.95</b> 163.43 164.31	7.23 <b>6.45</b> 7.22 7.22	6.64 <b>6.00</b> 6.63 6.62	8.28 <b>7.22</b> 8.28 8.30	9.85 <b>8.49</b> 9.86 9.90

**Table 17:** Ablation results for Dataset B Rothe et al. (2018) comparing different RISE-TRAIN configurations. The overall best result is in **bold**.

		L1 (	MAE) ↓			GM	IEAN↓		MSE ↓			
Method	All	Many	Median	Few	All	Many	Median	Few	All	Many	Median	Few
RISE (T1) RISE (T2)	7.94 <b>7.67</b>	7.46 <b>7.11</b>	12.66 <b>12.29</b>	22.75 <b>21.46</b>	4.50 <b>4.32</b>	4.33 <b>4.11</b>	6.89 <b>6.63</b>	14.66 <b>12.53</b>	139.69 <b>129.11</b>	118.35 <b>106.23</b>	339.35 <b>303.55</b>	829.19 <b>799.58</b>

**Table 18:** Ablation results for K=2 with varying upsampling rates (left) and for  $\alpha$ =3 with varying expert numbers (K) (right) on Dataset B( Rothe et al. (2018)). L1 (MAE) metric is shown. The overall best result is in **bold**.

		L1 (	MAE) ↓	
Config	All	Many	Median	Few
α=1	7.86	7.17	13.67	23.15
$\alpha=2$	7.81	7.15	13.35	22.71
$\alpha=3$	7.67	7.11	12.29	21.46
$\alpha$ =4	7.70	7.12	12.48	21.64
$\alpha$ =5	7.68	7.12	12.55	21.79

		L1 (	MAE) ↓	
Config	All	Many	Median	Few
K=2	7.67	7.11	12.29	21.46
K=3	7.69	7.17	11.89	21.20
K=4	7.78	7.20	12.61	22.27
K=5	8.28	7.43	15.67	25.46

**Table 19:** Ablation results comparing different RISE-INFERENCE configurations on Dataset B (Rothe et al. (2018)). The best baseline result for each metric and data subset is in red, and the overall best result is in **bold**.

	L1 (MAE) ↓				GMEAN ↓				$MSE\downarrow$			
Method	All	Many	Median	Few	All	Many	Median	Few	All	Many	Median	Few
Baseline Methods RankSIM	7.68	7.12	12.30	21.46	4.33	4.12	6.61	12.47	129.12	106.19	304.08	799.94
RISE Inference Strategies Expert average (I1) Train-based Route (I2) Held-out-based router (I3)	8.32 8.00 <b>7.67</b>	7.62 7.37 <b>7.11</b>	14.22 13.37 <b>12.29</b>	17.33 14.28 <b>12.53</b>	4.79 4.57 <b>4.32</b>	4.47 4.30 <b>4.11</b>	8.91 7.97 <b>6.63</b>	17.33 14.28 <b>12.53</b>	143.35 135.82 <b>129.11</b>	117.02 111.01 <b>106.23</b>	351.47 329.48 <b>303.55</b>	855.11 826.43 <b>799.58</b>

## E.4 RISE PERFORMANCE WITH BEST-PERFORMING ROUTER CONFIGURATION

To assess the robustness of our approach, we perform five independent experimental runs and report the mean and standard deviation for Dataset A, B & STS-B on each performance metric in Table 20 and the balanced metrics with error bars for Dataset A are reported in Table 21. This evaluation provides statistical insight into the consistency and reliability of the results. For each run, the router is trained and the backbone model achieving the highest routing accuracy on the validation set  $D_{\text{val}}$  is selected for reporting. Router with the SRL backbone is picked for the the Dataset A (Moschoglou et al. (2017)) dataset, while RankSim backbone is utilized for both IMDB and STS datasets.

Our proposed RISE paradigm consistently outperforms its corresponding baseline methods across multiple metrics, with particularly notable gains in medium- and few-shot regions—where imbalanced regression models typically underperform. These improvements are statistically significant, often exceeding standard error margins. For instance, on the Dataset A (Moschoglou et al. (2017)) dataset, SRL+RISE achieves a 13.7% reduction in Few-shot MAE (9.85  $\rightarrow$  8.50) and a 12.6% reduction in Medium-shot MAE (8.35  $\rightarrow$  7.30), alongside a 28.4% improvement in Few-shot GMEAN (6.34  $\rightarrow$  4.54). Similar trends are observed in Dataset B (Rothe et al. (2018)), where BalancedMSE+RISE lowers Few-shot MAE by 9.8% (23.24  $\rightarrow$  20.97), and in STS, where LDS+FDS+RISE improves Medium-shot MAE by 11.2% (0.98  $\rightarrow$  0.87).

While RISE generally maintains or improves performance in majority (Many-shot) regions, there are isolated instances where baseline models marginally outperform RISE. For example, in Dataset A (Moschoglou et al. (2017)), RankSIM achieves a slightly lower Many-shot MAE (6.48 vs. 6.56), and in Dataset B (Rothe et al. (2018)), LDS+FDS reports a marginally better Many-shot MSE (106.61 vs. 107.06). However, these differences are minor and fall within overlapping standard deviation intervals.

Importantly, RISE demonstrates strong generalization by significantly improving performance in minority regions while preserving accuracy on majority classes. This balance highlights the effectiveness of RISE in addressing the fundamental challenge of imbalanced regression, offering a scalable and principled solution for real-world settings.

## F BROADER IMPACT

RISE offers a practical and efficient alternative to end-to-end training by leveraging pre-trained models. Unlike typical deep learning approaches, it requires training only the expert heads and router network while keeping the backbone frozen. This lightweight design makes it feasible for large-scale models and suitable for scenarios where full retraining is impractical. While our experiments used the full training set, RISE can potentially be adapted for final-layer tuning using only a small validation set, as supported by recent adaptation methods (Kirichenko et al. (2023)).

RISE differs from standard fine-tuning by targeting specific regions of poor performance—often underrepresented or minority subgroups—through expert specialization. This targeted improvement enhances fairness, particularly in sensitive applications like healthcare or finance, where disparities in prediction can have serious consequences. By improving minority performance without sacrificing majority accuracy, RISE moves toward more equitable and efficient machine learning systems.

**Table 20:** Comparison of RISE-paired with the baseline methods across Dataset A (Moschoglou et al. (2017)), Dataset B Rothe et al. (2018), and STS datasets. Results show MAE, GMEAN, and MSE metrics for different data segments (All, Many-shot, Medium-shot, Few-shot). Values are reported as mean  $\pm$  standard deviation. Best results for each metric and data subset are in bold, we also report the router accuracy for each RISE configuration in parentheses.

		MA	Æ↓			GME	AN↓		MSE ↓				
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few	
Dataset A													
VANILLA	11.06	9.99	12.90	16.65	7.08	6.30	8.41	13.57	203.69	165.70	275.75	367.1	
	±0.01	±0.05	±0.14	±0.29	±0.03	±0.05	±0.15	±0.34	±1.13	±2.32	±7.13	±11.5	
VANILLA+RISE (0.60)	10.07	9.20	11.19	15.33	6.18	5.52	7.19	11.99	173.84	146.76	211.69	328.2	
	±0.04	±0.08	±0.14	±0.21	±0.04	±0.06	±0.11	±0.28	±0.77	±2.37	±6.80	±6.86	
BalancedMSE	8.71	8.45	9.02	10.30	5.59	5.45	5.94	6.07	127.28	118.71	133.87	191.2	
	±0.06	±0.04	±0.20	±0.14	±0.06	±0.07	±0.14	±0.17	±1.57	±1.30	±5.93	±3.86	
BalancedMSE+RISE (0.72)	7.62	7.53	7.90	8.79	4.58	4.56	4.63	4.74	106.40	103.18	111.68	158.9	
	±0.05	±0.04	±0.15	±0.15	±0.06	±0.06	±0.11	±0.19	±1.31	±1.89	±4.43	±3.92	
LDS+FDS	7.47	6.92	8.23	10.52	4.77	4.46	5.30	6.84	95.23	79.98	118.33	177.2	
	±0.08	±0.11	±0.13	±0.25	±0.06	±0.07	±0.10	±0.23	±1.85	±2.73	±4.26	±5.09	
LDS+FDS+RISE (0.56)	7.27	6.85	7.91	9.54	4.51	4.30	4.81	6.08	92.59	80.18	113.40	153.6	
	±0.08	±0.10	±0.13	±0.24	±0.06	±0.08	±0.10	±0.21	±1.81	±2.61	±4.31	±6.32	
 RankSIM	7.01	6.48	7.82	9.85	4.55	4.14	5.37	7.04	83.23	71.48	98.21	154.2	
	±0.04	±0.06	±0.08	±0.09	±0.04	±0.04	±0.05	±0.19	±0.77	±1.37	±3.40	±1.35	
RankSIM+RISE (0.55)	6.93	6.56	7.34	9.25	4.34	4.07	4.79	6.11	82.47	73.88	90.22	143.2	
	±0.04	±0.06	±0.08	±0.09	±0.04	±0.04	±0.03	±0.17	±0.72	±1.35	±3.48	±1.37	
	7.20	6.59	8.35	9.85	4.50	4.14	5.34	6.34	91.67	76.09	118.91	165.1	
	±0.02	±0.04	±0.08	±0.25	±0.03	±0.03	±0.12	±0.33	±0.64	±0.55	±1.70	±6.24	
SRL +RISE (0.87)	6.43	5.96	7.30	8.50	3.36	3.13	3.87	4.54	80.70	67.98	103.09	140.3	
	±0.02	±0.03	±0.09	±0.22	±0.02	±0.03	±0.12	±0.22	±0.65	±0.50	±1.60	±5.82	
D-44 D													
<b>Dataset B</b> VANILLA	8.04	7.20	15.18	26.20	4.51	4.11	10.69	18.81	137.96	108.17	366.46	972.0	
	±0.03	±0.03	±0.12	±0.15	±0.03	±0.02	±0.09	±0.25	±0.87	±0.63	±6.46	±7.5	
VANILLA+RISE (0.85)	7.91	7.22	13.65	24.73	4.45	4.15	8.38	16.59	135.34	108.90	333.20	925.6	
	±0.03	±0.03	±0.13	±0.13	±0.03	±0.03	±0.07	±0.22	±0.90	±0.67	±6.67	±7.13	
BalancedMSE	8.10	7.56	12.27	23.24	4.68	4.45	7.10	13.25	139.62	116.96	302.67	868.3	
	±0.03	±0.03	±0.17	±0.21	±0.01	±0.01	±0.11	±0.26	±1.55	±1.15	±9.31	±14.0	
BalancedMSE+RISE (0.81)	7.73	7.28	12.12	20.97	4.41	4.29	6.82	11.99	136.36	108.79	300.42	820.7	
	±0.03	±0.02	±0.13	±0.25	±0.01	±0.00	±0.07	±0.52	±1.50	±1.01	±8.45	±13.0	
	7.70	7.13	12.54	21.84	4.32	4.11	7.55	12.75	129.91	106.61	310.90	781.8	
	±0.01	±0.01	±0.06	±0.39	±0.01	±0.01	±0.10	±0.38	±0.85	±0.57	±2.90	±21.5	
LDS+FDS+RISE (0.81)	7.64	7.11	12.09	21.24	4.27	4.07	6.46	12.17	131.02	107.06	301.96	768.1	
	±0.02	±0.01	±0.06	±0.38	±0.01	±0.01	±0.09	±0.37	±0.86	±0.61	±3.08	±20.2	
	7.69	7.12	12.33	21.55	4.33	4.12	6.65	12.68	129.14	106.78	302.58	802.8	
	±0.02	±0.02	±0.12	±0.37	±0.01	±0.01	±0.08	±0.41	±0.72	±0.43	±5.76	±28.1	
RankSIM+RISE (0.8)	7.66	7.11	12.08	21.38	4.30	4.09	6.48	12.54	127.49	108.04	298.66	800.7	
. ,	±0.02	±0.02	±0.12	±0.37	±0.01	±0.01	±0.08	±0.41	±0.74	±0.33	±5.90	±29.0	
	7.70	7.13	12.66	 21.94	4.34	4.13	6.93	12.93	131.96	107.38	337.57	 768.8	
-	±0.02	±0.03	±0.09	±0.50	±0.01	±0.01	±0.04	±0.47	±1.11	±0.69	±5.73	±25.8	
SRL+RISE (0.79)	7.68	7.19	11.98	19.39	4.35	4.15	6.43	11.22	130.07	107.54	298.14	773.4	
	±0.02	±0.02	±0.09	±0.50	±0.02	±0.01	±0.04	±0.47	±1.11	±0.70	±5.63	±24.9	
STS-B													
LDS+FDS	0.77	0.72	0.98	0.76	0.38	0.33	0.67	0.45	0.91	0.81	1.06	0.94	
I De EDe Dice (0.51)	±0.00	±0.01	±0.02	±0.02	±0.01	±0.01	±0.02	±0.01	±0.01	±0.01	±0.05	±0.00	
LDS+FDS+RISE (0.51)	<b>0.75</b> ±0.00	0.73 ±0.00	<b>0.87</b> ±0.02	<b>0.66</b> ±0.02	<b>0.30</b> ±0.01	<b>0.25</b> ±0.01	<b>0.56</b> ±0.01	<b>0.34</b> ±0.02	0.92 ±0.01	<b>0.79</b> ±0.00	1.08 ±0.06	<b>0.76</b> ±0.04	
RankSIM	0.76	0.74	0.75	0.64	0.50	0.47	0.54	0.37	0.86	0.86	0.85	0.63	
RankSIM+RISE (0.55)	±0.00 <b>0.73</b>	±0.01 <b>0.73</b>	±0.01 <b>0.75</b>	±0.04 <b>0.63</b>	±0.02 <b>0.39</b>	±0.02 <b>0.37</b>	±0.01 <b>0.54</b>	±0.03 <b>0.36</b>	±0.01 <b>0.84</b>	±0.02 <b>0.84</b>	±0.01 <b>0.85</b>	±0.03	
KankonvitkioE (U.JJ)	±0.01	±0.01	±0.01	±0.05	±0.01	±0.02	±0.02	±0.04	±0.02	±0.02	±0.02	±0.09	
	0.00	0.04	1.07		0.62		0.70	0.60	1.17	1.07	157	1 20	
SRL	0.89 ±0.01	0.84 ±0.01	1.07 ±0.04	0.98 ±0.07	0.63 ±0.02	0.57 ±0.02	0.79 ±0.05	0.69 ±0.09	1.17 ±0.02	1.07 ±0.02	1.57 ±0.08	1.30 ±0.12	
SRL+RISE (0.57)	0.82	0.80	0.93	0.76	0.43	0.39	0.70	0.36	1.06	1.01	1.24	1.14	
. ,	±0.01	±0.00	±0.03	±0.07	±0.01	±0.02	±0.04	±0.05	±0.018	±0.01	±0.08	±0.1	

**Table 21:** Comparison of RISE with baseline methods for Dataset A (Moschoglou et al. (2017)) with balanced metrics. Values are reported as mean ± standard deviation. Best results for each metric and data subset are in bold.

		bMA	AE↓			bGMI	EAN↓		$bMSE \downarrow$				
Method	All	Many	Med	Few	All	Many	Med	Few	All	Many	Med	Few	
Dataset A													
VANILLA	13.18	9.99	12.94	19.84	7.30	6.30	8.41	13.57	271.42	165.70	276.55	483.80	
	±0.06	±0.05	±0.16	±0.25	±0.08	±0.05	±0.15	$\pm 0.34$	±3.33	$\pm 2.32$	±8.22	±10.75	
VANILLA+RISE	12.15	9.20	11.18	18.81	6.10	5.52	7.19	11.99	242.73	146.76	210.97	458.86	
	±0.04	±0.08	±0.15	±0.17	±0.08	±0.06	±0.11	±0.28	±2.29	±2.37	±7.28	±8.16	
BalancedMSE	9.35	8.45	8.99	11.40	6.44	5.45	5.94	6.07	153.18	118.71	132.34	236.83	
	±0.07	±0.04	±0.22	±0.14	±0.10	±0.07	±0.14	±0.17	±1.89	±1.30	±6.34	±3.38	
BalancedMSE+RISE	8.44	7.53	7.82	10.61	5.35	4.56	4.63	4.70	134.30	103.18	110.51	217.58	
	±0.07	±0.04	±0.17	±0.16	±0.11	±0.06	±0.11	±0.19	±1.72	±1.89	±4.80	±3.79	
LDS+FDS	9.35	8.45	8.99	11.40	5.74	5.45	5.94	6.07	153.18	118.71	132.34	236.83	
	±0.07	$\pm 0.04$	±0.22	±0.14	±0.11	±0.07	±0.14	$\pm 0.17$	±1.89	±1.30	±6.34	±3.38	
LDS+FDS+RISE	8.31	6.85	7.91	11.55	5.61	4.30	4.81	6.08	122.83	80.18	112.24	216.64	
	±0.07	±0.10	±0.14	±0.25	±0.09	±0.08	±0.10	±0.31	±1.05	±2.61	±4.48	±6.09	
RankSIM	8.07	6.48	7.81	11.46	6.14	4.14	5.37	7.04	111.49	71.48	97.46	202.09	
	±0.04	±0.06	$\pm 0.08$	±0.05	±0.07	±0.04	±0.05	±0.19	$\pm 0.47$	±1.37	±3.37	±0.88	
RankSIM+RISE	7.92	6.56	7.31	11.08	5.09	4.07	4.79	6.11	108.90	73.88	89.22	192.92	
	±0.03	±0.06	±0.08	±0.05	±0.06	±0.04	±0.03	±0.17	±0.43	±1.35	±3.45	±1.04	
SRL	8.28	6.59	8.41	11.65	5.15	4.14	5.34	6.34	121.11	76.09	118.76	214.45	
	±0.04	±0.04	±0.08	±0.18	±0.06	±0.03	±0.12	±0.33	±1.26	±0.55	±1.58	±4.64	
SRL+RISE	7.36	5.96	7.32	10.24	4.40	3.13	3.87	4.54	105.65	67.98	102.10	184.75	
	±0.04	±0.03	±0.09	±0.17	±0.05	±0.03	±0.12	±0.22	±1.26	±0.50	±1.50	±4.43	