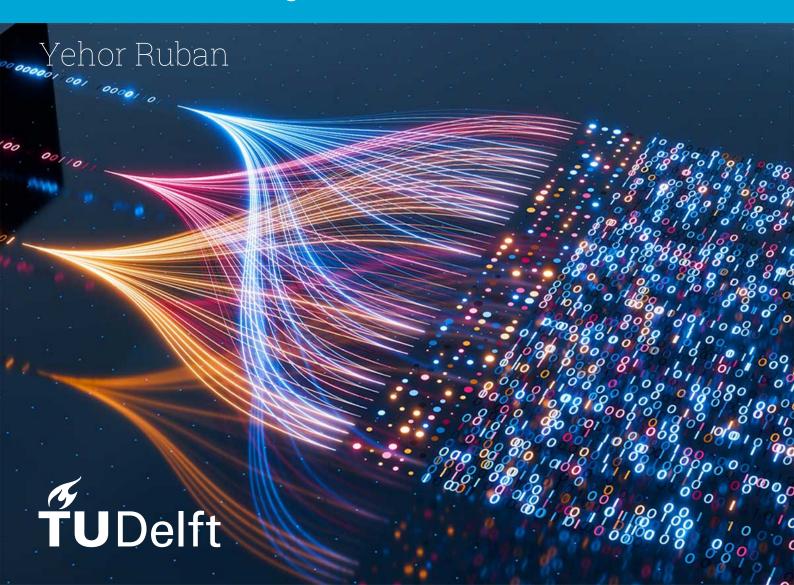
Ridge Regression and Random Neural Networks under High-Dimensional Conditions

Derivative-Based Calibration for Stable Learning



Ridge Regression and Random Neural Networks under High-Dimensional Conditions

Derivative-Based Calibration for Stable Learning

by

Yehor Ruban

Instructor: N. Parolya

Faculty: Electrical Engineering, Mathematics

and Computer Science



Summary

This thesis investigates the problem of selecting an optimal regularization parameter in high-dimensional ridge regression models with random features. The work is situated at the intersection of machine learning, statistical signal processing, and random matrix theory, and aims to improve the understanding and stability of regression-based learning in high-dimensional settings.

When the number of model parameters becomes comparable to the number of training samples, conventional statistical assumptions no longer hold, and model performance can exhibit sharp transitions or instability. In such regimes, ridge regularization plays a critical role in balancing model bias and variance. However, determining the optimal regularization parameter γ is non-trivial: existing techniques such as cross-validation are computationally costly, while analytical approaches based on deterministic equivalents—such as the method proposed by Couillet and Liao [6]—rely on asymptotic knowledge of population statistics that may be inaccessible in practice.

To address this challenge, this thesis introduces a new *derivative-based calibration rule* for ridge regression in random-feature neural networks. The proposed method identifies the optimal value of γ directly from data by analyzing the empirical derivative of the training loss with respect to γ . This criterion provides a fully data-driven, computationally efficient alternative to deterministic-equivalent or cross-validation-based calibration.

Theoretical development and analysis are supported by controlled numerical experiments. Synthetic data simulations confirm that the derivative-based calibration achieves predictive performance comparable to Couillet's deterministic-equivalent rule across a wide range of dimensional ratios. Furthermore, perturbation analysis demonstrates that the proposed method yields a more stable feature-weight vector $\boldsymbol{\beta}$ under small input variations, indicating improved robustness and reproducibility of the learned model.

Validation on real data is performed using the Fashion–MNIST dataset in a binary classification setup. Results show that both calibration rules produce similar test accuracies (above 99%), while the derivative-based approach consistently exhibits lower coefficient variance and better recall. This stability of β suggests that the model learns a more reliable internal representation, which can be advantageous in practical applications where interpretability, consistency, and downstream reuse of learned features are important.

Finally, the discussion situates these findings in the context of broader research on high-dimensional learning stability, linking them to studies in bioinformatics, neuroimaging, and industrial fault detection where reproducibility of learned representations is critical. The thesis concludes that derivative-based calibration provides a simple, effective, and theoretically grounded alternative to existing methods for ridge parameter selection, combining strong empirical performance with enhanced stability and interpretability.

Contents

Summary						
1	Intro	oduction	1			
	1.1	Background and Motivation	1			
	1.2		1			
	1.3	Thesis Structure	2			
2	Lite	rature Review	3			
			3			
			4			
		·	4			
	2.2		5			
		2.2.1 Random Weights and Efficient Feature Extractors	5			
		· ·	6			
	2.3		7			
			7			
			7			
	2.4		8			
		· · · · · · · · · · · · · · · · · · ·	8			
			9			
	2.5	·	0			
_		•				
3			11			
	3.1		11			
			11			
			3			
		,	6			
			8			
	3.3		9			
		· ·	9			
		· · · · · · · · · · · · · · · · · · ·	9			
			20			
			20			
	3.4	5	21			
		3.4.1 Potential cause of the problem	21			
4	Results 22					
	4.1	Research on the original Experiment	2			
		4.1.1 Mean Squared Error against the regularization parameter γ	22			
		4.1.2 Optimal γ value and error behaviour with increasing amount of data	24			
			25			
			25			
	4.2	,	28			
	4.3		30			
	4.4	·	33			
	4.5	· · · · · · · · · · · · · · · · · · ·	34			
	-		34			
		•	37			

Contents

 5 Discussion 5.1 Summary of Findings 5.1.1 Recap 5.1.2 Interpretations of the Results 5.2 Advantages of the Proposed Derivative-Based Calibration over Couillet's R 5.2.1 Conceptual Advantages - Theory 5.2.2 Applied Advantages - Practice 6 Conclusion 6.1 Summary and Final Remarks 6.2 Future Work References A Source Code A.1 Mean Squared Error against the regularization parameter γ A.2 Optimal γ value and error behaviour with increasing amount of data A.3 Normalized ridge error against regularization parameter γ 		6 6 7 8 8 2
6.1 Summary and Final Remarks	5	
A Source Code A.1 Mean Squared Error against the regularization parameter γ A.2 Optimal γ value and error behaviour with increasing amount of data	5	3
A.1 Mean Squared Error against the regularization parameter γ A.2 Optimal γ value and error behaviour with increasing amount of data	5	4
 A.3 Normalized higge error against regularization parameter γ		3 8 0 4

1

Introduction

1.1. Background and Motivation

In modern machine learning, models with a large number of parameters often outperform smaller ones, even in cases where the number of features approaches or exceeds the number of training samples. This high-dimensional regime challenges traditional statistical assumptions and motivates the need for new analytical tools to understand model behavior, stability, and generalization.

Neural networks, especially those with random or fixed hidden weights, provide a simple yet powerful framework for studying these questions. When the hidden-layer weights are drawn at random, the network acts as a nonlinear feature extractor, and the final linear layer reduces to a regression problem in the resulting random feature space. This connection allows techniques from classical statistics and random matrix theory to be used for analysis.

Among the most important components in such models is the choice of the regularization parameter in ridge regression. It controls the balance between fitting the data and maintaining stability of the learned coefficients. Traditional tuning methods, such as cross-validation or asymptotic formulas, can be computationally costly or rely on strong model assumptions. This motivates the search for alternative, data-driven calibration approaches that can perform reliably in high-dimensional settings.

1.2. Research Objectives

The main objective of this thesis is to develop and analyze a new method for selecting the ridge regularization parameter in random-feature regression models. The proposed *derivative-based calibration rule* determines the optimal regularization strength by analyzing the empirical derivative of the training loss with respect to the regularization parameter. The approach is entirely data-driven and avoids the need for repeated training or asymptotic parameter estimation.

This work builds upon the theoretical framework introduced by Romain Couillet and co-authors in their book *Random Matrix Methods for Machine Learning* [6] and related article [29], which provide deterministic-equivalent analyses of high-dimensional ridge regression and random-feature neural networks. Their framework offers a rigorous, asymptotic approach to understanding how regularization behaves when both the number of features and the number of samples grow large. Within this context, Couillet and colleagues proposed a deterministic-equivalent calibration method for choosing the optimal regularization parameter based on random matrix theory.

The present thesis extends this line of work by introducing an alternative, fully data-driven calibration scheme that does not rely on asymptotic knowledge of model parameters. Specifically, the goals are

1.3. Thesis Structure 2

to:

• Formulate the derivative-based regularization rule within the framework of random-feature ridge regression;

- Compare its performance to Couillet's deterministic-equivalent calibration method, both theoretically and empirically;
- Evaluate the method's stability under data perturbations and assess how it affects the learned coefficients β;
- Validate the findings through controlled simulations and real-data experiments using the Fashion— MNIST dataset.

Through these objectives, the work seeks to contribute both a practical calibration tool and a deeper understanding of how regularization influences model stability in high-dimensional random networks.

1.3. Thesis Structure

The remainder of this thesis is organized as follows. Chapter 2 reviews the theoretical background, including neural networks, random features, ridge regression, and random matrix theory. Chapter 3 develops the mathematical framework and introduces the proposed derivative-based calibration rule. Chapter 4 presents numerical experiments on both synthetic and real data, comparing the proposed method with existing approaches. Chapter 5 discusses the implications of the results, highlighting the theoretical and practical significance of coefficient stability. Finally, Chapter 6 concludes the thesis and outlines directions for future research.

Overall, the thesis combines theoretical analysis, numerical validation, and conceptual interpretation to provide a comprehensive view of regularization and stability in high-dimensional random-feature models.

2

Literature Review

Chapter 2 provides an overview of the main concepts and historical advances that underpin this thesis. We begin in Section 1 with a general and broad introduction to the concept of neural networks, highlighting both their historical origins and their progression into deep, multi-layer architectures. This sets the stage for why neural networks have become a dominant tool for complex learning problems and tasks, but also why it can be potentially challenging to analyze them in a theoretical way.

From there, Section 2 narrows the focus to random neural networks, emphasizing how random weights can be viewed as an efficient means of feature extraction, particularly in contexts with limited data or strict computational budgets. We also discuss the use of randomization in modern practice—ranging from dropout to "learning without backpropagation"—as these techniques share conceptual parallels with the simpler random networks studied in this thesis.

In Section 3, we demonstrate that regression can be seen as a natural counterpart to single-layer neural networks, especially through the lens of ridge regression on random features. Here, we connect the dots between "pure" linear regression on one hand and single-hidden-layer "extreme learning machines" on the other, underscoring how penalization (in the form of ridge regression) addresses overfitting and stabilizes parameter estimates in neural nets.

Following that, Section 4 addresses the high-dimensional regime and the corresponding role of random matrix theory (RMT). Modern applications often operate in scenarios where the number of parameters and/or features is comparable to (or larger than) the number of data points. We show that RMT, with its body of results on large random matrices and kernels, offers valuable insights into the behavior of large-scale regressions and random neural networks alike.

Finally, Section 5 synthesizes these threads—bridging random neural networks, ridge regression, and high-dimensional analysis—to foreshadow the technical work and contributions of this thesis. By understanding the interplay among these ingredients, we can pinpoint how random fixed-weight models can yield theoretical insights and practical benefits in high-dimensional settings.

2.1. Introduction to Neural Networks

Neural networks are one of the most important paradigms in machine learning, roughly inspired by the structure of actual biological neurons. The core idea involves organizing computational units (called 'neurons') into interconnected layers that in turn learn to approximate functions straight from the data. Over the past decades, this arrangement design has evolved from basic, single-layer perceptrons into deep, multi-layer architectures that can perform well in tasks such as image recognition, natural language processing, even content generation. All of the aforementioned has cemented neural networks

as the primary tool in modern understanding of artificial intelligence.

In this section, we first examine the historical perspective on neural networks, shedding light on how the early perceptrons established foundational ideas that still resonate today. Subsequently, we explore the journey from simple to deep networks, highlighting key developments and the shift toward highly expressive, large-scale models. This progression not only frames why advanced neural networks can be so powerful, but it also underscores the growing need to better understand simpler random-weight approaches both for computational practicality and for theoretical clarity.

2.1.1. Historical Perspective

The origins of artificial neural networks (ANNs) trace back to the mid-twentieth century, when researchers first began exploring simplified mathematical models of biological neurons. One of the earliest milestones was the perceptron model introduced by Rosenblatt in 1958 [41], which formalized a neuron as a linear threshold unit capable of learning simple decision boundaries. The perceptron embodied the promise of machine learning before the term was even widely used: a system that could adapt its parameters based on examples rather than explicit programming.

The optimism surrounding early neural networks was tempered by their theoretical and computational limitations. In 1969, Minsky and Papert [34] rigorously demonstrated that a single-layer perceptron could not represent non-linearly separable functions such as the XOR problem. This result, combined with limited computational resources, led to a decline in neural network research, often referred to as the first "AI winter." Nevertheless, the perceptron era established the mathematical and conceptual foundations for later advances in supervised learning.

Interest was reignited in the 1980s with the rediscovery and popularization of the *backpropagation* algorithm [42], which enabled efficient gradient-based training of multi-layer networks. Backpropagation, initially studied in control theory and later formalized in the context of neural computation, allowed the training of deep feedforward networks through layerwise error propagation. Around the same period, Hopfield [17] introduced recurrent neural networks as energy-based systems, and Kohonen [23] proposed self-organizing maps for unsupervised learning - broadening the scope of connectionist models beyond simple classification.

By the early 1990s, networks such as the multilayer perceptron (MLP) and convolutional neural network (CNN) [26] were already demonstrating strong performance on specific pattern-recognition tasks. However, the combination of limited computational power and small datasets again constrained progress. It was not until the 2000s, when large annotated datasets, increased computational resources (notably GPUs), and improved regularization methods such as dropout [45] became available - that neural networks began their modern resurgence. This period marked the advent of "deep learning" [25], characterized by very deep architectures and millions of parameters capable of learning highly complex, hierarchical feature representations.

In summary, the historical trajectory of neural networks can be viewed as an evolution from simple, interpretable linear classifiers to powerful but opaque high-dimensional systems. This progression contextualizes the present thesis: while modern deep networks achieve remarkable empirical success, understanding their behavior remains challenging. Studying simplified models, such as random-feature neural networks analyzed through the lens of ridge regression and random matrix theory - offers a principled way to recover theoretical tractability without abandoning the essential nonlinear mechanisms that make neural networks effective.

2.1.2. From Simple to Deep

While early neural networks such as the perceptron demonstrated the feasibility of data-driven learning, their representational power was limited to linearly separable problems. The introduction of hidden layers provided a crucial breakthrough: multi-layer networks could approximate arbitrary continuous

functions under mild assumptions, as formalized in the *universal approximation theorem* [8, 18]. This theoretical result established that even relatively shallow networks, when equipped with non-linear activation functions, possess the capacity to represent complex mappings between inputs and outputs.

In practical terms, however, training such networks remained difficult throughout the 1980s and 1990s. Deep architectures suffered from vanishing gradients, poor initialization, and overfitting in the absence of large datasets. Research therefore concentrated on specific network designs that exploited structure in the data. Notable among these were the convolutional neural networks (CNNs) of LeCun et al. [27], which leveraged spatial weight sharing for image data, and recurrent neural networks (RNNs) designed for sequential or temporal patterns [11]. These architectures anticipated modern deep learning by embedding strong inductive biases into network topology.

The modern "deep learning revolution" emerged in the mid-2000s, driven by three converging factors: increased computational power (notably GPU acceleration), large-scale datasets such as ImageNet [9], and methodological innovations in optimization and regularization. Layer-wise unsupervised pretraining using deep belief networks [15] and autoencoders [48] mitigated initialization problems, while activation functions such as the rectified linear unit (ReLU) [35] and regularization techniques like dropout [45] enabled stable gradient propagation and improved generalization. Collectively, these developments unlocked the practical potential of deep architectures containing tens or even hundreds of layers.

Subsequent milestones rapidly followed. Architectures such as AlexNet [24], VGG [44], ResNet [14], and Transformers [47] demonstrated that depth, residual connections, and attention mechanisms could be scaled to unprecedented levels, achieving state-of-the-art performance across vision, language, and multimodal tasks. Yet, as networks grew deeper and more parameter-rich, their theoretical understanding became more opaque. The same complexity that made deep networks powerful also made them difficult to analyze or interpret rigorously.

This realization has inspired a complementary line of research that seeks tractable, theoretically grounded models capturing the essential mechanisms of neural computation without the full complexity of deep learning. Among these are random-feature networks, kernel approximations, and mean-field analyses derived from random matrix theory. The present thesis situates itself within this latter tradition: by examining simplified random-weight neural networks through the lens of ridge regression, we aim to bridge the gap between deep learning practice and high-dimensional statistical theory.

2.2. Random Neural Networks

The study of random neural networks (RNNs)—also known as *random-feature models* or *extreme learning machines* (ELMs)—represents a significant simplification of traditional neural networks. In these architectures, the weights of certain layers (often the hidden or feature-extraction layers) are not trained, but instead initialized randomly and kept fixed throughout learning. Only the final layer, typically a linear readout, is optimized. This randomization transforms the learning problem from a non-convex optimization task into a tractable linear regression, enabling both analytical treatment and rapid computation.

2.2.1. Random Weights and Efficient Feature Extractors

The conceptual foundation of random-weight networks can be traced to early work on fixed-feature representations and kernel methods. Neal [36] first demonstrated that infinitely wide neural networks with random weights converge to Gaussian processes, establishing a deep connection between random networks and Bayesian kernel machines. Later, Rahimi and Recht [39] proposed the *random kitchen sinks* approach, showing that a wide class of shift-invariant kernels can be efficiently approximated by

random feature maps of the form

$$\phi(\mathbf{x}) = \sqrt{\frac{2}{N}} \begin{bmatrix} \cos(\mathbf{w}_{1}^{\top} \mathbf{x} + b_{1}) \\ \cos(\mathbf{w}_{2}^{\top} \mathbf{x} + b_{2}) \\ \vdots \\ \cos(\mathbf{w}_{N}^{\top} \mathbf{x} + b_{N}) \end{bmatrix},$$
(2.1)

where each \mathbf{w}_i is sampled independently from a distribution proportional to the Fourier transform of the kernel, and b_i is a random phase. This approach effectively replaces expensive kernel computations with low-dimensional random projections followed by a linear model, enabling scalable kernel learning in high dimensions.

A closely related stream of research explored randomization directly within neural architectures. In the *Extreme Learning Machine* framework [19], the hidden-layer weights W are drawn at random and fixed, while the output weights β are obtained in closed form through ridge regression:

$$\boldsymbol{\beta} = (\boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma} + \gamma I_N)^{-1} \boldsymbol{\Sigma}^{\mathsf{T}} \mathbf{Y}, \tag{2.2}$$

where $\Sigma = \sigma(\mathbf{W}\mathbf{X})$ denotes the matrix of hidden-layer activations and $\sigma(\cdot)$ is a nonlinearity such as ReLU or sigmoid. This simple construction can approximate complex mappings while avoiding the heavy computational cost of backpropagation. Huang and colleagues showed that even with completely random hidden weights, the ELM achieves strong generalization and extremely fast training times, laying the foundation for the modern random-feature viewpoint.

From a theoretical standpoint, random-weight networks can be understood as finite-dimensional approximations of kernel machines. The random mapping $\phi(\mathbf{x}) = \sigma(\mathbf{W}\mathbf{x})$ induces an implicit kernel

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}_{\mathbf{w}}[\sigma(\mathbf{w}^{\top} \mathbf{x}) \, \sigma(\mathbf{w}^{\top} \mathbf{x}')], \tag{2.3}$$

which for specific choices of activation functions yields well-known kernels such as the *arc-cosine kernel* [4]. The random network thus serves as a Monte Carlo estimator of this kernel, bridging the gap between neural and kernel-based learning.

Randomization also provides a form of implicit regularization. Fixed random features prevent the model from overfitting by constraining its expressiveness and decoupling feature learning from weight optimization. Furthermore, randomization simplifies theoretical analysis, allowing tools from random matrix theory and statistical physics to be applied directly to study generalization, eigenvalue spectra, and limiting distributions [38, 7]. These analyses have shown that even simple random networks exhibit complex phase transitions between under- and over-parameterized regimes—insights that are difficult to obtain in fully trained deep networks.

2.2.2. Random Neural Networks in Practice

Despite their simplicity, random-feature models have found wide use in practice. Rahimi and Recht's work [40] demonstrated that random Fourier features can scale kernel learning to millions of data points without loss in accuracy. The same principle underlies several modern architectures: reservoir computing [21, 30], echo-state networks, and certain forms of dropout regularization can all be interpreted as imposing structured randomness within neural dynamics.

In computer vision and signal processing, random convolutional features have been used as lightweight alternatives to deep learned filters when computational efficiency or interpretability is required [43]. In reinforcement learning and continual learning, random projections help stabilize representation drift by maintaining a fixed feature subspace [37]. Moreover, recent theoretical works [20, 32] have demonstrated that the behavior of gradient-trained wide networks can often be approximated by their random-feature counterparts in the so-called *neural tangent kernel* (NTK) regime.

The appeal of random neural networks lies in this balance between expressiveness and tractability. They are expressive enough to model nonlinear relationships, yet simple enough to admit precise

theoretical characterizations. In the context of this thesis, this property is essential: by freezing the random features and focusing on the analytical form of the ridge regression solution (2.2), we can study high-dimensional learning dynamics directly, derive deterministic equivalents, and explore new data-driven regularization strategies such as the derivative-based calibration proposed in later chapters.

2.3. Regression in a Neural Networks Context

The connection between regression and neural networks is more than superficial: at its core, a feedforward neural network implements a composition of affine transformations and nonlinearities, followed by a linear readout. When the nonlinear mapping is fixed or random, the learning task in the final layer reduces to a regression problem on transformed features. This observation underpins both theoretical analyses of neural networks and practical algorithms such as the extreme learning machine and random feature regression used throughout this thesis.

2.3.1. Regression-Neural Network Problem Equivalence

Consider a single-hidden-layer neural network with input $\mathbf{x} \in \mathbb{R}^p$, hidden weights $\mathbf{W} \in \mathbb{R}^{N \times p}$, activation function $\sigma(\cdot)$, and output weights $\boldsymbol{\beta} \in \mathbb{R}^N$. The network output for a given input is

$$f(\mathbf{x}) = \boldsymbol{\beta}^{\top} \sigma(\mathbf{W}\mathbf{x}). \tag{2.4}$$

Given a dataset $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, stacking all activations into a matrix

$$\Sigma = \begin{bmatrix} \sigma(\mathbf{W}\mathbf{x}_1)^{\top} \\ \sigma(\mathbf{W}\mathbf{x}_2)^{\top} \\ \vdots \\ \sigma(\mathbf{W}\mathbf{x}_n)^{\top} \end{bmatrix} \in \mathbb{R}^{n \times N},$$
(2.5)

the empirical training objective under a squared-loss criterion becomes

$$\min_{\boldsymbol{\beta}} \frac{1}{n} \| \Sigma \boldsymbol{\beta} - \mathbf{Y} \|_2^2, \tag{2.6}$$

where $\mathbf{Y} = [y_1, y_2, \dots, y_n]^{\top}$. If \mathbf{W} and σ are fixed, the optimization in (2.6) is identical to a linear regression on the nonlinear feature vectors $\phi(\mathbf{x}_i) = \sigma(\mathbf{W}\mathbf{x}_i)$. In this sense, the neural network acts as a nonlinear feature generator, while the training process in the final layer corresponds to ordinary least squares (OLS) estimation:

$$\hat{\boldsymbol{\beta}}_{\mathsf{OLS}} = (\boldsymbol{\Sigma}^{\mathsf{T}} \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{\mathsf{T}} \mathbf{Y}. \tag{2.7}$$

This equivalence highlights the dual view of single-layer networks as regressors operating in a random feature space. From a functional perspective, $f(\mathbf{x})$ approximates an underlying mapping $f^*(\mathbf{x})$ via linear combination of basis functions $\{\sigma(\mathbf{w}_i^{\top}\mathbf{x})\}_{i=1}^N$, analogously to classical regression models that approximate functions through polynomial or spline bases. When N is large, the random basis becomes highly expressive, and the regression solution (2.7) provides an efficient projection of the target function onto the subspace spanned by the random features.

2.3.2. Ridge Regression (Penalization)

In high-dimensional settings, where N (the number of features) is comparable to or larger than n (the number of samples), the matrix $\Sigma^{\top}\Sigma$ in (2.7) is often ill-conditioned or singular. This makes the OLS solution unstable and highly sensitive to small perturbations in the data. To mitigate this, a regularization term is introduced, leading to the *ridge regression* or *Tikhonov regularization* estimator:

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\boldsymbol{\Sigma}^{\top} \boldsymbol{\Sigma} + \gamma I_N)^{-1} \boldsymbol{\Sigma}^{\top} \mathbf{Y}, \tag{2.8}$$

where $\gamma>0$ is the regularization parameter controlling the trade-off between data fidelity and weight magnitude.

Ridge regression can be interpreted from several complementary perspectives:

- Statistical: it shrinks coefficient estimates toward zero, reducing variance at the cost of a small bias [16]. This stabilizes predictions in the presence of multicollinearity or limited data.
- **Bayesian:** it corresponds to a maximum a posteriori (MAP) estimator under a Gaussian prior $\beta \sim \mathcal{N}(0, \frac{1}{\alpha}I_N)$, linking regularization to prior belief on model complexity [46].
- Numerical: it improves the conditioning of $\Sigma^{\top}\Sigma$, ensuring invertibility and controlling the amplification of noise.

In the neural-network setting, ridge regularization plays a role analogous to weight decay or ℓ_2 -regularization in deep learning. The effect of γ on the spectrum of $(\Sigma^\top \Sigma)$ determines both the smoothness of the fitted function and the stability of the learned weights [12]. As $\gamma \to 0$, the estimator approaches the minimum-norm interpolator; as $\gamma \to \infty$, it collapses toward zero.

Formally, ridge regression defines a kernel machine with implicit kernel

$$k(\mathbf{x}, \mathbf{x}') = \frac{1}{N} \sigma(\mathbf{W}\mathbf{x})^{\top} (\sigma(\mathbf{W}\mathbf{x}')), \tag{2.9}$$

so that the predictor may be expressed equivalently as

$$f(\mathbf{x}) = \mathbf{k}(\mathbf{x})^{\top} (\mathbf{K} + \gamma I_n)^{-1} \mathbf{Y}, \tag{2.10}$$

where $\mathbf{K} = \Sigma \Sigma^{\top}/N$ is the empirical kernel matrix. This dual formulation makes the connection between ridge regression, kernel methods, and random-feature neural networks explicit, establishing a unified mathematical framework that later sections of this thesis leverage for high-dimensional analysis and regularization calibration.

2.4. High-Dimensional Statistics and Random Matrix Theory

The practical and theoretical challenges of modern machine learning often arise in the so-called *high-dimensional regime*, where the number of model parameters or features is comparable to, or even exceeds, the number of available samples. Classical statistical theory, which assumes a fixed feature dimension p and an increasing sample size p0, no longer provides reliable approximations in this regime. Instead, understanding the behavior of estimators and algorithms requires asymptotic frameworks that allow both p1 and p2 to grow together at a comparable rate.

2.4.1. High-Dimensional Regime

In high-dimensional statistics, quantities that were once negligible—such as correlations among features or small eigenvalues of sample covariance matrices—become dominant. A central parameter characterizing this setting is the *aspect ratio*

$$c = \frac{p}{n},\tag{2.11}$$

which measures the dimensionality of the problem relative to the number of observations. When $c\ll 1$, the system is overdetermined, and classical intuition holds: estimators are consistent, and sample covariance matrices approximate their population counterparts. However, when c approaches or exceeds unity ($c\geq 1$), the sample covariance matrix becomes singular, and naive estimators such as ordinary least squares (2.7) become ill-posed. This is precisely the regime where regularization methods like ridge regression become essential.

The study of this high-dimensional behavior has revealed a series of nonintuitive phenomena. For example, ridge regression and related estimators exhibit a *double descent* risk curve [3], where increasing model complexity initially worsens generalization (the classical bias–variance trade-off) but later improves it again once p>n. These effects stem from spectral properties of random design matrices and have motivated the development of asymptotic tools to characterize estimator performance when both p and p are large.

Modern theoretical analyses therefore treat the design matrix $X \in \mathbb{R}^{p \times n}$ as random, often assuming entries that are independent and identically distributed (i.i.d.) with zero mean and variance scaled as 1/n. This scaling ensures nontrivial limiting spectra as $n, p \to \infty$. The statistics of the eigenvalues of sample covariance matrices $S = \frac{1}{n}XX^{\top}$ form the cornerstone of this analysis and are governed by results from *random matrix theory* (RMT).

2.4.2. Random Matrix Theory Essentials

Random matrix theory provides a framework for studying the spectral behavior of large random matrices—such as covariance matrices, kernel matrices, or weight matrices in neural networks—when their dimensions grow without bound. Originating in nuclear physics with Wigner's work on energy spectra [50], RMT has since become a central tool in high-dimensional statistics [2] and modern machine learning [7].

One of the most fundamental results in this field is the *Marchenko–Pastur law* [31], which describes the limiting empirical distribution of eigenvalues of the sample covariance matrix $S=\frac{1}{n}XX^{\top}$. If the entries of X are i.i.d. with mean zero and variance 1/n, then as $p,n\to\infty$ with $p/n\to c$, the spectral density of S converges almost surely to

$$\rho(\lambda) = (1 - \frac{1}{c})_{+} \delta(\lambda) + \frac{\sqrt{(\lambda_{+} - \lambda)(\lambda - \lambda_{-})}}{2\pi c\lambda}, \qquad \lambda \in [\lambda_{-}, \lambda_{+}],$$
 (2.12)

where $\lambda_{\pm}=(1\pm\sqrt{c})^2$ define the spectral support and $(x)_{+}=\max(0,x)$. This result implies that, even in the absence of structure, the eigenvalues of high-dimensional covariance matrices exhibit systematic, nontrivial distributions that deviate from their population counterparts.

RMT also provides tools such as the Stieltjes transform, defined as

$$m(z) = \int \frac{1}{\lambda - z} d\rho(\lambda), \qquad (2.13)$$

which serves as a key analytical function for characterizing the spectral properties of random matrices and for deriving deterministic equivalents in high-dimensional limits. These equivalents approximate random quantities—such as the generalization error or the trace of resolvents—with deterministic limits as $n, p \to \infty$ [5].

In the context of ridge regression and random-feature networks, random matrix theory enables closed-form asymptotic expressions for training and test errors [10, 28, 13]. These results make it possible to predict model behavior and regularization effects without Monte Carlo simulations, by solving fixed-point equations involving the spectral density (2.12) or its transform (2.13). Crucially, they reveal how phenomena such as double descent, bias—variance trade-offs, and stability depend on spectral properties of the random design matrix.

Overall, RMT provides the mathematical backbone for analyzing high-dimensional learning systems, including random-feature neural networks. It supplies the deterministic equivalents and asymptotic limits upon which the theoretical developments and calibration strategies of this thesis are built.

In this thesis, random matrix theory underpins not only the analysis of the ridge estimator but also the construction of the random-feature model itself. Specifically, the hidden-layer weight matrix **W** is drawn with i.i.d. Gaussian entries, a choice justified by RMT results showing that Gaussian ensembles yield well-defined and analytically tractable spectral distributions in the large-dimensional limit. This property

ensures that the empirical covariance of the random features $\Sigma = \sigma(\mathbf{W}\mathbf{X})$ concentrates around its deterministic equivalent, enabling precise asymptotic analysis. Consequently, the theoretical framework of RMT directly supports both the modeling assumptions and the derivation of the derivative-based calibration rule proposed in this thesis, as it allows one to characterize how the choice of γ and the feature dimensionality jointly affect stability and generalization in high-dimensional regimes. That is to be discussed in details later in Chapter 3.

2.5. Synthesis and Relevance to Our Work

The preceding sections have outlined the key developments that converge to form the conceptual foundation of this thesis. Neural networks have evolved from simple perceptrons to deep, highly expressive architectures capable of learning complex representations. Alongside this empirical success, simplified variants such as random-feature networks and extreme learning machines have emerged as tractable models that retain much of the representational richness of deep networks while admitting analytical treatment.

The random-weight paradigm provides a crucial simplification: by fixing the hidden-layer parameters and focusing learning on a single linear output layer, the training problem reduces to a convex regression task. This structure allows the use of closed-form estimators such as ridge regression, linking neural computation directly to classical statistical theory. The equivalence between neural networks and regression models, made explicit in (2.8), reveals that many stability and generalization properties of networks are ultimately governed by the spectral characteristics of the random feature matrix $\Sigma = \sigma(\mathbf{WX})$.

High-dimensional statistics and random matrix theory extend this connection further by providing the analytical machinery to characterize such spectral properties. Results such as the Marchenko–Pastur law and its deterministic equivalents describe how the eigenvalue distributions of large random matrices govern estimator variance, bias, and stability. Through this lens, neural networks—and random-feature models in particular—can be studied as high-dimensional statistical systems whose performance depends not only on their architecture or activation function, but also on the asymptotic interactions between sample size, feature dimensionality, and regularization strength.

This theoretical synthesis motivates the central focus of the present thesis. Traditional calibration methods for ridge-type estimators, such as cross-validation or Couillet's deterministic-equivalent approach [6], provide systematic ways to tune the regularization parameter γ . However, these methods rely either on computationally intensive resampling or on asymptotic approximations that require full knowledge of model parameters. By contrast, the derivative-based calibration rule developed in this work offers a fully data-driven alternative: it determines the optimal γ by analyzing the empirical derivative of the training loss with respect to the regularization parameter itself.

Viewed in this broader context, the proposed approach embodies the same philosophy that underlies much of modern high-dimensional learning: to design estimators that remain theoretically interpretable, computationally efficient, and empirically robust even when operating in regimes far beyond the reach of classical statistical assumptions. The subsequent chapters build upon this synthesis, developing both the analytical framework and the empirical validation that demonstrate how derivative-based calibration can serve as a practical and theoretically grounded regularization principle for random-feature neural networks.

Methodology

In the methodology section, the setting of the study will first be described, detailing the high-dimensional setting and the corresponding parameters arrangement involved, as well as the mathematical formulation of the single-layer perceptron model with random Gaussian weights. Following this, the problem inherent in the conventional approach to ridge regression will be articulated, with a particular focus on the discrepancies observed in the optimal regularization parameter γ found through minimizing traditional for this setting error metrics as opposed to different, direct error measures. Finally, potential solutions to address these discrepancies will be proposed and explored, with the aim of refining the ridge regression method to achieve more reliable and consistent results in high-dimensional contexts.

In recent decades, random neural networks have gained prominence as an effective solution to various challenges in machine learning. These networks are particularly advantageous in scenarios with limited training data, as they help mitigate the constraints of computational and memory resources. Furthermore, random neural networks serve as efficient random feature extractors, enhancing their utility across a range of applications. In this context, the application of ridge regression within these networks, is explored, aiming to leverage their strengths while addressing the unique challenges they present.

3.1. Model Setting

As was stated before, focus of this study lies within researching behaviour of ridge regression methods in application with relatively simple neural network setting in high dimensions, and the issues of the conventional approach caused by such an arrangement. In this section, model setting as well as the toolbox of the original, traditional approach to high-dimensional ridge regression is being discussed.

3.1.1. Ridge Regression With Random Neural Network

A structure of the single-hidden-layer neural network that is being used for the studies is depicted on the figure 3.1. 2

Hence, the setting is the following:

Given the input data matrix $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$, where each column $\mathbf{x}_i \in \mathbb{R}^p$ represents an input vector, the neural network processes this input through a randomly initialized weight matrix $\mathbf{W} \in \mathbb{R}^{N \times p}$. The entries of \mathbf{W} are independent and identically distributed (i.i.d.), following a standard Gaussian distribution.

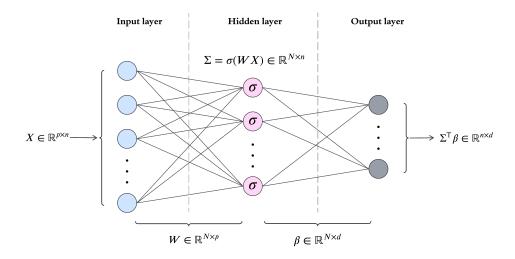


Figure 3.1: A diagram illustrating the structure of the single-layer random neural network.

The output of the first layer is denoted by $\Sigma = \sigma(\mathbf{W}\mathbf{X}) \in \mathbb{R}^{N \times n}$, where $\sigma : \mathbb{R} \to \mathbb{R}$ is some activation function applied entry-wise to the product $\mathbf{W}\mathbf{X}$. This results in a feature matrix Σ , where each column $\sigma(\mathbf{W}\mathbf{x}_i)$ represents the transformed feature vector for the input \mathbf{x}_i . These columns can be interpreted as random features (nonlinear random features in case function σ is nonlinear itself) of the input data.

The second layer of the neural network involves a weight matrix $\beta \in \mathbb{R}^{N \times d}$, which is learned to map the feature matrix Σ to the target output $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n] \in \mathbb{R}^{d \times n}$, where each column $\mathbf{y}_i \in \mathbb{R}^d$ represents the target vector associated with the input \mathbf{x}_i . The learning of β is performed by minimizing a regularized loss function.

Therefore, the overall architecture of the neural network involves transforming the input data through a random weight matrix, applying some activation function to obtain random nonlinear features, and then learning a weight matrix to map these features to the desired output. This setup leverages the randomness in the weights to extract useful features from the input data, while the learning process focuses on minimizing the prediction error through a regularized optimization problem.

The output weight matrix β is learned to minimize the regularized mean squared error:

$$L(\beta) = \frac{1}{n} \sum_{i=1}^{n} \|\mathbf{y}_i - \beta^{\top} \sigma(\mathbf{W} \mathbf{x}_i)\|^2 + \gamma \|\beta\|_F^2$$

The aim is to find the $\hat{\beta}$ that minimizes this loss. The loss function can be rewritten using matrix notation. The loss function can be written as:

$$L(\beta) = \frac{1}{n} \mathsf{Tr}((\mathbf{Y} - \beta^{\top} \mathbf{\Sigma})(\mathbf{Y} - \beta^{\top} \mathbf{\Sigma})^{\top}) + \gamma \mathsf{Tr}(\beta^{\top} \beta).$$

We then differentiate $L(\beta)$ with respect to β :

$$\frac{\partial L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \frac{\partial}{\partial \boldsymbol{\beta}} \left(\frac{1}{n} \mathsf{Tr}((\mathbf{Y} - \boldsymbol{\beta}^{\top} \boldsymbol{\Sigma}) (\mathbf{Y} - \boldsymbol{\beta}^{\top} \boldsymbol{\Sigma})^{\top}) + \gamma \mathsf{Tr}(\boldsymbol{\beta}^{\top} \boldsymbol{\beta}) \right).$$

Using matrix calculus, we have:

$$\frac{\partial}{\partial \beta} \left(\frac{1}{n} \mathsf{Tr}((\mathbf{Y} - \beta^{\top} \boldsymbol{\Sigma}) (\mathbf{Y} - \beta^{\top} \boldsymbol{\Sigma})^{\top}) \right) = \frac{2}{n} \boldsymbol{\Sigma} (\mathbf{Y} - \beta^{\top} \boldsymbol{\Sigma})^{\top},$$

and

$$\frac{\partial}{\partial \boldsymbol{\beta}} \left(\gamma \mathsf{Tr}(\boldsymbol{\beta}^{\top} \boldsymbol{\beta}) \right) = 2 \gamma \boldsymbol{\beta}.$$

Combining these results, we get:

$$\frac{\partial L(\beta)}{\partial \beta} = \frac{2}{n} \mathbf{\Sigma} (\mathbf{Y} - \beta^{\top} \mathbf{\Sigma})^{\top} + 2\gamma \beta = \frac{2}{n} \mathbf{\Sigma} (\mathbf{\Sigma}^{\top} \beta - \mathbf{Y}^{\top}) + 2\gamma \beta.$$

Setting the gradient to zero for minimization, we have:

$$0 = \frac{2}{n} \mathbf{\Sigma} (\mathbf{\Sigma}^{\top} \hat{\beta} - \mathbf{Y}^{\top}) + 2\gamma \hat{\beta}.$$

From there, rearranging the terms:

$$0 = \frac{2}{n} \mathbf{\Sigma} \mathbf{\Sigma}^{\top} \hat{\beta} + \frac{2}{n} \mathbf{\Sigma} \mathbf{Y}^{\top} + 2\gamma \hat{\beta}$$
$$0 = \left(\frac{1}{n} \mathbf{\Sigma} \mathbf{\Sigma}^{\top} + \gamma \mathbf{I}_{N}\right) \hat{\beta} + \frac{1}{n} \mathbf{\Sigma} \mathbf{Y}^{\top}$$
$$\left(\frac{1}{n} \mathbf{\Sigma} \mathbf{\Sigma}^{\top} + \gamma \mathbf{I}_{N}\right) \hat{\beta} = \frac{1}{n} \mathbf{\Sigma} \mathbf{Y}^{\top}$$

Thus we obtain the closed-form expression for the ridge-regressor:

$$\hat{\beta} = \frac{1}{n} \mathbf{\Sigma} \left(\frac{1}{n} \mathbf{\Sigma} \mathbf{\Sigma}^{\top} + \gamma \mathbf{I}_n \right)^{-1} \mathbf{Y}^{\top}.$$
 (3.1)

3.1.2. Error Approximation in High-Dimensional Setting: a Conventional Approach

In the realm of high-dimensional settings, traditional approaches to error estimation have been well-documented and provide foundational insights into the behavior of models. This conventional approach is detailed extensively in the book [6].

Indeed, it is not difficult to observe that the error terms can be computed directly from the data. The training mean squared error (MSE) on the given training set (X, Y) is given by:

$$E_{\mathsf{train}} = \frac{1}{n} \left\| \mathbf{Y}^{\top} - \mathbf{\Sigma}^{\top} \boldsymbol{\beta} \right\|_{F}^{2} = \frac{\gamma^{2}}{n} \mathsf{tr} \mathbf{Y} \mathbf{Q}^{2}(\gamma) \mathbf{Y}^{\top}, \tag{3.2}$$

Please note, that here β is as in the formula (3.1), whilst $\mathbf{Q}(\gamma)$ is the resolvent of $\frac{1}{n}\mathbf{\Sigma}^{\top}\mathbf{\Sigma}$, defined as:

$$\mathbf{Q}(\gamma) \equiv \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_n\right)^{-1}.$$
 (3.3)

Similarly, the test MSE on a test set $(\hat{\mathbf{X}}, \hat{\mathbf{Y}}) \in \mathbb{R}^{p \times \hat{n}} \times \mathbb{R}^{d \times \hat{n}}$ of test sample size \hat{n} is given by:

$$E_{\text{test}} = \frac{1}{\hat{n}} \left\| \hat{\mathbf{Y}}^{\top} - \hat{\mathbf{\Sigma}}^{\top} \boldsymbol{\beta} \right\|_{F}^{2} = \frac{1}{\hat{n}} \text{tr} \hat{\mathbf{Y}} \hat{\mathbf{Y}}^{\top} - \frac{2}{n\hat{n}} \text{tr} \mathbf{Y} \mathbf{Q} \mathbf{\Sigma}^{\top} \hat{\mathbf{\Sigma}} \hat{\mathbf{Y}}^{\top} + \frac{1}{n^{2}\hat{n}} \text{tr} \mathbf{Y} \mathbf{Q} \mathbf{\Sigma}^{\top} \hat{\mathbf{\Sigma}} \hat{\mathbf{\Sigma}}^{\top} \mathbf{\Sigma} \mathbf{Q} \mathbf{Y}^{\top}$$
(3.4)

where $\hat{\Sigma} = \sigma(\mathbf{W}\hat{\mathbf{X}})$.

The traditional approach, as mentioned before, opts for effectively estimating these error 3.2 & 3.4 rather than using the original expressions. The question arises: why estimate errors when direct computation of E_{train} and E_{test} is possible? The challenge lies in the high-dimensional regime, where parameters $p,\,n,$ and N (dimensionality of the data, size of the dataset and number of neurons in the hidden layer, respectively) grow large simultaneously. Under these conditions, computing E_{train} , for example, becomes computationally intensive due to the inversion and the inherent randomness of Σ through W.

The randomness, while not a direct computational hurdle, complicates the assessment of the asymptotic behavior of E_{train} . Small changes in input can lead to significant variations in output, complicating the understanding of error behavior. To address this, the resolvent $\hat{\mathbf{Q}}$ is employed to approximate the asymptotic behavior of E_{train} . This method leverages the resolvent to simplify the complexity associated with the randomness of \mathbf{W} . This logic is applicable to E_{test} as well, providing a comprehensive framework for error approximation in high-dimensional settings.

The approximations \hat{E}_{train} and \hat{E}_{test} are based on the assumption that the matrix \mathbf{W} is sub-Gaussian and that the function σ is Lipschitz continuous [29]. These assumptions are necessary to apply concentration inequalities and derive deterministic equivalents for the resolvent and the error terms.

Assumptions and Lemmas To rigorously justify the use of \hat{E}_{train} and \hat{E}_{test} , we rely on several key results from random matrix theory.

Assumption 1. (Sub-Gaussian W): The matrix $\mathbf W$ is defined as $\mathbf W = \phi(\tilde{\mathbf W})$, where $\tilde{\mathbf W}$ has i.i.d. $\mathcal N(0,1)$ entries and $\phi(\cdot)$ is λ_ϕ -Lipschitz. For $a=\phi(b)\in\mathbb R$, ϕ maps $b\sim\mathcal N(0,I_p)$ to $a\sim\mathcal N_\phi(0,I_p)$.

Assumption 2. (Function σ): The function σ is Lipschitz continuous with parameter λ_{σ} . This assumption holds for many activation functions used in neural networks, such as the rectified linear unit (ReLU).

Assumption 3. (Growth Rate): As $n \to \infty$,

$$0 < \liminf_{n \to \infty} \min \left\{ \frac{p}{n}, \frac{N}{n} \right\} \le \limsup_{n \to \infty} \max \left\{ \frac{p}{n}, \frac{N}{n} \right\} < \infty$$

while $\gamma, \lambda_{\sigma}, \lambda_{\phi} > 0$ and d are kept constant. Additionally,

$$\limsup_{n\to\infty} \|\mathbf{X}\| < \infty, \quad \limsup_{n\to\infty} \max_{ij} |\mathbf{Y}_{ij}| < \infty.$$

Under these assumptions, we can state the following key results:

Lemma 1. (Concentration of Quadratic Forms): Let Assumptions 1 and 2 hold. For $\mathbf{X} \in \mathbb{R}^{p \times n}$ and $w \sim \mathcal{N}_{\phi}(0, I_p)$, define the random vector $\sigma \equiv \sigma(w^{\top}\mathbf{X})$. For $\mathbf{A} \in \mathbb{R}^{n \times n}$ with $\|\mathbf{A}\| \leq 1$,

$$P\left(\left|\frac{1}{N}\boldsymbol{\sigma}^{\top}\mathbf{A}\boldsymbol{\sigma} - \frac{1}{N}\mathsf{tr}(\boldsymbol{\Phi}\mathbf{A})\right| > t\right) \leq Ce^{-cN\min\left(\frac{t^2}{t_0^2},t\right)}$$

for some constants C, c > 0, where $\Phi = \mathbb{E}[\sigma(w^{\top}\mathbf{X})\sigma(w^{\top}\mathbf{X})^{\top}]$ [29].

This lemma ensures that the quadratic forms involving $\sigma(w^{\top}\mathbf{X})$ concentrate around their expectations, allowing us to replace random terms with deterministic equivalents.

Theorem 3.1.2.1. (Asymptotic Equivalent for $\mathbb{E}[\mathbb{Q}]$): Let Assumptions 1-3 hold and define \mathbb{Q} as

$$\bar{\mathbf{Q}} = \left(\frac{N}{n} \frac{\mathbf{\Phi}}{1+\delta} + \gamma I_n\right)^{-1},\tag{3.5}$$

where δ is the unique positive solution to $\delta = \frac{1}{N} \text{tr}(\Phi \bar{\mathbf{Q}})$. Then, for all $\epsilon > 0$, there exists c > 0 such that

$$\|\mathbb{E}[\mathbf{Q}] - \bar{\mathbf{Q}}\| \le cn^{-1/2 + \epsilon}$$

[29].

This theorem provides an asymptotic equivalent for the expectation of the resolvent \mathbf{Q} , allowing us to approximate $\mathbb{E}[\mathbf{Q}]$ with $\bar{\mathbf{Q}}$ in high-dimensional settings.

Using these results, we can derive the asymptotic approximations for the training and test mean squared errors:

Theorem 3.1.2.2. (Asymptotic Training Mean-Square Error): Let Assumptions 1-3 hold. Then, for all $\epsilon > 0$,

$$n^{1/2-\epsilon}(E_{\mathsf{train}} - \hat{E}_{\mathsf{train}}) \to 0$$
 almost surely,

where

$$E_{\text{train}} = \frac{1}{N} \left\| \mathbf{Y}^{\top} - \mathbf{\Sigma}^{\top} \hat{\beta} \right\|_{F}^{2} = \frac{\gamma^{2}}{N} \text{tr} \mathbf{Y} \mathbf{Q}^{2}(\gamma) \mathbf{Y}^{\top},$$

$$\hat{E}_{\text{train}} = \frac{\gamma^{2}}{N} \text{tr} \left(\mathbf{Y}^{\top} \mathbf{Y} \bar{\mathbf{Q}} \left(\frac{1}{n} \text{tr} (\mathbf{\Phi} \bar{\mathbf{Q}}^{2}) (I_{N} - \frac{1}{n} \text{tr} (\mathbf{\Phi}^{2} \bar{\mathbf{Q}}^{2}))^{-1} \right) \bar{\mathbf{Q}} \right)$$
(3.6)

[29].

Similarly, for the test mean squared error:

Conjecture 3.1.2.1. Deterministic Equivalent for E_{test} : Let Assumptions 1-2 hold. Then, for all $\epsilon > 0$,

$$n^{1/2-\epsilon}(E_{\textit{test}} - \hat{E}_{\textit{test}}) o 0$$
 almost surely,

where

$$E_{test} = \frac{1}{\hat{n}} \left\| \hat{\mathbf{Y}}^{\top} - \hat{\mathbf{\Sigma}}^{\top} \hat{\beta} \right\|_{F}^{2},$$

$$\hat{E}_{test} = \frac{1}{\hat{n}} \left\| \hat{\mathbf{Y}}^{\top} - \hat{\mathbf{\Phi}}^{\top} \bar{\mathbf{Q}} \mathbf{Y}^{\top} \right\|_{F}^{2} + \frac{1}{n} tr(\mathbf{Y}^{\top} \mathbf{Y} \bar{\mathbf{Q}} \mathbf{\Phi} \bar{\mathbf{Q}}) \left(1 - \frac{1}{n} tr(\mathbf{\Phi} \bar{\mathbf{Q}}) \right)^{-1}$$
(3.7)

[29].

Kernel-Based Reformulation of Asymptotic Errors Based on our previous considerations, we can express the results in terms of the limiting kernel ${\bf K}$, which provides a deterministic framework for understanding the behavior of the model. The kernel ${\bf K}$ effectively captures the feature mapping induced by the activation function σ and is pivotal in characterizing the asymptotic properties of the model's performance. This approach leverages **Theorem 3.1.2.1** and **Conjecture 3.1.2.1**, which demonstrate how the effective kernel $\bar{{\bf K}}$ influences the training and test errors.

To compute the effective kernel $\bar{\mathbf{K}}$, we employ the fixed-point equation derived from **Theorem 3.1.2.1**, which is iteratively solved to obtain:

$$\bar{\mathbf{K}} = \frac{N}{n} \frac{\mathbf{K}}{1+\delta},\tag{3.8}$$

where δ is defined as:

$$\delta = \frac{1}{N} \text{tr}(\mathbf{K}\bar{\mathbf{Q}}).$$

Derivation of Asymptotic Errors with {\bf K} The asymptotic training mean squared error E_{train} is given by:

$$E_{\mathsf{train}} = \frac{1}{N} \left\| \mathbf{Y}^{\top} - \mathbf{\Sigma}^{\top} \boldsymbol{\beta} \right\|_{F}^{2} = \frac{\gamma^{2}}{N} \mathsf{tr}(\mathbf{Y} \mathbf{Q}^{2} \mathbf{Y}^{\top}), \tag{3.9}$$

which can be rewritten using the kernel as:

$$\bar{E}_{\mathsf{train}} = \frac{\gamma^2}{n} \mathsf{tr} \left(\mathbf{Y} \bar{\mathbf{Q}} \left(\frac{\frac{1}{N} \mathsf{tr}(\bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}})}{1 - \frac{1}{N} \mathsf{tr}(\bar{\mathbf{K}} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}})} \bar{\mathbf{K}} + \mathbf{I}_n \right) \bar{\mathbf{Q}} \mathbf{Y}^{\top} \right)$$
(3.10)

where:

$$\bar{\mathbf{Q}} = \left(\frac{N}{n} \frac{\bar{\mathbf{K}}}{1+\delta} + \gamma \mathbf{I}_n\right)^{-1}.$$

Similarly, the test mean squared error E_{test} is expressed as:

$$E_{\text{test}} = \frac{1}{\hat{n}} \left\| \hat{\mathbf{Y}}^{\top} - \hat{\mathbf{\Sigma}}^{\top} \boldsymbol{\beta} \right\|_{F}^{2}, \tag{3.11}$$

and can be reformulated in terms of the kernel:

$$\bar{E}_{\mathsf{test}} = \frac{1}{\hat{n}} \left\| \hat{\mathbf{Y}}^{\top} - \bar{\mathbf{K}}_{\mathbf{X}\hat{\mathbf{X}}}^{\top} \bar{\mathbf{Q}} \mathbf{Y}^{\top} \right\|_{F}^{2} + \frac{\frac{1}{N} \mathsf{tr} \left(\mathbf{Y} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}} \mathbf{Y}^{\top} \right)}{1 - \frac{1}{N} \mathsf{tr} \bar{\mathbf{K}} \bar{\mathbf{Q}} \bar{\mathbf{K}} \bar{\mathbf{Q}}} \left(\frac{1}{\hat{n}} \mathsf{tr} \bar{\mathbf{K}}_{\hat{\mathbf{X}}\hat{\mathbf{X}}} - \frac{1}{\hat{n}} \mathsf{tr} (\mathbf{I}_{n} + \gamma \bar{\mathbf{Q}}) (\bar{\mathbf{K}}_{\mathbf{X}\hat{\mathbf{X}}}^{\top} \bar{\mathbf{K}}_{\mathbf{X}\hat{\mathbf{X}}} \bar{\mathbf{Q}}) \right)$$
(3.12)

where the effective kernels are defined as:

$$\mathbf{K}_{\mathbf{X}\mathbf{X}} = \mathbb{E}[\sigma(\mathbf{X}^{\top}\mathbf{w})\sigma(\mathbf{w}^{\top}\mathbf{X})],$$

$$\mathbf{K}_{\mathbf{X}\hat{\mathbf{X}}} = \mathbb{E}[\sigma(\mathbf{X}^{\top}\mathbf{w})\sigma(\mathbf{w}^{\top}\hat{\mathbf{X}})],$$

$$\mathbf{K}_{\hat{\mathbf{X}}\hat{\mathbf{X}}} = \mathbb{E}[\sigma(\hat{\mathbf{X}}^{\top}\mathbf{w})\sigma(\mathbf{w}^{\top}\hat{\mathbf{X}})],$$

and their normalized counterparts:

$$\bar{\mathbf{K}} = \frac{N}{n} \frac{\mathbf{K}}{1+\delta},$$

$$\bar{\mathbf{K}}_{\mathbf{X}\hat{\mathbf{X}}} = \frac{N}{n} \frac{\mathbf{K}_{\mathbf{X}\hat{\mathbf{X}}}}{1+\delta},$$

$$\bar{\mathbf{K}}_{\hat{\mathbf{X}}\hat{\mathbf{X}}} = \frac{N}{n} \frac{\mathbf{K}_{\hat{\mathbf{X}}\hat{\mathbf{X}}}}{1+\delta}.$$

The kernel ${\bf K}$ replaces the need for direct computation of the matrix ${\bf \Sigma}$ by leveraging deterministic equivalents, which simplify the complexity of high-dimensional analysis. This methodology aligns with the results derived in [6], where the effective kernel $\bar{{\bf K}}$ is shown to have a significant impact on regression performance, providing a reliable approximation of the model's behavior in the asymptotic regime. By applying these insights, we can gain a deeper understanding of the model's performance and address potential issues such as overfitting by examining the interaction between p/n, N/n, and the choice of activation functions.

In summary, these results provide a rigorous foundation for approximating the training and test mean squared errors in high-dimensional settings. By using the deterministic equivalents \hat{E}_{train} and \hat{E}_{test} , we can effectively handle the complexity introduced by the randomness of Σ and gain insights into the asymptotic behavior of the error terms [29].

3.1.3. Variability and Extensions of the Model

In this subsection different parameters and model arrangement are going to be considered in order to explore the possible variability of the model. Exploring various aspects of the neural network model

is crucial for understanding its performance and behavior. Modifying different parameters and functions, such as activation functions, kernel choices, and ratios like $\frac{p}{n}$ and $\frac{N}{n}$, can significantly impact the model's ability to generalize and avoid overfitting. Moreover, it will be later shown in the results section that our analysis leads to different results, which vary substantially with changes in the parameters.

Different Activation Functions and Kernels Activation functions play a pivotal role in defining the nonlinear transformations within the neural network. Commonly used activation functions include ReLU, sigmoid, tanh, and others, each impacting the network's behavior differently. These functions lead to different kernel functions, influencing how data is mapped and processed. For instance, the ReLU activation function corresponds to a specific kernel, while other functions like sigmoid and tanh result in distinct kernels.

Different activation functions lead to different kernels, impacting the behavior and performance of the neural network. For a given dataset \mathbf{X} , it is possible to compute the "limiting" kernel \mathbf{K} for the listed activation functions $\sigma(\cdot)$ using theoretical results. By iterating the fixed-point equation presented in Theorem 3.1.2.1, one can derive the effective kernel $\mathbf{K} \equiv \frac{N}{n} \frac{\mathbf{K}}{1+\delta}$ in practical settings where n,p,N are large. This effective kernel provides insights into the regression performance and is advantageous as it applies to deterministic input data \mathbf{X} rather than relying on randomly modeled data. Table 3.1 from [6] lists the limiting kernels for various activation functions, highlighting the diverse impact these functions have on the model:

$\sigma(t)$	$\kappa(\mathbf{x},\mathbf{y})$
t	$\mathbf{x}^{\top}\mathbf{y}$
t	$\frac{2}{\pi} \ \mathbf{x}\ \cdot \ \mathbf{y}\ \left(\mathcal{L} \cdot \arcsin(\mathcal{L}) + \sqrt{1 - \mathcal{L}^2}\right)$
$ReLU(t) \equiv \max(t,0)$	$\frac{1}{2\pi} \ \mathbf{x}\ \cdot \ \mathbf{y}\ \left(\mathcal{L} \cdot \arccos(-\mathcal{L}) + \sqrt{1 - \mathcal{L}^2} \right)$
$a_{+} \max(t,0) + a_{-} \max(-t,0)$	$\frac{1}{2}(a_+^2 + a^2)\mathbf{x}^{T}\mathbf{y} + \frac{1}{2\pi}\ \mathbf{x}\ \cdot \ \mathbf{y}\ \left((a_+ + a)^2 \left(-\mathcal{L} \cdot \arccos(\mathcal{L}) + \sqrt{1 - \mathcal{L}^2} \right) \right)$
$a_2t^2 + a_1t + a_0$	$a_2^2 \left(2(\mathbf{x}^\top \mathbf{y})^2 + \ \mathbf{x}\ ^2 \ \mathbf{y}\ ^2 \right) + a_1^2 \mathbf{x}^\top \mathbf{y} + a_2 a_0 \left(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2 \right) + a_0^2$
erf(t)	$\frac{2}{\pi}\arcsin\left(\frac{2\mathbf{x}^{\top}\mathbf{y}}{\sqrt{(1+2\ \mathbf{x}\ ^2)(1+2\ \mathbf{y}\ ^2)}}\right)$
$1_{t>0}$	$\frac{1}{2} - \frac{1}{2\pi} \arccos(\mathcal{L})$
sign(t)	$\frac{2}{\pi}\arcsin(\mathcal{L})$
$\cos(t)$	$\exp\left(-\frac{1}{2}(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2)\right)\cosh(\mathbf{x}^{\top}\mathbf{y})$
$\sin(t)$	$\exp\left(-\frac{1}{2}(\ \mathbf{x}\ ^2 + \ \mathbf{y}\ ^2)\right)\sinh(\mathbf{x}^\top\mathbf{y})$
$\exp(-t^2/2)$	$\frac{1}{\sqrt{(1+\ \mathbf{x}\ ^2)(1+\ \mathbf{y}\ ^2)-(\mathbf{x}^{\top}\mathbf{y})^2}}$

Table 3.1: Limiting kernel $\kappa(\mathbf{x}, \mathbf{y}) = \mathbb{E}[\sigma(\mathbf{w}^{\top}\mathbf{x})\sigma(\mathbf{w}^{\top}\mathbf{y})]$ for standard Gaussian \mathbf{w} , with $\mathcal{L} \equiv \frac{\mathbf{x}^{\top}\mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|}$.

Variability with Increasing Layers and Loss Functions Beyond activation functions and kernel selections, the variability of the model's architecture itself, such as the number of layers, plays a crucial role in determining performance. Adding more layers can introduce additional complexity and depth, allowing the network to capture intricate patterns and interactions within the data. However, this complexity must be balanced with the risk of overfitting, particularly in high-dimensional settings where the ratio p/n and N/n can exacerbate model instability.

In the upcoming results section, we will present experiments that demonstrate the impact of different architectural configurations on the model's performance. This will include analyses of varying the number of layers and other hyperparameters to provide insight into how such modifications can be optimally tuned for specific tasks.

3.2. Experimental Setup

For the initial testing, we consider a one-hidden-layer neural network with Gaussian weights and a linear activation function for simplicity. The goal is to study the behavior of this network in relation to the training and test errors, particularly when regularizing the model using the ridge regression penalty.

In our setup, the input matrix X (both for training and test sets) is generated as a Gaussian random matrix with independent standard normal entries. We use the following specific steps for data generation:

- The weight matrix \mathbf{W} is Gaussian, meaning each entry of $\mathbf{W} \in \mathbb{R}^{N \times p}$ is drawn from $\mathcal{N}(0,1)$, where N represents the number of neurons in the hidden layer and p is the input dimensionality.
- The target vector b is generated randomly with binary entries, where each entry takes a value of 0 or 1 according to a certain density parameter. For instance, a density of 90% means approximately 90% of the entries of b are set to 1 and the rest to 0.
- The output vector \mathbf{Y} is generated according to the relationship $\mathbf{Y} = \mathbf{\Sigma}^{\top} \mathbf{b} + \epsilon$, where $\mathbf{\Sigma} = \sigma(\mathbf{W}\mathbf{X})$ is the activation matrix, \mathbf{X} represents the input matrix, and ϵ is a Gaussian noise term with a certain noise level.

For the ridge regression model, we aim to find the optimal regularization parameter γ that minimizes both the training and test errors. To estimate the training error E_{train} , we calculate it as the squared difference between the observed training outputs $\mathbf{Y}_{\text{train}}$ and the predictions of the model. Similarly, the test error E_{test} is computed on the unseen test data. Additionally, deterministic equivalents of these errors, \hat{E}_{train} and \hat{E}_{test} , are obtained using approximations from random matrix theory, which simplify the analysis of the model's behavior.

For this setup, we use the following values for the key parameters:

- p/n ratio (a parameter under investigation), denoted as c_1 , where p is the input dimension and n is the number of training data points.
- N/n ratio (another parameter under investigation), denoted as c_2 , where N is the number of neurons in the hidden layer.
- A range of values for n (number of data points) for the initial experiment taken as values from 100 to 3000, increasing in increments of 200.
- A binary target vector b generated with a certain density value (parameter under investigation).
- A range of regularization parameters γ explored, from 10^{-7} to 10^2 .

We explore these settings using both a Monte Carlo method (random-based computations of E_{train} and E_{test}) and deterministic equivalents based on approximations for the kernel matrix \mathbf{K} .

The choice of the regularization parameter γ was made by examining the range in which the most interesting behavior occurs. As shown in the next section (Figure X.X), the test and training errors exhibit significant changes in behavior when γ is between 10^{-7} and 10^2 . This range captures the transition where the errors are minimized and gives us insights into the model's performance. Consequently, we focus on this interval to explore the optimal value of γ .

Regarding the choice of data point values n, we selected the range of 100 to 3000 to fit within the computational resources available on the Delft hypercomputer cluster. Larger values of n were not feasible due to resource limitations, and this range was found to be sufficient to observe the desired trends in the model's behavior.

3.3. Algorithm Optimization

This section describes how we transform a straightforward ("naive") implementation of the pipeline into an optimized one that is both faster and numerically more stable. We assume the reader is familiar with the experimental setup and notation from the preceding sections: given standardized data $(X_{\rm tr}, X_{\rm te})$, labels $y_{\rm tr}, y_{\rm te} \in \{-1, +1\}$, a random-features matrix $W \in \mathbb{R}^{N \times p}$, and an activation $\sigma(\cdot)$, we evaluate two one–dimensional calibration rules for the ridge penalty γ -Coulliet's deterministic-equivalent test proxy $\overline{E}_{\rm test}(\gamma)$ (to be discussed in later chapters, cf. Eq. (3.12)) and our direct-loss rule based on $\partial_{\gamma} \mathcal{L}_{\rm direct}^{\rm (RSS)}(\gamma)$ (cf. Eq. (4.3b))—and then form $\widehat{\beta}(\gamma)$ via the closed-form ridge expression (Eq. (3.1)).

3.3.1. Naive baseline: what is expensive

In a direct implementation, for *each* grid value γ one typically:

- 1. Recomputes kernel blocks $K_{tr} = K(X_{tr}, X_{tr}), K_{xX} = K(X_{tr}, X_{te}), K_{XX} = K(X_{te}, X_{te});$
- 2. Solves the fixed point for $\delta = \delta(\gamma)$ by repeatedly forming and inverting

$$Q(\gamma, \delta) = \frac{N}{n} \frac{K_{\text{tr}}}{1+\delta} + \gamma I_n;$$

- 3. Evaluates $\overline{E}_{\mathrm{test}}(\gamma)$ and/or $\partial_{\gamma}\mathcal{L}_{\mathrm{direct}}^{\mathrm{(RSS)}}(\gamma)$ using dense $n \times n$ linear algebra;
- 4. Forms the ridge solution $\widehat{\beta}(\gamma) = \frac{\Sigma}{n} \left(\frac{1}{n} \Sigma^{\top} \Sigma + \gamma I \right)^{-1} y$, where $\Sigma = \sigma(WX_{\rm tr})$.

Steps (2)-(4) perform an $n \times n$ factorization/inversion $per \gamma$. In perturbation studies, repeating the same work per scenario quickly dominates runtime and accumulates numerical error.

3.3.2. Optimized pipeline: what we change

We refactor the computation so that all operations that do not depend on γ (and, where possible, that do not depend on the scenario) are computed once and then reused. The concrete changes are:

(O1) Precompute kernels once and reuse. Compute the Gram blocks

$$K_{\rm tr}, K_{xX}, K_{XX}$$

once per dataset (or once per scenario when the training inputs are perturbed). The arc–cosine kernel implementation is numerically stabilized via norm products clamped away from 0 and cosine arguments clipped to (-1,1).

(O2) Eigendecompose $K_{\rm tr}$ once; diagonalize the fixed–point map. Write $K_{\rm tr} = V\Lambda V^{\top}$ with $\Lambda = {\rm diag}(\lambda_i)$. The fixed–point iteration

$$\delta \mapsto \frac{1}{n} \sum_{i=1}^{n} \frac{\lambda_i}{a(\delta)\lambda_i + \gamma}, \qquad a(\delta) = \frac{N/n}{1+\delta},$$

no longer requires forming Q^{-1} ; it becomes a cheap diagonal update using only $\{\lambda_i\}$. All $\overline{E}_{\text{test}}(\gamma)$ terms that involve traces with $Q(\gamma,\delta)^{-1}$ reduce to sums over $d_i(\gamma,\delta)=a(\delta)\lambda_i+\gamma$.

(O3) SVD trick for random features; avoid matrix inverses. Let $A = \Sigma/\sqrt{n}$ with thin SVD $A = USV^{\top}$, $S = \operatorname{diag}(s_i)$ (i = 1, ..., r). Then

$$\left(\frac{1}{n}\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma} + \gamma\boldsymbol{I}\right)^{-1} \; = \; \boldsymbol{V} \operatorname{diag}\!\!\left(\frac{1}{s_i^2 + \gamma}\right) \! \boldsymbol{V}^{\top} \; + \; \frac{1}{\gamma} P_{\mathrm{null}},$$

so that both $\widehat{\beta}(\gamma)$ and the direct derivative

$$\partial_{\gamma} \mathcal{L}_{\text{direct}}^{(\text{RSS})}(\gamma) = \gamma \left(\sum_{i=1}^{r} \frac{(v_i^{\top} y)^2}{(s_i^2 + \gamma)^3} + \frac{\|y - P_{\text{span}} y\|_2^2}{\gamma^3} \right) - \widehat{\sigma}_{\text{RSS}}^2(\gamma) \left(\sum_{i=1}^{r} \frac{1}{(s_i^2 + \gamma)^2} + \frac{n-r}{\gamma^2} \right)$$

are computed by *elementwise* operations on $\{s_i^2\}$. The RSS variance is computed consistently from the same SVD:

$$\widehat{\sigma}_{\text{RSS}}^{2}(\gamma) = \frac{1}{n} \left\| \frac{\gamma}{S^{2} + \gamma} V^{\top} y \right\|_{2}^{2} + \frac{1}{n} \|y - P_{\text{span}} y\|_{2}^{2}.$$

- **(O4) Factor once, solve many.** Cholesky/triangular solves replace explicit inverses everywhere (when we operate in the $n \times n$ domain). In the RKHS route, all per- γ work lives on the diagonal $\{d_i\}$; in the random features route, all per- γ work lives on $\{s_i^2\}$.
- (O5) Grid evaluation by broadcasting. We evaluate $\overline{E}_{\mathrm{test}}(\gamma)$ and $\partial_{\gamma}\mathcal{L}_{\mathrm{direct}}^{(\mathrm{RSS})}(\gamma)$ on the full log–grid using vectorized broadcasts over $d_i(\gamma,\delta)$ or $s_i^2+\gamma$, thereby avoiding Python loops.
- (O6) Caching across scenarios. We keep W fixed and reuse all W-dependent quantities (e.g. shapes, normalization). For Coulliet's route, only $(K_{\rm tr},K_{xX},K_{XX})$ and their eigenpairs change when training inputs are perturbed.

3.3.3. Complexity at a glance

Let n be the train size, $|\mathcal{G}|$ the grid length, and $r = \text{rank}(A) \le \min\{N, n\}$.

- Naive: $O(|\mathcal{G}| n^3)$ for repeated inversions in the δ loop and in $\beta(\gamma)$ per grid point, plus kernel recomputation per γ .
- · Optimized:
 - One eigendecomposition of K_{tr} : $O(n^3)$.
 - One SVD of A: $O(nr^2)$ or $O(Nr^2)$ (economy SVD).
 - Per γ : O(n) for Coulliet (diagonal arithmetic on $\{d_i\}$), O(r) for Direct (on $\{s_i^2\}$).
 - Per scenario: recompute kernels and one eigendecomposition of $K_{\rm tr}$; reuse all grid–wise vectorization and SVD formulas.

In practice this reduces wall-clock time by one to two orders of magnitude for the grids and scenario counts considered, while improving numerical stability by avoiding explicit inverses.

3.3.4. Practical notes

- Storage. We store baseline and per–scenario betas/gammas in compressed NPZs. Reusing ${\cal W}$ ensures fair, paired comparisons across scenarios.
- **Determinism.** A single seed controls scenario design and W; results reported in the main text do not depend on the particular seed choice in any qualitative way.

Takeaway. The optimized pipeline moves all γ -dependence to diagonal arithmetic on spectral quantities, eliminates explicit inverses, and restricts per scenario searches to local grids. This yields substantial speedups and improves numerical behavior, especially critical in the perturbation analysis, where many closely related problems must be solved repeatedly.

3.4. Problem Arising

To briefly recap: in the preceeding sections we have explicitly derived the expressions for E_{train} (3.2) and E_{test} (3.4), which represent training and test mean squared errors (MSE) respectively, in the high-dimensional regime. By using results from random matrix theory we formulated deterministic equivalents \hat{E}_{train} (3.10) and \hat{E}_{test} (3.12) for E_{train} and E_{test} . These equivalents allow for the approximation of these error metrics in scenarios where the number of data points n, the number of features p, and the number of neurons N are large. Approximations themselves depend on the resolvent $\hat{\mathbf{Q}}$ (3.5), which captures the effects of randomness inherent in the data through the Gram matrix and the kernel matrix. From there we have found an optimal value of hyperparameter γ that allows to optimize expressions for both \hat{E}_{train} (3.10) and \hat{E}_{test} (3.12).

However, a critical issue arises when comparing the optimal γ values obtained through this approach with those obtained through direct minimization of the ridge loss function, which in our case we dxefined as:

$$\mathcal{L}_{ridge} = \frac{\|\hat{\beta} - b\|_F^2}{\|b\|_F^2},\tag{3.13}$$

where b is the target vector, and $\hat{\beta}$ is the solution obtained from ridge regression: (3.1):

$$\hat{\beta} = \frac{1}{n} \mathbf{\Sigma} \left(\frac{1}{n} \mathbf{\Sigma} \mathbf{\Sigma}^{\top} + \gamma \mathbf{I}_n \right)^{-1} \mathbf{Y}^{\top},$$

This metric evaluates how well the predicted coefficients match the true target vector, with the goal of minimizing the relative Frobenius norm error.

The core of the problem lies in the discrepancy between the optimal γ values obtained by minimizing \hat{E}_{train} and \hat{E}_{test} (derived from random matrix theory approximations) and the optimal γ obtained by minimizing the ridge loss function directly. Ideally, one would expect that these values converge or align to some degree, as both approaches aim to minimize the model error. However, the experimental results suggest otherwise, showing significant differences in the γ values that yield the best performance for each case. We are going to see the details of this discrepancy in the results sections further.

3.4.1. Potential cause of the problem

One potential cause of this discrepancy could be rooted in the assumptions and approximations used in deriving \hat{E}_{train} and \hat{E}_{test} . The use of random matrix theory allows us to simplify the problem by averaging over the randomness in the data, which can smooth out some variability that is captured more directly in the ridge loss function. As a result, the deterministic equivalents may not fully capture all aspects of the data structure or noise present in real-world data.

Another contributing factor could be the simplified setup used for the target vector b, which is binary (0s and 1s) with a fixed density. While this setup is useful for initial theoretical testing, it may not reflect the complexity of more general real-world data, where target vectors can exhibit more variability and noise. The binary nature of b could limit the flexibility of the model, potentially influencing the optimization of γ in different ways when using the ridge loss function versus the random matrix theory approximations.

Finally, the inherent differences between the MSE-based training and test errors and the ridge loss function itself may also play a role. The ridge loss function focuses on minimizing the difference between the predicted coefficients $\hat{\beta}$ and the true target vector b, while the MSE errors focus on minimizing the prediction errors for the outputs. This distinction may lead to divergent results in terms of what constitutes the "optimal" value for γ .

4

Results

In the following chapter we will present the results of our experiments and findings regarding the relationship between the regularization parameter γ , Mean Squared Error , ridge loss function and various model parameters.

4.1. Research on the original Experiment

In this section we interrogate the simplest controllable setting in which our theory can be tested: data generated exactly according to the "synthetic law" $\mathbf{Y} = \mathbf{\Sigma}^{\top} \mathbf{b} + \varepsilon$ with a *known* ground-truth coefficient vector b. The benchmark serves three purposes. First, it allows us to *reproduce* the Gaussian-W baseline of Couillet & Liao [6] and recover the random-matrix-theory optimal regularisation γ_{RMT}^* . Second, by comparing this predictive optimum with the direct alignment loss $\mathcal{L}_{\text{ridge}}$, we expose the *causal gap*—the systematic misalignment between γ_{RMT}^* and the value that best recovers b (Section 4.1.4). Finally, sweeping three aspect ratios $p/n \in \{0.7, 1.0, 2.0\}$ and three sparsity levels of b (10 %, 50 %, 90 %), we chart the regimes where this gap widens or narrows, providing a roadmap for the more complex real-data experiments that follow.

4.1.1. Mean Squared Error against the regularization parameter γ

We will begin with replicating the initial experiment from the book [6], investigating the Mean Squared Error values for both E_{train} (3.9) , E_{test} (3.11) as well as \hat{E}_{train} (3.10) , \hat{E}_{test} (3.12) against a ridge regularization parameter γ . The goal is to observe the trends in the MSE for both the training and test sets as γ changes, and to identify the optimal γ that minimizes the test error and whether it is possible in general.

In order to illustrate the impact of different activation functions (linear, ReLU, and signum) and varying levels of the target vector density (15% versus 90%), we present in Figure 4.1 six separate plots. Each sub-figure shows the training and test MSE as a function of the ridge parameter γ . As expected, the linear, ReLU, and signum activations exhibit similar qualitative "U-shapes" in the test-error curves, though at different scales depending on density and nonlinearities. The graph presented here shows a characteristic behavior: the training error increases as γ increases, while the test error initially decreases, reaches a minimum, and then begins to increase again as γ becomes too large. This result aligns with the conclusion from the book that the hyperparameter γ can be fine-tuned to optimize test performance, given that n, p, and N are not too small and comparable in magnitude.

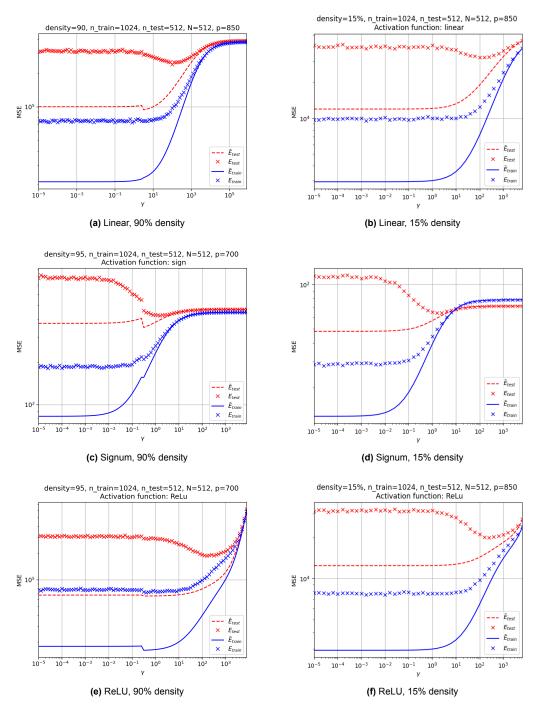


Figure 4.1: Training and test MSE versus the regularization parameter γ for three activation functions (linear, signum, ReLU) under two different densities of the target vector (90% and 15%). Each subfigure exhibits the typical U-shape in test error, alongside a rising training error as γ increases.

4.1.2. Optimal γ value and error behaviour with increasing amount of data

In this subsection, we investigate the behavior of the mean squared errors ($E_{\rm test}$ and $\hat{E}_{\rm test}$) and their corresponding optimal γ values as the number of data points n increases, while keeping the ratios p/n and N/n fixed. The parameters have been chosen such that p, N and n are sufficiently large and comparable in magnitude. This choice is motivated by the fact that in the high-dimensional regime, where the number of features p is approximately of the same order as the number of samples n, the theoretical results from random matrix theory hold.

The experiment is designed by keeping p/n and N/n fixed, while increasing n from 100 to 3500 in increments. For each value of n, 80% of the data is used for training and 20% for testing. At each step, we search for the optimal regularization parameter γ from a grid ranging from 10^{-7} to 10^1 (in increments of 0.1), which minimizes E_{test} and \hat{E}_{test} . These optimal γ values are then plotted against n, along with the corresponding values for the errors.

The following figures present the results of this experiment:

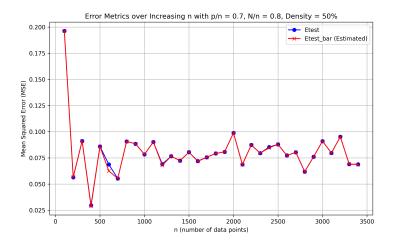


Figure 4.2: Error metrics (E_{test} and \hat{E}_{test}) over increasing n with p/n=0.7, N/n=0.8, and density = 50%.

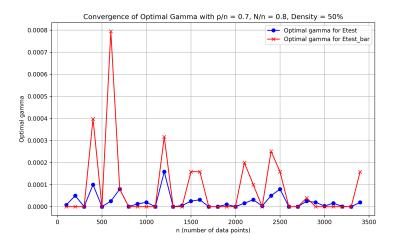


Figure 4.3: Convergence of optimal γ for E_{test} and \hat{E}_{test} over increasing n with p/n = 0.7, N/n = 0.8, and density = 50%.

From the results in Figures 4.2 and 4.3, we observe that the errors themselves remain mostly stable as the number of data points increases, with slight fluctuations that can be attributed to the random nature of the data generation. However, the optimal values for γ show more variability. While there is

some fluctuation, there is also a noticeable trend suggesting that the optimal γ values for E_{test} and \hat{E}_{test} show converging behaviour as n increases.

This empirical result aligns with the theoretical conjecture 3.1.2.1 shown earlier, which asserts that the difference between E_{test} and \hat{E}_{test} should diminish as n increases.

4.1.3. Optimal γ Values study for Different Loss Functions

In this subsection, we address a critical issue that arises in our investigation: the discrepancy between the optimal γ values needed to minimize different error metrics. The setting for this analysis remains consistent with the previous sections, where we explore the behavior of error metrics and optimal regularization parameters (γ) while increasing the number of data points (n), and keeping the ratios p/n and N/n fixed. As mentioned earlier, the experiment generates data according to the law $Y = \Sigma^\top b + \epsilon$, where b is a sparse vector, and Σ represents the network's internal operations. For each increment in n, we determine the optimal γ that minimizes both E_{test} and \hat{E}_{test} .

However, when we compare the optimal γ values required for minimizing the ridge loss function we defined earlier with those minimizing the test errors (E_{test} and \hat{E}_{test}), we observe a noticeable difference. As shown in the graphs below, the optimal γ for the ridge loss function does not always align with the optimal values for E_{test} and \hat{E}_{test} .

In this experiment, we also consider the "direct" loss function, given by Equation (3.13), which measures the Frobenius norm difference between the predicted coefficients $\hat{\beta}$ and the target vector b. As seen in the results, the optimal γ values obtained for minimizing this direct loss function differ from those obtained when optimizing the test errors, E_{test} and \hat{E}_{test} . The graphs clearly illustrate this difference, with the values for γ needed to minimize the direct loss function consistently diverging from the values required for minimizing the test errors. These results highlight the distinct behavior of the ridge loss function compared to the test errors derived from random matrix theory.

Effect of Regularization Strength on Ridge Regression Error (Different Densities). Figure 4.4 compares three plots of the normalized ridge error $\|\hat{\beta} - b\|_F^2 / \|b\|_F^2$ (3.13) versus γ under different densities for b. In the **middle panel** (approximately 50% density), we see a pronounced "U-shaped" curve with a clear global minimum. By contrast, the **left** and **right** panels (very high or very low density) show flatter or less sharply defined minima, suggesting that in those regimes, the ridge-regularized solution does not exhibit as strong a trade-off between underfitting and overfitting.

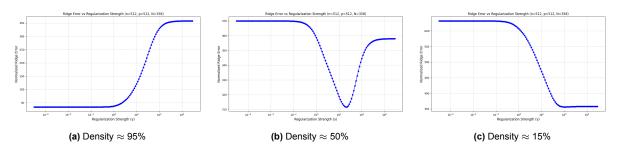


Figure 4.4: Normalized ridge error $\|\hat{\beta} - b\|^2/\|b\|^2$ as a function of the regularization strength γ (log scale) for three different densities. The center plot (50% density) exhibits a clear "U-shape," whereas the other two densities do not show a single, pronounced minimum.

4.1.4. Diagnosing the Causal Gap

Having reproduced the MSE-based optimum γ_{RMT}^* , we now turn to the question of **causal alignment**: does the predictive optimum also recover the ground-truth mechanism b?¹ Concretely, we contrast three objective curves as a function of γ :

 $^{^{1}}$ Causal faithfulness here is operationalised by the direct alignment loss \mathcal{L}_{ridge} introduced in (3.13).

- the **direct causal loss** $\mathcal{L}_{\text{ridge}}(\gamma)$ (3.13) (red dotted),
- the **RMT proxy** $\hat{E}_{\text{test}}(\gamma)$ (blue dotted), and
- the **empirical test MSE** $E_{\text{test}}(\gamma)$ (purple dotted).

The graphs below show the difference between the value of each objective curve once minimized by γ against the increasing number of points n. Different plots demonstrate different effects according to various values of p/n ratio as well as the density of an underlying true vector β .

Note that the grid of γ values from 10^{-7} to 10^5 through 0.1 increment was used for the optimization in this experiment. For each objective curve, the value of the corresponding loss function has been calculated for all values of γ , then the minimum has been chosen.

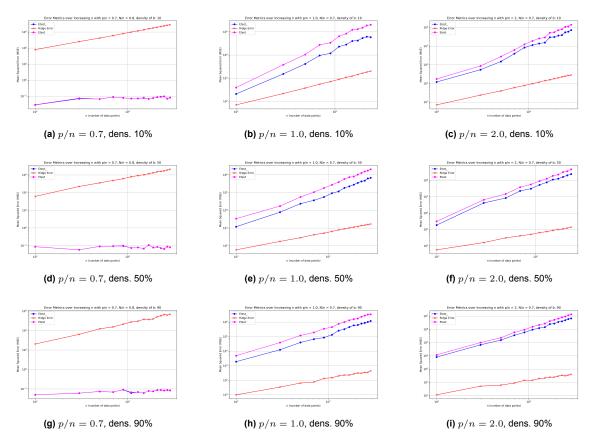


Figure 4.5: Direct causal loss $\mathcal{L}_{\text{ridge}}$ (3.13) (black), RMT-predicted \hat{E}_{test} (blue dashed) and empirical E_{test} (red dotted) versus γ across nine synthetic regimes. Stars mark the minimisers of each curve.

We sweep nine synthetic regimes (Figure 4.5): three aspect ratios $p/n \in \{0.7, 1.0, 2.0\}$ (left–to–right columns) and three sparsity levels of b (10 %, 50 %, 90 %; top–to–bottom rows). Each panel marks with a star the γ that minimises the corresponding curve.

It also does seem prudent to look at the actual optimal γ values found through loss minimization. Below is the graph of all such γ values.

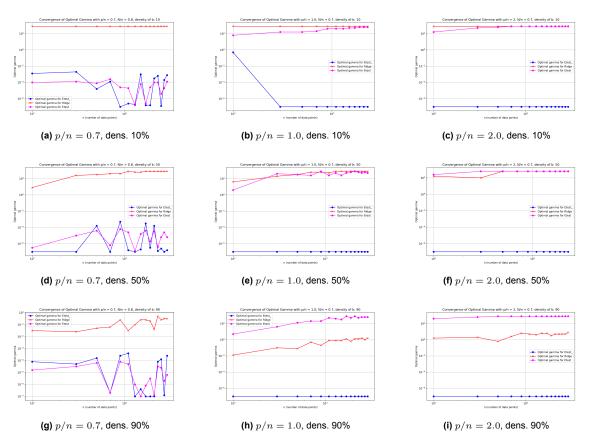


Figure 4.6: Optimal regularisation strength γ^* as a function of sample size n for nine synthetic regimes. Each panel reports: γ_{Etest} (blue), γ_{Etest} (magenta), and the causal γ_{Ridge} (red) that minimises $\mathcal{L}_{\mathsf{ridge}}$. Log–log axes highlight stability or drift across two orders of magnitude in n.

Sample–size sweep: predictive error vs. causal alignment. Figure 4.5 (error curves) and Figure 4.6 (corresponding γ^* 's) convey a consistent message across all nine regimes. First, the *predictive* quantities E_{test} and \hat{E}_{test} track each other almost perfectly once $n \gtrsim 200$, confirming the accuracy of the RMT surrogate. Second – and crucial for our causality narrative – the *causal* loss $\mathcal{L}_{\text{ridge}}$ (red curve) behaves very differently: it remains one to three orders of magnitude lower than the predictive errors for p/n < 1 and sparse b, and, even in the over-parameterised cases (p/n=1,2), it grows far more slowly.

The optimal regularisation strengths in Figure 4.6 expose why. The γ that minimises $\mathcal{L}_{\text{ridge}}$ (red \times) sits consistently around $10^{-1}-10^0$ and varies little with n, whereas the predictive optima (blue and magenta dots) collapse toward $\sim 10^{-5}$ as n grows. Put plainly, the model needs a substantially stronger penalty to recover the ground-truth coefficients than it does to minimise prediction error. This divergence widens with higher sparsity and larger aspect ratio, exactly the regimes where Figure 4.5 showed the largest causal gap. Hence, tuning γ on E_{test} (or its RMT proxy) is not merely sub-optimal for causal recovery—it systematically pushes the solution toward a coefficient vector that mis-aligns with the data-generating mechanism. Bridging this gap is therefore essential, motivating the gradient-based correction developed in the next section.

4.2. Loss-Landscape Analysis

The causal gap identified in Section 4.1.4 raises a natural question: how should the regularisation strength γ be adjusted to minimise the alignment loss $\mathcal{L}_{\text{ridge}}$ itself? In this section, we take a differential view, deriving closed-form (or nearly closed-form) expressions for the gradient $\partial \mathcal{L}_{\text{ridge}}/\partial \gamma$ under the synthetic law. We first inspect the oracle gradient – available when the true variance components are known. And then introduce an empirical variance estimator that yields a practical, data-driven surrogate.

Let us start first with deriving a closed-form expression for $\mathcal{L}(\gamma)$ as in 3.13

Note first that synthetic data generation law has two representations: $\mathbf{Y} = \mathbf{\Sigma}^{\top}b + \epsilon^{\top}$ version, or $\mathbf{Y} = b^{\top}\mathbf{\Sigma} + \epsilon$). β regressor from 3.1 has then two equivalent representations according to a data-generation law:

$$\hat{\beta} = \frac{1}{n}(\frac{1}{n}\mathbf{\Sigma}\mathbf{\Sigma}^\top + \gamma\mathbf{I}_N)^{-1}\mathbf{\Sigma}\mathbf{Y}^\top, \text{ and}$$

$$\hat{\beta} = \frac{1}{n} \mathbf{\Sigma} \left(\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_n \right)^{-1} \mathbf{Y}^{\top}.$$

We are going to proceed with the latter variant in our calculations. For the sake of simplicity, we are also going to omit the normalizing factor ||b|| and assume $\mathcal{L}(\gamma) = ||\hat{\beta} - b||_F^2$.

Then the derivation goes as follows:

$$\mathcal{L}(\gamma) = ||\hat{\beta} - b||_F^2 = (\hat{\beta} - b)^\top (\hat{\beta} - b) = \hat{\beta}^\top \hat{\beta} - 2\hat{\beta}^\top b + b^\top b.$$

Uncovering the first term of that, we get:

$$\hat{\beta}^{\top}\hat{\beta} = \frac{1}{n^2} \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma I_n)^{-1} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma I_n)^{-1} \mathbf{Y}^{\top}$$

Then, let us rearrange the terms a little and then add and subtract γI_n from/to the central side. We thus obtain:

$$\begin{split} \hat{\beta}^{\top} \hat{\beta} &= \frac{1}{n} \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} ((\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n}) - \gamma \mathbf{I}_{n}) (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} \mathbf{Y}^{\top} \\ &= \frac{1}{n} \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} (\mathbf{Y}^{\top} - \gamma (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} \mathbf{Y}^{\top}) \\ &= \frac{1}{n} (\mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} \mathbf{Y}^{\top} - \gamma \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-2} \mathbf{Y}^{\top})) \\ &= \frac{1}{n} \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-1} \mathbf{Y}^{\top} - \frac{\gamma}{n} \mathbf{Y} (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_{n})^{-2} \mathbf{Y}^{\top}. \end{split}$$

Therefore we can write down a full expression for $\mathcal{L}(\gamma)$ as such:

$$\mathcal{L}(\gamma) = \frac{1}{n}\mathbf{Y}(\frac{1}{n}\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma} + \gamma\mathbf{I}_n)^{-1}\mathbf{Y}^{\top} - \frac{\gamma}{n}\mathbf{Y}(\frac{1}{n}\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma} + \gamma\mathbf{I}_n)^{-2}\mathbf{Y}^{\top} - 2\frac{1}{n}\mathbf{Y}(\frac{1}{n}\boldsymbol{\Sigma}^{\top}\boldsymbol{\Sigma} + \gamma\mathbf{I}_n)^{-1}\boldsymbol{\Sigma}^{\top}b + b^{\top}b.$$

Next up, let us take the derivative of the expression above with respect to γ :

$$\frac{\partial \mathcal{L}(\gamma)}{\partial \gamma} = -\frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-2} \mathbf{Y}^{\mathsf{T}} - \left(\frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-2} \mathbf{Y}^{\mathsf{T}} - 2\gamma \frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-3} \mathbf{Y}^{\mathsf{T}} \right)
+ 2 \frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-2} \mathbf{\Sigma}^{\mathsf{T}} b
= -2 \frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-2} \mathbf{Y}^{\mathsf{T}} + 2\gamma \frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-3} \mathbf{Y}^{\mathsf{T}} + 2 \frac{1}{n} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\mathsf{T}} \mathbf{\Sigma} + \gamma \mathbf{I}_{n} \right)^{-2} \mathbf{\Sigma}^{\mathsf{T}} b.$$

Setting the derivative to zero and multiplying by n/2 yields:

$$\hat{\gamma} \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \hat{\gamma} \mathbf{I}_{n} \right)^{-3} \mathbf{Y}^{\top} - \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \hat{\gamma} \mathbf{I}_{n} \right)^{-2} \mathbf{Y}^{\top} + \mathbf{Y} \left(\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \hat{\gamma} \mathbf{I}_{n} \right)^{-2} \mathbf{\Sigma}^{\top} b = 0.$$

Among these terms the only unknown quantity is $\Sigma^{\top}b$, which under the synthetic law $Y = \Sigma^{\top}b + \varepsilon$ can be rewritten as $\mathbf{Y} - \varepsilon$.

In order to simplify further derivations, we will introduce some notation changes at this point. Note that in the experimental setting used in the central part of this thesis, target output \mathbf{Y} is of the size $1 \times n$, and therefore is a vector. We will refer to a vector-shaped target output variable as \mathbf{y} onward. Furthermore, let us introduce the following resolvent notation:

$$\mathbf{Q} \equiv (\frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_n)$$

Considering the above notation, the expression for the loss function derivative takes the following form:

$$\frac{\partial \mathcal{L}(\gamma)}{\partial \hat{\gamma}} = \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \mathbf{y}^{\top} + \mathbf{y} \mathbf{Q}^{-2} \Sigma^{\top} b$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \mathbf{y}^{\top} + \mathbf{y} \mathbf{Q}^{-2} (\mathbf{y}^{\top} - \epsilon^{\top})$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \mathbf{y}^{\top} + \mathbf{y} \mathbf{Q}^{-2} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \epsilon^{\top}$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \epsilon^{\top} = 0.$$

Note that the only out-of-sample part of the expression above is $\mathbf{y}\mathbf{Q}^{-2}\epsilon^{\top}$, due to the unknown noise parameter ϵ . To find the optimal parameter $\hat{\gamma}$ we are to estimate that part. In order to do that, let us take the expectation of that expression with respect to ϵ , taking into account that by definition, ϵ is a mean-zero noise factor with some variance σ^2 :

$$\mathbb{E}\left[\frac{\partial \mathcal{L}(\gamma)}{\partial \hat{\gamma}}|\epsilon\right]$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbb{E}[\mathbf{y} \mathbf{Q}^{-2} \epsilon^{\top} | \epsilon]$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbb{E}[(b^{\top} \mathbf{\Sigma} + \epsilon) \mathbf{Q}^{-2} \epsilon^{\top} | \epsilon]$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbb{E}[b^{\top} \mathbf{\Sigma} \mathbf{Q}^{-2} \epsilon^{\top} + \epsilon \mathbf{Q}^{-2} \epsilon^{\top} | \epsilon]$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbb{E}[\epsilon \mathbf{Q}^{-2} \epsilon^{\top} | \epsilon]$$

$$= \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \sigma^{2} \text{tr}(\mathbf{Q}^{-2}) = 0.$$

We thus obtain the "oracle" expression for direct loss derivative. But in order to extract the optimal parameter from the formula - we would need to know the real variance value. Therefore, the question now is – how do we approximate the real out-of-sample noise variance σ^2 . The most straightforward approach would be to consider the Residual Sum of Squares estimator:

$$\hat{\sigma^2}_{RSS} = \frac{(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})}{n} \quad ,$$

Where $\hat{\mathbf{y}} = \hat{\beta}^{\top} \mathbf{\Sigma}$.

Hence, Loss-Landscape Analysis yields three expression forms for the loss function derivative:

1.
$$\frac{\partial \mathcal{L}(\gamma)}{\partial \hat{\gamma}} = \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^{\top} - \mathbf{y} \mathbf{Q}^{-2} \epsilon^{\top} \quad \longleftarrow \text{real}$$

$$2. \quad \left. \frac{\partial \mathcal{L}(\gamma)}{\partial \hat{\gamma}} \right|_{\mathbb{E}[\cdot|\epsilon]} = \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^\top - \sigma^2 \text{tr}(\mathbf{Q}^{-2}) \quad \longleftarrow \text{ under expectation (the "oracle")}$$

$$3. \quad \left. \frac{\partial \mathcal{L}(\gamma)}{\partial \hat{\gamma}} \right|_{\hat{\sigma}} = \hat{\gamma} \mathbf{y} \mathbf{Q}^{-3} \mathbf{y}^\top - \hat{\sigma}_{RSS}^2 \mathrm{tr}(\mathbf{Q}^{-2}) \quad \longleftarrow \text{ with estimated variance}$$

In the following sections we are going to investigate the behavior of parameter γ found through the oracle and the expression with the estimated variance, and compare that γ with the one found through optimizing \hat{E}_{train} and \hat{E}_{test} .

4.3. Robustness study

In this subsection we are going to conduct a series of experiments aimed at studying values of regularization parameter γ related to optimizing each of the above-derived expressions and the behavior of the direct loss computed using those values.

We begin the robustness analysis by asking a simple question: do the derivative–based losses admit a (practical) minimum when measurement noise is present? Concretely, we plot the derivative of the causal loss $\partial \mathcal{L}/\partial \gamma$ as a function of γ for several noise levels and examine whether we can drive it arbitrarily close to zero on a sufficiently fine logarithmic grid.

Derivative formulas. We consider two expressions (derived in the previous section) that differ only in the variance term. We write them for $\gamma > 0$ and $\mathbf{Q}(\gamma) \equiv (\mathbf{A}^{\top}\mathbf{A} + \gamma\mathbf{I})$ in the following way:

$$\left. \frac{\partial \mathcal{L}(\gamma)}{\partial \gamma} \right|_{\text{real variance}} = \gamma \mathbf{y}^{\top} \mathbf{Q}(\gamma)^{-3} \mathbf{y} - \sigma^{2} \operatorname{tr}(\mathbf{Q}(\gamma)^{-2}).$$
 (4.1)

$$\left. \frac{\partial \mathcal{L}(\gamma)}{\partial \gamma} \right|_{\text{estimated variance}} = \gamma \mathbf{y}^{\top} \mathbf{Q}(\gamma)^{-3} \mathbf{y} - \hat{\sigma}_{RSS}^{2} \operatorname{tr}(\mathbf{Q}(\gamma)^{-2}), \tag{4.2}$$

where $\hat{\sigma}_{RSS}^2$ denotes a data-dependent variance plug-in (estimated variance).

Notice that with slight abuse of notation we put a σ subscript in 4.1 instead of $\mathbb{E}[\cdot|\epsilon]$ for the sake of simplicity. That meant to denote that the oracle derivative uses the true value of variance σ of the noise term $\varepsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$, The estimated variant replaces σ^2 with the residual estimate $\hat{\sigma}^2_{\text{RSS}}$.

Experimental setup. We fix $n=200,\ p/n=1.0,\ N/n=0.8$ and a 50 %-dense b. For each standard deviation level $\sigma\in\{1,2,3,4,5\}$ we:

- 1. generate synthetic data $Y = b^{T} \Sigma + \varepsilon$ with $\varepsilon \sim \mathcal{N}(0, \sigma^{2}\mathbf{I})$;
- 2. sweep γ on logarithmic grids $10^{-5} \le \gamma \le 10^3$ (estimated) and $10^{-3} \le \gamma \le 10^2$ (oracle);
- 3. evaluate (4.1) and (4.2) at each γ .

Operationally, we call a loss *minimisable on the grid* if there exists a grid point γ for which the (nonnegative) derivative is closest to zero, i.e. it minimizes $\left|\partial \mathcal{L}/\partial \gamma\right|$ (or equivalently attains the smallest positive value).

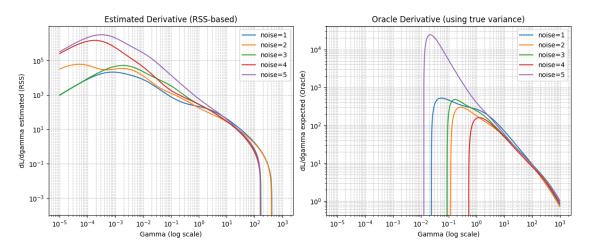


Figure 4.7: Derivative of the causal loss under five noise levels. Left: RSS \Box based estimate on $\gamma \in [10^{-5}, 10^3]$. Right: oracle derivative on $\gamma \in [10^{-3}, 10^2]$ (both axes log–log).

Observation. Across all tested noise levels, the curves of $\partial \mathcal{L}/\partial \gamma$ (both oracle and estimated) decrease with γ and admit a unique zero crossing within the displayed range. Consequently, with a sufficiently fine logarithmic grid, we can always find a γ whose derivative value lies arbitrarily close to zero, i.e. both losses are *minimisable* in the above grid sense. This justifies using the derivative root (or its closest grid proxy) as a robust, noise-aware choice of the regulariser, which we will compare directly against \hat{E}_{train} and \hat{E}_{test} -based selections in the next subsection.

Fixed vs. nonfixed RSS in the derivative

We now specify two practical RSS–based plug-ins for (4.2). Let $\hat{\sigma}^2_{RSS}(\cdot)$ be the residual–variance estimator as a function of the argument indicated in parentheses. Then:

fixed-RSS:
$$\frac{\partial \mathcal{L}}{\partial \gamma}\Big|_{\text{RSS, fixed}} (\gamma; \gamma_{\text{rss}}) = \gamma \mathbf{y}^{\top} \mathbf{Q}(\gamma)^{-3} \mathbf{y} - \hat{\sigma}_{\text{RSS}}^2(\gamma_{\text{rss}}) \operatorname{tr}(\mathbf{Q}(\gamma)^{-2}),$$
 (4.3a)

Both (4.3a)-(4.3b) are concrete instances of the estimated-variance derivative (4.2), differing only in whether the RSS variance is evaluated at a fixed reference $\gamma_{\rm rss}$ ("fixed–RSS") or at the current γ ("nonfixed–RSS").

where $\mathbf{Q}(\gamma) = \frac{1}{n} \mathbf{\Sigma}^{\top} \mathbf{\Sigma} + \gamma \mathbf{I}_n$. In the *fixed* variant, the variance is computed once at a reference value γ_{rss} (here we take $\gamma_{\mathrm{rss}} = \arg\min_{\gamma} \hat{E}_{\mathrm{test}}(\gamma)$), and then held fixed while we scan γ in the derivative. In the *nonfixed* variant, the same γ used in the derivative also enters the RSS computation, effectively defining a self-consistent, γ -dependent variance plug-in.

Experimental setting. We use: $n=100,\ p/n=1.0,\ N/n=0.8,\ a$ 50%-dense b, noise levels $\sigma\in\{1,2,\dots,14\}$, and a logarithmic grid $\gamma\in[10^{-6},10^3]$ with step $\Delta\log_{10}\gamma=0.01.$ For robustness, all curves are averaged over six independent runs. At each noise level we select four candidates: $\gamma_{\hat{E}_{\text{train}}}$, γ_{oracle} (root of the oracle derivative), and γ_{RSS} (root of the corresponding RSS derivative). We then evaluate $\mathcal{L}_{\text{ridge}}$ at each selected γ .

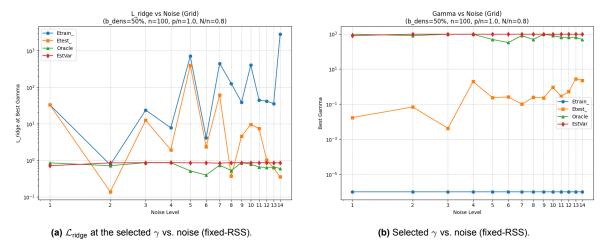


Figure 4.8: Fixed-RSS variant: variance $\hat{\sigma}^2_{\text{RSS}}(\gamma_{\text{rss}})$ computed at $\gamma_{\text{rss}} = \arg\min_{\gamma} \hat{E}_{\text{test}}(\gamma)$, then held constant while scanning γ in the derivative. Points are averages over six runs.

Fixed-RSS observations (Figure 4.8). Across noise levels, the *causal* selections ($\gamma_{\text{oracle}}, \gamma_{\text{RSS}}$) yield $\mathcal{L}_{\text{ridge}}$ that remains near $\mathcal{O}(1)$ and changes smoothly with noise, while the prediction–driven selections ($\gamma_{\hat{E}_{\text{train}}}, \gamma_{\hat{E}_{\text{test}}}$) produce larger and more volatile $\mathcal{L}_{\text{ridge}}$. The corresponding γ values show that causal tuning prefers substantially stronger regularisation.

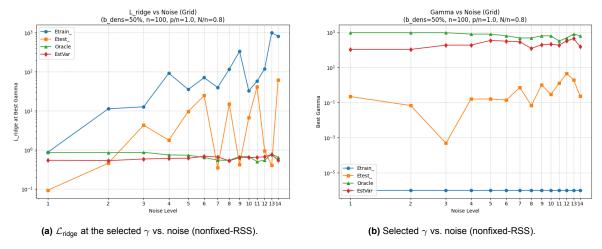


Figure 4.9: Nonfixed-RSS variant: variance $\hat{\sigma}^2_{\text{RSS}}(\gamma)$ recomputed at each grid point and fed back into the derivative. Points are averages over six runs.

Nonfixed-RSS observations (Figure 4.9). Allowing the variance estimate to follow the scanned γ produces a self-consistent root and, in several noise regimes, an even lower $\mathcal{L}_{\text{ridge}}$ than fixed-RSS while retaining the same qualitative behaviour (stable loss and stronger γ). Both variants confirm the central robustness claim of this section: a causally aligned regulariser can be obtained for a wide range of noise levels, and it consistently outperforms predictive selections in direct parameter recovery.

4.4. Convergence Study

The robustness results above showed that a causally aligned regulariser can be recovered across a wide range of noise levels. We now ask a complementary question: how do the corresponding choices of γ and their achieved alignment error behave as the sample size n grows? In other words, do the derivative-based prescriptions stabilise (and agree) as we move toward the high–dimensional asymptotic?

Experimental set-up. We fix the aspect ratios and ground-truth sparsity and sweep n, such that:

$$p/n = 1.0,$$
 $N/n = 0.8,$ density(b) = 50%, $\sigma = 8.$

As usual, for each $n \in \{100, 200, 300, 400, 500\}$ we generate synthetic data $Y = b^\top \Sigma + \varepsilon$ and compute four candidates on the same logarithmic grid $\gamma \in [10^{-6}, 10^3]$ with step $\Delta \log_{10} \gamma = 0.01$: (i) $\gamma_{\widehat{E}_{\text{train}}}$ minimising $\widehat{E}_{\text{train}}$, (ii) $\gamma_{\widehat{E}_{\text{train}}}$ minimising $\widehat{E}_{\text{test}}$, (iii) γ_{oracle} as the root of the oracle derivative $\partial \mathcal{L}/\partial \gamma$ in (4.1), and (iv) γ_{estVar} as the root of the empirical variance version (4.2). To smooth finite-sample fluctuations we average every curve over ten independent repeats.² For convergence of the *alignment* itself we report $\mathcal{L}_{\text{ridge}}(\gamma)$ evaluated at each selected γ .

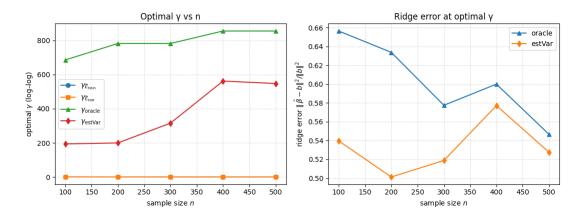


Figure 4.10: Behaviour of the regularization parameter and direct loss expression with growing sample size n (averaged over 10 runs; $p/n=1.0,\ N/n=0.8,\ \sigma=8,\ 50\%$ density of b). Left panel: the four γ selections $\gamma_{\widehat{E}_{\text{train}}},\ \gamma_{\widehat{E}_{\text{test}}},\ \gamma_{\text{oracle}},\ \text{and}\ \gamma_{\text{estVar}}.$ Right panel: direct alignment loss $\mathcal{L}_{\text{ridge}}$ at γ_{oracle} and γ_{estVar} .

Findings

- Derivative-based gap shrinks with n. As n increases from 100 to 500, the two causal selections move closer: γ_{oracle} remains in the high hundreds (about 7×10^2 - 9×10^2), while γ_{estVar} rises from roughly 2×10^2 toward 5×10^2 - 6×10^2 , narrowing the gap. In contrast, the predictive optima $\gamma_{\widehat{E}_{\text{train}}}$ and $\gamma_{\widehat{E}_{\text{test}}}$ stay near zero (orders of magnitude smaller). Reading the log axis: the tick labels show the actual γ values (spacing is logarithmic only). Thus a point around "600" means $\gamma\approx6\times10^2$; "850" means $\gamma\approx8.5\times10^2$ (not 10^{600} , etc.).
- Alignment gap shrinks with n. The right panel of Figure 4.10 shows $\mathcal{L}_{\text{ridge}}(\gamma_{\text{oracle}})$ and $\mathcal{L}_{\text{ridge}}(\gamma_{\text{estVar}})$ both decreasing as n grows, and their difference narrows–clear finite-sample convergence behaviour.
- Empirical edge persists. Across the inspected n, the estVar choice attains $\mathcal{L}_{\text{ridge}}$ that is consistently comparable to, and often lower than, the oracle echoing the advantage observed in the robustness study.

 $^{^{2}}$ Runs were executed in batches on the cluster, which is why the displayed grid of n stops at 500.

Takeaway (revised). Due to runtime limits on the TU Delft cloud cluster and averaging over repeated runs, we report $n \in \{100, 200, 300, 400, 500\}$. This range is too short to *prove* asymptotic convergence, but Figure 4.10 already exhibits stabilisation of the derivative-based γ in the 10^2 - 10^3 decade and a steady reduction of the alignment error with n. We fixed the noise level to $\sigma = 8$ (no special theoretical status); this choice was guided empirically by earlier experiments showing that for $\sigma \gtrsim 6$ -7, the estVar rule tracks or slightly improves upon the oracle. Within these constraints, the evidence supports the practical conclusion: the data-driven derivative selection $\gamma_{\rm estVar}$ is stable in n, operates on the same scale as the oracle, and delivers improving $\mathcal{L}_{\rm ridge}$ as sample size grows, while predictive selections remain far smaller and misaligned with the causal objective.

4.5. Real-Data Validation: Fashion-MNIST

We validate the proposed estimator on real images using a two-class subset of Fashion-MNIST (labels 1 vs. 2) under the same ReLU random-features ridge classification pipeline and preprocessing used throughout the paper, therefore, mimicking the author's experiment. The study proceeds in two steps. First, we assess *predictive robustness* in a head-to-head comparison with the Coulliet $E_{\rm test}$ calibration: over 50 Monte–Carlo runs (new random feature matrices W each run), we train both methods on identical train/test splits and aggregate accuracies and confusion matrices to verify that our approach does not underperform Coulliet's method on average. Second, holding the baseline split fixed, we conduct a controlled *perturbation analysis*: we introduce localized edits to the *training* inputs (single pixel and square patch brightness shifts of varying sizes), refit both models on the perturbed data using their respective γ selection rules, and compare the resulting coefficient vectors β against their baseline counterparts learned on the original data. This design lets us quantify not only changes in test performance but also how each method's inferred importance pattern *moves under perturbations* (stability vs. reactivity and concentration), providing empirical evidence about the interpretability and robustness of the learned representations.

4.5.1. Robustness analysis - fashion MNIST

We begin by conducting a robustness study to verify that our proposed method performs at least as well as the Coulliet's $E_{\rm test}$ calibration in terms of predictive accuracy. Specifically, we run 50 Monte-Carlo trials on the two-class Fashion MNIST split under identical preprocessing and ReLU random-features ridge settings, aggregate accuracies and confusion matrices, and compare the two methods head-to-head.

Setting and data. We evaluate the proposed estimator on a two-class subset of Fashion MNIST (labels 1 vs. 2), cast as a binary classification problem with labels $y \in \{-1, +1\}$. Images are 28×28 ; we vectorize each image to $x \in \mathbb{R}^p$ with p = 784. We adopt the same random-features ridge pipeline used throughout the paper, with a ReLU nonlinearity:

$$\sigma(t) = \max\{t, 0\} \quad ,$$

and the corresponding ReLU (first-order arc-cosine) kernel $K(\cdot,\cdot)$ used to build the deterministic-equivalent criterion.

Data Preprocessing. Before training, the Fashion-MNIST data are transformed to match the theoretical assumptions of the random-features ridge model and ensure consistent global scaling across Monte-Carlo runs. The procedure proceeds as follows.

In terms of class selection and vectorization, we retain only two classes (labels 1 and 2), corresponding to a binary classification task with $y \in \{-1, +1\}$. Each 28×28 grayscale image is vectorized into $x \in \mathbb{R}^p$, p = 784. Let $X = [x_1, \dots, x_M] \in \mathbb{R}^{p \times M}$ denote the full dataset of all selected samples.

Then, we perform global normalization. All pixel intensities are first scaled to the unit interval by

dividing by the global maximum

$$X \leftarrow \frac{X}{\max(X)}.$$

This ensures all entries of X lie in [0, 1].

We then perform a global centering and rescaling step so that the average squared norm of the centered vectors equals p. Let

$$\mu_g = \frac{1}{M} \sum_{i=1}^{M} x_i, \qquad X_c = X - \mu_g \mathbf{1}_M^{\mathsf{T}},$$

and define the global scaling factor

$$s_g = \frac{\sqrt{p}}{\sqrt{\frac{1}{M} \sum_{i=1}^{M} \|x_i - \mu_g\|_2^2}}.$$

The globally standardized data are then

$$X_{\mathsf{std}} = s_a X_c.$$

After global standardization, we perform a second normalization across the pooled subset of the two chosen classes to ensure comparability of within-class magnitudes. Let

$$X^{(1)}, X^{(2)} \in \mathbb{R}^{p \times M_j}$$

denote the subsets for classes 1 and 2, respectively. We form their pooled matrix

$$X_{\text{pool}} = [X^{(1)} \ X^{(2)}]$$
 ,

and compute

$$\mu_p = \frac{1}{M_1 + M_2} \sum_i x_i^{\text{(pool)}} \quad , \qquad s_p = \frac{\sqrt{p}}{\sqrt{\frac{1}{M_1 + M_2} \sum_i \|x_i^{\text{(pool)}} - \mu_p\|_2^2}} \quad .$$

Each class subset is then re-centered and rescaled via

$$X^{(j)} \leftarrow s_p \left(X^{(j)} - \mu_p \right) \quad , \qquad j \in \{1, 2\} \quad .$$

Next step is the train–test partitioning. With a prescribed aspect ratio $c_1 = p/n$ and training fraction ρ_{tr} , the number of training and test samples per class are

$$n = \left\lfloor rac{p}{c_1}
ight
floor, \qquad n_{\mathsf{te}} = \left\lfloor n \, rac{1 -
ho_{\mathsf{tr}}}{
ho_{\mathsf{tr}}}
ight
floor.$$

For each class j, a random permutation π_j (with fixed seed) selects disjoint subsets for training and test:

$$X_{\mathrm{tr}}^{(j)} = X_{[:,\,\pi_j(1:n_j)]}^{(j)}, \qquad X_{\mathrm{te}}^{(j)} = X_{[:,\,\pi_j(n_j+1:n_j+n_{\mathrm{te},j})]}^{(j)} \quad .$$

The final training and test matrices are concatenations

$$X_{\mathsf{tr}} = [\, X_{\mathsf{tr}}^{(1)} \ X_{\mathsf{tr}}^{(2)} \,], \qquad X_{\mathsf{te}} = [\, X_{\mathsf{te}}^{(1)} \ X_{\mathsf{te}}^{(2)} \,] \quad ,$$

with corresponding binary label vectors

$$y_{\mathsf{tr}} = \begin{bmatrix} -\mathbf{1}_{n_1} \\ +\mathbf{1}_{n_2} \end{bmatrix}, \qquad y_{\mathsf{te}} = \begin{bmatrix} -\mathbf{1}_{n_{\mathsf{te},1}} \\ +\mathbf{1}_{n_{\mathsf{te},2}} \end{bmatrix}.$$

In summary, this preprocessing pipeline enforces deterministic centering and scaling of all images, guarantees $\mathbb{E}\|x_i\|_2^2 \approx p$, and produces balanced train–test splits under fixed seeds. The resulting $(X_{\rm tr}, X_{\rm te}, y_{\rm tr}, y_{\rm te})$ pairs form the standardized inputs to the ReLU random-features ridge classifier used throughout the validation experiments.

Model and hyperparameters. Similarly to the Coulliet's experiment We consider a random-features ridge classifier with hidden width (inner neuron layer) $N = \lfloor c_2 n \rceil$ and i.i.d. Gaussian features $W \in \mathbb{R}^{N \times p}$, $W_{ij} \sim \mathcal{N}(0,1)$. Given inputs $X \in \mathbb{R}^{p \times n}$ and labels $y \in \{-1,+1\}^n$, let

$$\Sigma = \sigma(WX)$$

with $\sigma(t) = \max\{t,0\}$ (ReLU). For any ridge parameter $\gamma > 0$, the coefficient vector is obtained *in closed form* from the standard ridge expression (see (3.1)):

$$\widehat{\beta}(\gamma)$$
 as in (3.1).

Prediction uses the score $f(x) = \Sigma(x)^{\top} \widehat{\beta}(\gamma)$ and the decision rule $\operatorname{sign}(f(x))$. The *choice* of γ is determined by the calibration rules described next; $\widehat{\beta}$ is then formed by plugging that γ into (3.1).

Coulliet calibration and our direct calibration. We compare two one–parameter selection rules for γ evaluated on the same logarithmic grid $\gamma \in \{10^t : t \in [-6, 4], \Delta t = 0.005\}$:

• Coulliet (\overline{E}_{test}) calibration. We select

$$\gamma_{\rm C} = \arg\min_{\gamma} \ \overline{E}_{\sf test}(\gamma),$$

where $\overline{E}_{\text{test}}(\gamma)$ is the deterministic–equivalent test proxy from the theory (cf. Eq. (3.12), which depends on the fixed–point $\delta(\gamma)$).

· Direct calibration (nonfixed RSS; ours). We select

$$\gamma_{\rm D} = \arg\min_{\gamma} |\partial_{\gamma} \mathcal{L}_{\rm direct}^{\rm (RSS)}(\gamma)|,$$

where $\mathcal{L}_{\mathrm{direct}}^{\mathrm{(RSS)}}(\gamma)$ is the nonfixed RSS–based direct loss from our loss–landscape analysis (Eq. (4.3b)). In other words, we pick the grid point whose derivative magnitude is minimal.

In both cases, once γ is chosen (either $\gamma_{\rm C}$ or $\gamma_{\rm D}$), we compute $\widehat{\beta}$ via the closed–form ridge expression (3.1) and classify by ${\rm sign}(f)$.

Monte–Carlo protocol. We run 50 Monte–Carlo trials to assess predictive robustness. The training/test splits $(X_{\text{train}}, X_{\text{test}})$ and labels are held fixed across trials. The Coulliet selection $\gamma_{\rm C}$ is computed once from $(X_{\text{train}}, X_{\text{test}})$ via $\overline{E}_{\text{test}}(\gamma)$ and is common to all trials. For each trial $r=1,\ldots,50$, we draw an independent $W^{(r)}$, recompute the direct selection $\gamma_{\rm D}^{(r)}$ on the grid (the direct criterion depends on W), form

$$\widehat{\beta}_{\mathrm{C}} = \widehat{\beta}(\gamma_{\mathrm{C}}), \qquad \widehat{\beta}_{\mathrm{D}}^{(r)} = \widehat{\beta}(\gamma_{\mathrm{D}}^{(r)}),$$

and evaluate train/test predictions by thresholding at zero. We record per–run accuracies and confusion vectors (TP, TN, FP, FN), and we report run–wise average accuracies together with confusion totals summed over the 50 runs (normalized by 50 when per–run rates are required).

Concrete hyperparameters (this experiment).

activation = ReLU,
$$(L_-, L_+) = (1, 2), \quad c_1 = \frac{p}{n} = 1.0, \quad c_2 = \frac{N}{n} = 1.0, \quad \text{train fraction } \tau = 0.8,$$
 $\gamma \text{ grid: } \gamma \in \{10^t: t \in [-6, 4], \ \Delta t = 0.005\}.$

Here c_1 sets $n\approx p$ to retain high-dimensional setting, and c_2 fixes $N\approx n$. The test proxy $\overline{E}_{\text{test}}(\gamma)$ uses the ReLU (first–order arc–cosine) kernel as in the theory section, through Eq. (3.12); our direct rule uses Eq. (4.3b). In all cases, the final coefficients are obtained from the ridge closed form (3.1).

Outputs. For transparency and downstream analysis, we store (i) the per–run accuracy arrays for both methods on train and test, (ii) the summed confusion counts across runs, and (iii) the per–run direct selections $\{\gamma_{\rm D}^{(r)}\}_{r=1}^{50}$ together with the common Coulliet $\gamma_{\rm C}$, enabling paired head–to–head statistical comparisons and dispersion summaries across the 50 trials.

Predictive results on Fashion–MNIST. Under the ReLU random–features ridge model with p=n=N=784 and 50 Monte–Carlo trials, both calibration rules achieve near–perfect prediction. The Coulliet selection (via $\overline{E}_{\text{test}}$) yields $\gamma_{\text{C}}=70.79$, whereas our direct rule (nonfixed RSS derivative) concentrates at a larger scale (median $\gamma_{\text{D}}\approx 295.12$). Averaged across runs, test accuracy is 0.9956 for the direct rule versus 0.9930 for Coulliet; train accuracy is 0.9943 (direct) versus 1.0000 (Coulliet). The test–set confusion totals (summed over runs) are reported below.

	Coulliet ($\gamma_{ m C}$)		Direct ($\gamma_{ m D}$)	
	Positive	Negative	Positive	Negative
True	4833	4898	4861	4896
False	2	67	4	39
Accuracy	0.9930		0.9956	

Table 4.1: Test confusion matrices (totals over 50 Monte—Carlo trials) for the two calibration methods. Entries correspond to counts of true and false predictions across all runs.

Discussion of findings. (i) Both methods are highly accurate; the direct calibration attains a small but consistent improvement on test (+0.27 percentage points). (ii) The gain is driven primarily by fewer false negatives (FN: $67 \rightarrow 39$), with a negligible increase in false positives (FP: $2 \rightarrow 4$). This translates into higher recall for the positive class under our rule. (iii) The larger γ selected by the direct method implies stronger regularization; correspondingly, its train accuracy is slightly below perfect (0.9943 vs. 1.0000 for Coulliet), yet generalizes marginally better on test-consistent with reduced overfitting. (iv) Overall, the direct rule matches or improves predictive performance while selecting a regularization scale aligned with the loss–landscape analysis, reinforcing its practical viability alongside Coulliet's $\overline{E}_{\rm test}$ calibration.

4.5.2. Perturbation analysis: Fashion--MNIST

Having established that our direct-loss calibration for the regularization parameter γ attains accuracy comparable to Coulliet's method, we now examine the behaviour and structure of the learned coefficient vector β . We probe the *stability vs. reactivity* of the two calibration rules under localized input edits. Starting from a fixed train/test split on the Fashion–MNIST 1 vs. 2 task and a fixed random-features matrix W, we apply controlled brightness shifts to selected pixels (singletons and square patches) *in the training inputs only*, refit each method with its own γ -selection rule, and compare the learned coefficient vectors to their baseline counterparts. This design isolates how each rule's inferred importance pattern moves under small, spatially localized changes. The test data remain unchanged throughout the analysis, ensuring that any observed differences stem solely from modifications in the training set. This allows us to attribute changes in the learned coefficients to the learning dynamics of each calibration rule, rather than to shifts in the evaluation distribution.

Baseline. With activation $\sigma(t) = \max\{t, 0\}$ (ReLU), aspect ratios $c_1 = p/n = 1.0$ and $c_2 = N/n = 1.0$ (p = n = N = 784), and train fraction $\tau = 0.8$, we construct

$$\Sigma = \sigma(WX_{\rm tr}), \qquad \Sigma(x) = \sigma(Wx),$$

just like before, and then select γ by two one–dimensional rules evaluated on the common logarithmic grid

$$\gamma \in \{10^t : t \in [-6, 3.5], \Delta t = 0.005\}.$$

Coulliet's calibration picks

$$\gamma_{
m C} = \arg\min_{\gamma} \; \overline{E}_{
m test}(\gamma) \;\;$$
 (deterministic equivalent; cf. Eq. (3.12)),

while our direct rule picks

$$\gamma_{\mathrm{D}} = \arg\min_{\gamma} \left| \partial_{\gamma} \mathcal{L}_{\mathrm{direct}}^{(\mathrm{RSS})}(\gamma) \right| \quad \text{(nonfixed RSS; cf. Eq. (4.3b))}.$$

Given either γ , the coefficients are formed by the closed–form ridge expression (Eq. (3.1)), yielding baseline vectors

$$\beta_{\mathrm{C}}^{\mathrm{base}} = \widehat{\beta}(\gamma_{\mathrm{C}}), \qquad \beta_{\mathrm{D}}^{\mathrm{base}} = \widehat{\beta}(\gamma_{\mathrm{D}}).$$

Perturbations (training-only, localized). Let $X_{\rm tr} \in \mathbb{R}^{p \times n}$ and $X_{\rm te} \in \mathbb{R}^{p \times n_{\rm te}}$ denote the standardized training and test matrices, respectively, constructed as described in Section 4.5.1. Recall that in our controlled high-dimensional setting, we fix p=n=784 and $n_{\rm te}\approx 196$, corresponding to 784 synthetic training and 196 test images (each of size 28×28) drawn from the two selected Fashion-MNIST classes. For a chosen subset of feature indices $S\subset\{1,\ldots,p\}$, a perturbation magnitude $\Delta>0$, and a sign $s\in\{-1,+1\}$, we define a *training-only* perturbation

$$X'_{\text{tr}} = X_{\text{tr}} + s \Delta \mathbf{1}_S \mathbf{1}_n^{\top}, \qquad X'_{\text{te}} = X_{\text{te}},$$

where $\mathbf{1}_S \in \mathbb{R}^p$ is an indicator vector equal to 1 on S and 0 elsewhere. The test matrix X_{te} remains unaltered, ensuring that performance differences can be attributed solely to modifications of the training set and not to shifts in the evaluation distribution.

We take $\Delta=0.5$ in the standardized feature space, which corresponds to a moderate brightness adjustment in the original image domain. The perturbation therefore increases or decreases the brightness of selected pixels across all training images (one entire pixel column of $X_{\rm tr}$) while leaving the test images intact. The following localized perturbation scenarios are considered:

- Single pixel: |S| = 1 (six random pixel locations; random sign).
- 5×5 patch: |S| = 25 (six random top-left locations; random sign).
- 10×10 patch: |S| = 100 (six random locations; random sign).
- 20×20 patch: |S| = 400 (six random locations; random sign; only if the 28×28 image size permits).

Each perturbation thus modifies the brightness of either a single pixel or a contiguous block of pixels across all n=784 training images, allowing us to evaluate how small, spatially localized changes in the training data affect the learned coefficients.

For each perturbed training matrix X'_{tr} we recompute the two γ 's on the *same* grid:

$$\gamma_{\mathrm{C}}^{(\mathrm{scen})} = \arg\min_{\gamma} \; \overline{E}_{\mathrm{test}}^{(\mathrm{scen})}(\gamma), \qquad \gamma_{\mathrm{D}}^{(\mathrm{scen})} = \arg\min_{\gamma} \; \big| \partial_{\gamma} \mathcal{L}_{\mathrm{direct}}^{(\mathrm{RSS})}(\gamma; X_{\mathrm{tr}}') \big|,$$

where $\overline{E}_{\mathrm{test}}^{(\mathrm{scen})}$ uses kernels built from $(X'_{\mathrm{tr}}, X_{\mathrm{te}})$ and the direct derivative is evaluated with the same fixed W but the edited X'_{tr} . We then form

$$\beta_{\mathrm{C}}^{(\mathrm{scen})} = \widehat{\beta} \Big(\gamma_{\mathrm{C}}^{(\mathrm{scen})} \Big) \,, \qquad \beta_{\mathrm{D}}^{(\mathrm{scen})} = \widehat{\beta} \Big(\gamma_{\mathrm{D}}^{(\mathrm{scen})} \Big) \,,$$

and evaluate train/test predictions on (X'_{tr}, X_{te}) by thresholding at zero.

For each scenario we store $(\gamma_C^{(scen)}, \gamma_D^{(scen)})$, the perturbed coefficients $(\beta_C^{(scen)}, \beta_D^{(scen)})$, and the corresponding train/test metrics. In the analysis, we will compare each perturbed vector to its baseline counterpart to quantify how concentrated or diffuse the response is under localized edits, and whether the direct rule's stronger regularization yields more stable importance patterns than the Coulliet selection.

Coefficient sensitivity under localized training edits. For each perturbation scenario s, we compare the perturbed coefficients $\beta^{(s)}$ to the baseline coefficients β^{base} via the relative squared ℓ_2 —deviation

$$\Delta_{\ell_2}(s) = 100 \times \frac{\|\beta^{(s)} - \beta^{\text{base}}\|_2^2}{\|\beta^{\text{base}}\|_2^2} [\%].$$

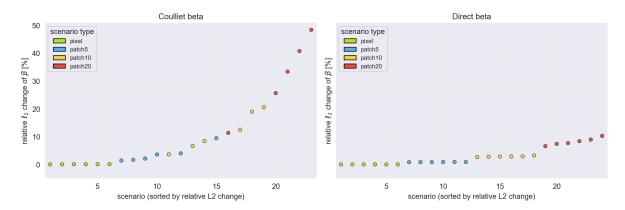


Figure 4.11: Relative ℓ_2 change of the learned coefficients under training-only, localized brightness edits. Left: Coulliet beta. Right: Direct beta. Scenarios are sorted by Δ_{ℓ_2} ; colors indicate perturbation type.

We report $\Delta_{\ell_2}(s)$ for all scenarios and for both calibration rules (Coulliet and Direct), sorting scenarios along the horizontal axis by increasing Δ_{ℓ_2} . Markers are color–coded by perturbation type (single pixel, 5×5 , 10×10 , 20×20). For readability, extreme outliers (top 1% of Δ_{ℓ_2} across all panels) are omitted from the plot; these correspond to boundary selections of γ and are discussed in the text.

Main takeaway Figure 4.11 summarizes the relative change of the learned coefficients under localized perturbations of the training data. Each marker corresponds to a single perturbation scenario, with the vertical axis showing the relative ℓ_2 —distance between the perturbed and baseline coefficients. Overall, Coulliet's calibration exhibits substantially larger coefficient shifts than the direct-loss rule, particularly for larger $(10 \times 10 \text{ and } 20 \times 20)$ patch edits. Quantitatively, the top Coulliet scenarios reach relative changes between 25% and 50%, with one extreme case exceeding $2.9 \times 10^6\%$ due to an instability in γ -selection near the boundary of the grid. In contrast, the corresponding top-five direct-loss scenarios remain in the range of roughly 7-10%. These results confirm that the direct-loss method produces significantly more stable coefficient vectors under identical data perturbations, suggesting that its calibration rule imposes a smoother and more robust regularization response. The extreme outlier in the Coulliet set (the 20×20 patch with a negative brightness shift) will be examined separately below to illustrate the mechanism behind such a large deviation.

Overlay of baseline and perturbed coefficients by perturbation type. To better understand how each calibration rule reacts to localized changes in the training data, we visualize the learned coefficient vectors β before and after perturbation. For each perturbation type (single pixel, 5×5 , 10×10 , and 20×20 patch), Figure 4.12 compares the baseline coefficients $\beta^{\rm (base)}$ with their perturbed counterparts $\beta^{(s)}$. Each subplot shows the coefficients sorted by their baseline order (ascending $\beta^{\rm (base)}$ values), allowing us to observe how localized brightness edits in the training set shift the learned importance pattern across the feature dimension. The left column corresponds to Coulliet's γ -selection rule, and the right column to our direct-loss calibration.

The blue line represents the baseline $\beta^{(\text{base})}$, while the orange line shows the coefficients obtained after the perturbation. Higher transparency of the orange curve allows the blue baseline to remain visible beneath it, highlighting where the perturbed β departs most from its reference. Together, these overlays provide an intuitive measure of model *stability* under input perturbations.

While Figure 4.12 allows direct inspection of how the learned coefficients shift along the feature index, the magnitude of those changes can sometimes be obscured when small and large coefficients coexist on the same scale. To make these deviations more apparent, we sort the coefficients by their absolute baseline magnitude in Figure 4.13, which emphasizes where the most influential features in the baseline model experience the strongest or weakest reweighting after perturbation.

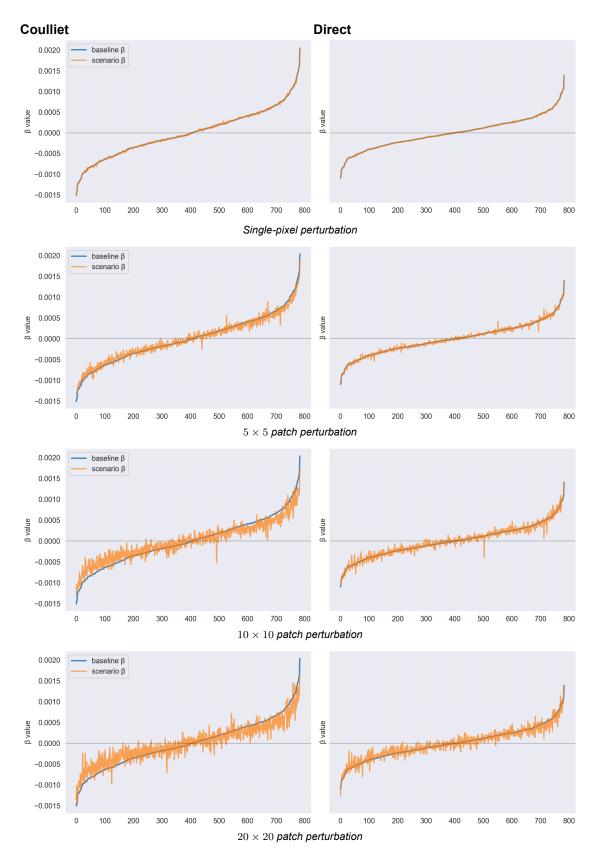


Figure 4.12: Overlay of baseline and perturbed coefficient vectors β for four representative perturbation types. Each row shows a two-panel image with Coulliet (left) and Direct (right). Blue: baseline coefficients $\beta^{(\mathrm{base})}$; orange: perturbed coefficients $\beta^{(s)}$.

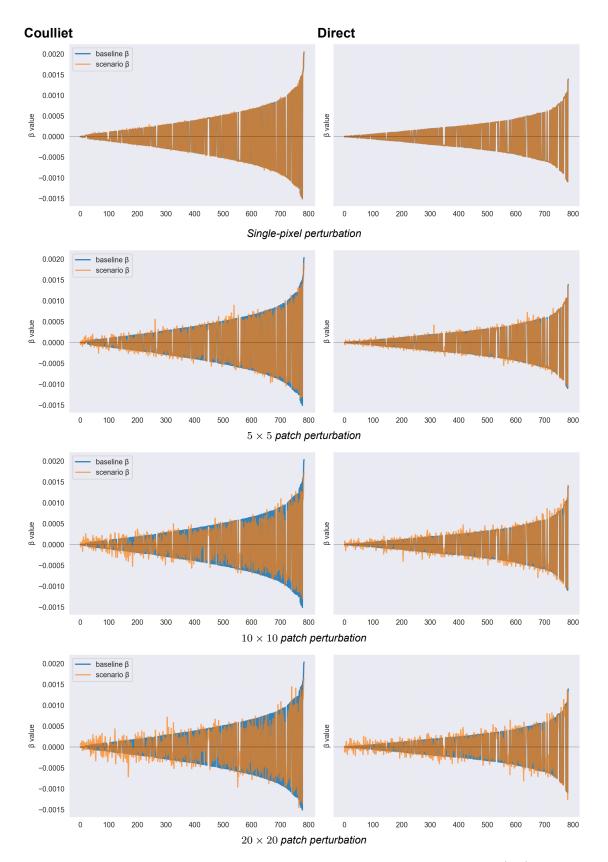
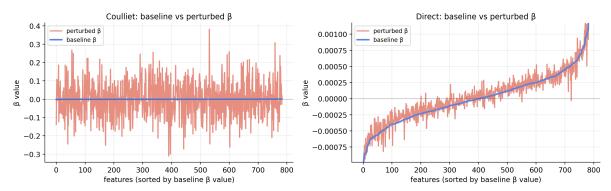


Figure 4.13: Same as Figure 4.12, but with coefficients sorted by their absolute baseline magnitude $|\beta^{(\mathrm{base})}|$. Sorting by absolute value highlights how perturbations affect the largest-magnitude coefficients most strongly.



Extreme 20×20 patch: baseline vs. perturbed β (features sorted by baseline order).

Figure 4.14: Stress case where Coulliet's calibration collapses to the grid floor $(\gamma_{\rm C}^{({\rm scen})} \approx 10^{-6})$, while the direct-loss rule stays in a moderate range. *Note*: the two panels use different vertical scales (we zoom the Direct panel to reveal its small variation). The blue baseline curves are essentially identical across panels; the visual discrepancy is solely due to the differing y-axis ranges.

Focused overlays on the most influential coefficients. Complementing the full–vector overlays, we zoom in on the extreme tails of the baseline coefficient distribution. For each perturbation type (single pixel, 5×5 , 10×10 , 20×20), we select the 100 coefficients from a chosen tail of the baseline vector and overlay the corresponding segments of $\beta^{(\mathrm{base})}$ and $\beta^{(s)}$:

$$\mathcal{I}_{\mathrm{tail}} \ = \ \begin{cases} \mathrm{indices\ of\ the\ 100\ smallest\ entries\ of\ } \beta^{\mathrm{(base)}}, & \text{if\ tail\ =\ bottom}, \\ \mathrm{indices\ of\ the\ } 100\ \mathrm{largest\ entries\ of\ } \beta^{\mathrm{(base)}}\ \mathrm{(ordered\ by\ } |\beta^{\mathrm{(base)}}|), & \mathrm{if\ tail\ =\ top}. \end{cases}$$

Plotting $\left\{\beta_i^{(\mathrm{base})},\,\beta_i^{(s)}\right\}_{i\in\mathcal{I}_{\mathrm{tail}}}$ against the rank within the selected tail emphasizes how the most negative/most positive baseline coefficients are reweighted by training perturbations. As before, the left column corresponds to Coulliet's calibration and the right to our direct-loss rule.

You can see the tail baseline-vs-scenario plots in figures 4.15 and 4.16 for top and bottom tail respectively.

Interpretation of tail-slice overlays. Examining the coefficient tails in Figures 4.15 and 4.16 reveals a clear structural difference between the two calibration rules. For Coulliet's method, the perturbations tend to *amplify* the absolute magnitude of the coefficients: in the bottom tail, the perturbed β 's are on average less negative (shifted upward), while in the top tail they are systematically lower than the baseline curve. This pattern indicates that Coulliet's calibration adjusts γ in a way that consistently rebalances the relative strength of extreme coefficients under localized training edits.

By contrast, the proposed *direct-loss* rule produces markedly more stable coefficient profiles. Across all perturbation types and both tails, the perturbed $\beta^{(s)}$ curves remain close to their baseline $\beta^{(\mathrm{base})}$ counterparts, with substantially reduced variance and no systematic bias in either direction. This suggests that the direct-loss calibration maintains a more consistent scaling of the learned weights and exhibits stronger robustness to localized input modifications.

A stress case where Coulliet's calibration collapses to vanishing regularization. To illustrate the behaviour behind the largest relative change omitted in the scatter 4.11 (for plotting purposes), we examine a 20×20 patch scenario that produces an extreme deviation for Coulliet's rule (top entry in the "Top-5 Coulliet changes"). In this instance the Coulliet grid search selects a *near-zero* penalty, $\gamma_{\rm C}^{\rm (scen)} \approx 10^{-6}$ (i.e., $\Delta \log_{10} \gamma \ll 0$ relative to the baseline), effectively turning off regularization. The direct-loss rule, by contrast, remains in a moderate range. We display two overlays of the coefficient vectors-unsorted (by the baseline order) and sorted by $|\beta^{\rm (base)}|$ to make the pattern visible under both views.

What happens and why? Two effects jointly explain the blow-up in the relative-change metric for Coulliet. (i) *Tiny denominator:* in this scenario the baseline Coulliet vector has a very small norm,

$$100 \cdot \frac{\|\beta^{(\text{scen})} - \beta^{(\text{base})}\|_2^2}{\|\beta^{(\text{base})}\|_2^2}$$

is therefore magnified even for moderate absolute differences. (ii) Near-zero regularization with rank deficiency: when $\gamma \to 0$, the ridge solution becomes highly sensitive to small singular directions of $A = \Sigma/\sqrt{n}$; with $\mathrm{rank}(A) = r < n$, the closed form adds a large null-space component (scaling like $1/\gamma$). Although this component is annihilated by Σ (so many coordinates stay near zero in the overlay), the span component can inflate along ill-conditioned directions, yielding a large relative change. In contrast, the direct-loss rule keeps γ away from the boundary, producing a smoother, better-conditioned update with much smaller variance and no catastrophic swing.

For reference, below are the top 5 scenario feature vector changes for both Coulliet and Direct method:

Largest relative changes.

Coulliet: top-5 $100 \cdot \ \beta^{(s)} - \beta^{base}\ _2^2 / \ \beta^{base}\ _2^2$			Direct: top-5 $100 \cdot \ \beta^{(s)} - \beta^{base}\ _2^2 / \ \beta^{base}\ _2^2$		
tag	type	change [%]	tag	type	change [%]
patch20_04_r7c3_d+0.50sgn-1	patch20	2,977,344.967	patch20_03_r0c0_d+0.50sgn1	patch20	10.275
patch20_05_r0c5_d+0.50sgn1	patch20	48.425	patch20_06_r0c8_d+0.50sgn1	patch20	8.966
patch20_06_r0c8_d+0.50sgn1	patch20	40.770	patch20_02_r1c2_d+0.50sgn1	patch20	8.384
patch20_02_r1c2_d+0.50sgn1	patch20	33.395	patch20_04_r7c3_d+0.50sgn-1	patch20	7.703
patch20_03_r0c0_d+0.50sgn1	patch20	25.682	patch20_05_r0c5_d+0.50sgn1	patch20	7.404

The Coulliet side shows markedly larger relative deviations; the first row is the outlier we analyze in the section above.

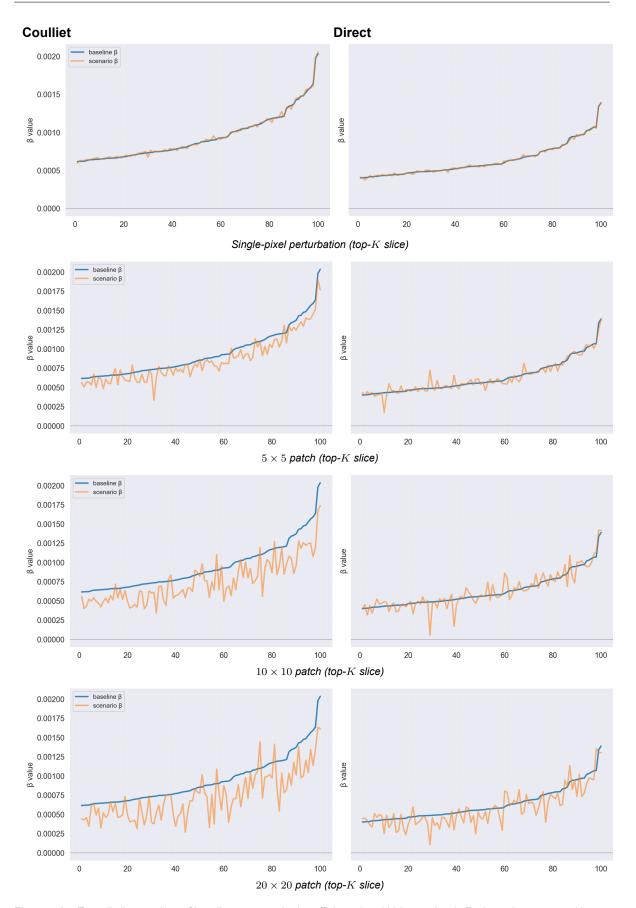


Figure 4.15: Top-tail-slice overlays of baseline vs. perturbed coefficients (top 100 beta values). Each row is a two-panel image with Coulliet (left) and Direct (right); the horizontal axis is the rank within the selected tail.

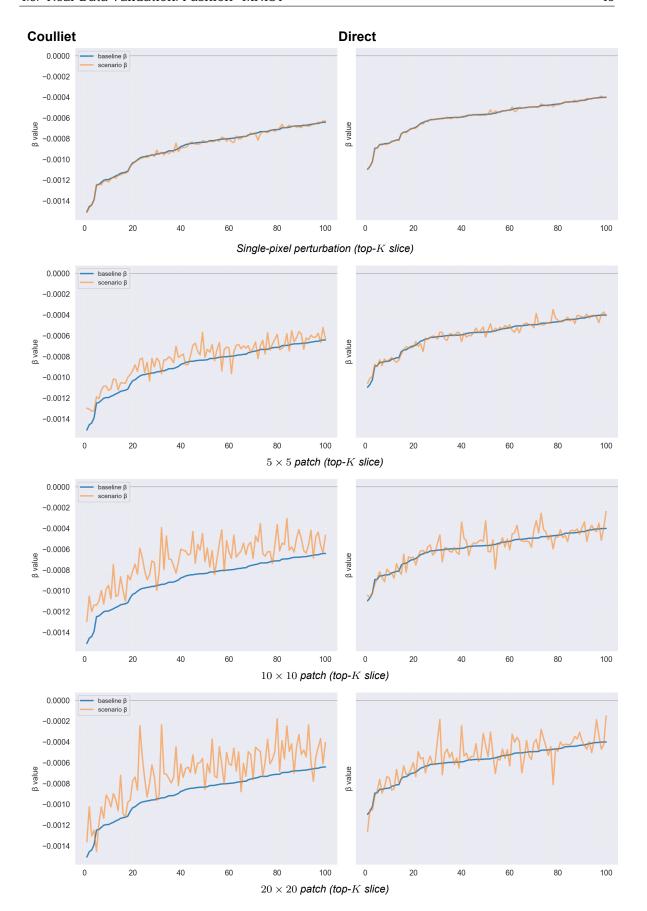


Figure 4.16: Bottom-tail-slice overlays of baseline vs. perturbed coefficients (bottom 100 beta values). Each row is a two-panel image with Coulliet (left) and Direct (right); the horizontal axis is the rank within the selected tail.

5

Discussion

The preceding chapters have established both the theoretical formulation and empirical validation of the proposed derivative-based regularization parameter selection rule for ridge regression in random-feature neural networks. Having presented the numerical and real-data results in Chapter 4, we now turn to a broader discussion of their implications.

This chapter aims to synthesize the key findings, interpret their meaning in light of the high - ridge regression framework. We further consider why the observed stability of the learned feature vectors β constitutes a practical advantage, in which contexts such stability may be desirable, and how it relates to interpretability and downstream reusability of the model. The discussion concludes with reflections on limitations, potential application domains, and directions for future research.

5.1. Summary of Findings

Let us start off by providing a short recap of the main experiments and their key findings once more.

5.1.1. Recap

The empirical evaluation conducted in Section 4.5 (Real Data Validation: Fashion-MNIST) provides a strong indication that the proposed derivative-based regularization parameter selection rule performs on par with, and in several aspects slightly better than, the deterministic-equivalent calibration method of Couillet [6].

On the binary Fashion-MNIST experiment, both calibration rules achieved nearly identical predictive performance, with test accuracies exceeding 99%. The derivative-based rule, however, yielded a marginally higher median test accuracy (99.56% vs. 99.30% for Couillet's rule) while exhibiting a reduced number of false negatives. This points to a slightly stronger recall and suggests that the proposed method does not sacrifice predictive power despite introducing an alternative criterion for γ selection.

More importantly, the perturbation analysis revealed a qualitative difference between the two approaches. Under small, controlled perturbations of the input data, models trained using the derivative-based calibration exhibited a notably more stable feature vector $\boldsymbol{\beta}$. In other words, the learned regression coefficients demonstrated lower sensitivity to localized variations in the training data, while remaining responsive to semantically meaningful structure within the input images. This stability implies a form of regularization that constrains model variance without suppressing relevant discriminative information.

Together, these findings highlight two key outcomes. First, the proposed method maintains competitive accuracy in a realistic, high-dimensional setting. Second, and perhaps more significantly, it yields solutions that are less sensitive to data perturbations, producing feature representations that are both robust and interpretable. This combination - predictive equivalence and enhanced stability - forms the basis for the broader discussion of the method's advantages, potential applications, and implications that follows in the subsequent sections.

5.1.2. Interpretations of the Results

The comparative analysis between the two calibration rules - Couillet's deterministic-equivalent criterion and the proposed derivative-based rule - reveals an interesting distinction not in terms of predictive performance, but in the qualitative behavior of the learned model parameters. Specifically, while both methods converge to similar predictive accuracies, the derivative-based direct loss calibration produces a more stable and consistent feature vector β under data perturbations.

This stability can be interpreted as a reduction in the variance component of the estimator. In high-dimensional ridge regression, the regularization parameter γ serves to balance the bias-variance trade-off: smaller values of γ allow the model to fit the data more closely but increase sensitivity to noise, whereas larger values enforce smoother, more regularized solutions. The proposed derivative-based rule tends to select slightly larger γ values compared to Couillet's approach (e.g., $\gamma_{\rm D}\approx 295$ versus $\gamma_{\rm C}\approx 71$ in the conducted Fashion-MNIST experiment), leading to a coefficient vector that varies less sharply with changes in the input data. Consequently, the resulting model exhibits reduced parameter variance, thereby enhancing robustness to perturbations and numerical instabilities without compromising generalization.

The perturbation experiments conducted in subsection 4.5.2 support this interpretation. When small localized changes were introduced to the training samples, the model calibrated using the derivative-based rule showed less adjustments in β , whereas Couillet's calibration produced noticeably more fluctuation in the same coefficients. Yet, this increased stability did not translate into underfitting: the model remained sensitive to meaningful regions of the input space (for instance, the discriminative edges in the Fashion-MNIST digits) and maintained comparable test performance. This suggests that the proposed rule enforces a selective form of regularization - one that suppresses sensitivity to noise and irrelevant variations while preserving responsiveness to the underlying signal. Nevertheless, it is important to take into account a rather minimalistic limited-complexity setting of the experiment - performing a ReLU-activated binary classification using ridge regression on fashion-MNIST dataset.

In summary, the results suggest that the main distinction between the two calibration strategies lies not in predictive accuracy but in their regularization dynamics. The derivative-based rule biases the estimator toward smoother, more reproducible solutions-a property that may become particularly advantageous in high-dimensional problems where interpretability, reproducibility, or robustness to measurement noise are of primary concern.

5.2. Advantages of the Proposed Derivative-Based Calibration over Couillet's Rule

In this subsection we will discuss the benefits that a more stable learned feature vector gives us, settings in which these benefits might be crucial and any potential real-world applications of the developed method.

5.2.1. Conceptual Advantages - Theory

Interpretability. In the context of random-feature ridge regression, "interpretability" should not be understood in the classical sense of direct feature attribution, since the random mapping $\sigma(WX)$ obscures the original input-space coordinates. Instead, interpretability arises in a relative sense: a model whose coefficients remain consistent under repeated experiments or mild perturbations can be more reliably analyzed, visualized, and reasoned about.

A stable β implies that the model's internal representation - the effective filter applied to the random features - is not an artifact of random initialization or particular noise realizations. This consistency makes the model's decision mechanism more transparent and trustworthy: one can meaningfully inspect the average structure of β across runs, or study which directions in the feature space are systematically emphasized. By contrast, if β changes substantially with each retraining, any attempt to interpret or visualize the learned representation becomes unreliable. Thus, stability becomes a prerequisite for interpretability in this high-dimensional random-feature setting.

Reusability of the Learned Coefficients. A related advantage concerns the potential for *down-stream reuse* of the learned coefficients β . In many applications, the ridge-regression layer serves as an intermediate component within a larger processing pipeline - for instance, as a feature extractor whose output is later used for clustering, transfer learning, or anomaly detection. When β is stable, it can be reused across slightly different datasets or experimental conditions without the need for retuning the regularization parameter or retraining from scratch. This property enables more reproducible and modular system design: the same model component can be integrated into different downstream tasks with predictable behavior.

In contrast, coefficients obtained through the Couillet calibration tend to vary more across retrainings, which limits their reuse beyond the specific dataset or random feature realization. The derivative-based rule, by reducing variance in β , implicitly supports the notion of a reusable "filter" that captures the core discriminative structure of the data rather than overfitting to local peculiarities.

5.2.2. Applied Advantages - Practice

Having established the conceptual benefits of stability in the preceding subsection, it is instructive to examine how this property manifests in practical, high-dimensional learning problems. In various real-world domains, researchers have identified coefficient or feature stability as a crucial component of model reliability, reproducibility, and interpretability. Stability ensures that the insights drawn from a model are not artifacts of random sampling or noise, but rather reflect persistent structure within the data itself. In this subsection, several representative studies are discussed to illustrate how the stability of learned representations, filters, or feature selections directly impacts the quality and trustworthiness of machine learning systems in applied contexts.

Feature Stability in High-Dimensional Learning. The empirical perspective on stability is thoroughly examined by Kalousis *et al.* [22], who conducted one of the earliest systematic analyses of feature selection stability across high-dimensional biological datasets. Their experiments covered several proteomics and microarray classification tasks, all characterized by a small number of samples (n < 100) and a very large number of features (p ranging from thousands to tens of thousands) - a

regime directly analogous to the high-dimensional random-feature setting studied in this thesis.

The authors evaluated a range of popular feature selection methods, including information gain, χ^2 , symmetrical uncertainty, RELIEF, and support-vector-based ranking techniques (SVM and SVM-RFE). To quantify the sensitivity of each method to sampling variations, they introduced three complementary similarity measures for feature preferences: (i) S_W for correlation of feature weights, (ii) S_R for rank consistency via Spearman's ρ , and (iii) S_S for overlap between top-k selected feature subsets. These measures collectively define a "stability profile" for each algorithm, summarizing how much the selected features change when the model is retrained on different resampled subsets of the same data.

The results were striking. Even among algorithms that achieved nearly identical classification accuracy, stability varied dramatically. For instance, on the colon cancer dataset ($p=2000,\ n=62$), SVM-RFE exhibited overlap values S_S below 0.25 for top-k=30 features, while simpler univariate methods such as information gain reached $S_S>0.9$ under the same setup. Similar patterns appeared across all five datasets: in the ovarian and leukemia studies, univariate filters consistently achieved overlap above 0.85, while multivariate SVM-based selectors fluctuated heavily with small data perturbations, dropping to as low as $S_S=0.1$. These results demonstrated that identical predictive scores can conceal highly unstable internal mechanisms.

From the perspective of this thesis, the analogy is clear. In the random-feature ridge regression studied here, two regularization calibration rules (Couillet's and the derivative-based one) may yield comparable prediction errors, yet differ sharply in the variance of their learned coefficients β . The derivative-based calibration thus plays a role analogous to preferring the more stable feature selector in Kalousis *et al.*: it favors solutions whose internal structure, the effective "filter" or feature weighting vector remains consistent under mild resampling or noise perturbations. Just as Kalousis and colleagues argued that stability is a prerequisite for reproducible and interpretable conclusions in biomedical settings, a stable β ensures reproducibility and trust in the learned representations in high-dimensional regression.

In short, the empirical message of Kalousis *et al.* reinforces the same principle that underlies this work: when dimensionality is high and data are scarce, predictive accuracy alone is not a sufficient measure of reliability. The internal stability of the learned coefficients becomes a distinguishing marker of models that capture genuine, reproducible structure rather than transient artifacts of random variation.

Stability and Redundancy in Neuroimaging Analysis. A second, domain-specific example of the importance of model stability is offered by Wang et~al.~ [49], who addressed the problem of stable feature selection in extremely high-dimensional functional MRI (fMRI) data. This example is particularly interesting for our case as it also works with images. In such neuroimaging tasks, the number of voxels or network connections ($p\sim10^3-10^4$) far exceeds the number of available subjects ($n\lesssim100$), and the features are often highly correlated or redundant. The study proposed an algorithm combining stability selection with the elastic~net, aiming not only to maintain high predictive accuracy but also to obtain reproducible and interpretable "biomarker" features that remain consistent across perturbations, label noise, and acquisition sites.

Wang *et al.* demonstrated that purely sparse methods such as ℓ_1 -regularized logistic regression or SVM yield unstable and overly sparse solutions: when applied to fMRI voxels, they select scattered and inconsistent sets of discriminative features, sensitive to even minor variations in the data. By integrating stability selection (which aggregates results over subsamples of both samples and features) with elastic-net regularization (which enforces group-wise feature retention), their model achieved markedly higher robustness and interpretability. For example, in the synthetic dataset with injected label noise (up to ten wrong labels), their method maintained a nearly constant voxel-selection accuracy, while alternatives such as the ℓ_1 -logistic model degraded sharply. Quantitatively, even with 10% label corruption, their approach preserved the correct identification of discriminative subregions, whereas univariate t-tests and standard SVMs produced false-positive activations or lost key regions entirely.

The advantages extended to real neuroimaging data. In a face-recognition fMRI experiment involving 26 subjects, Wang et al. successfully recovered five core regions implicated in visual and motor

aspects of facial processing - the occipital face area (OFA), fusiform face area (FFA), posterior superior temporal gyrus (pSTG), supplementary motor area (SMA), and sensorimotor cortex (SMC)- while competing methods either missed some of these regions or produced fragmented, spatially inconsistent activations. Furthermore, in a multi-center ADHD dataset where training and test data were collected at separate institutions (Peking University and New York University), the proposed method achieved a cross-center classification accuracy of 79.0% (AUC = 0.77), outperforming standard elastic net (72.6%) and randomized ℓ_1 -logistic models (67.7%). The ability to generalize across centers highlights the robustness of the learned feature weighting vector under substantial data variation.

The parallels to the present thesis are clear. In both cases, stability of the learned coefficients - whether interpreted as voxel importances in fMRI or as the regression filter β in random-feature ridge regression serves as a safeguard against spurious structure and noise-driven variability. Just as Wang et al. showed that stable, redundant feature selection improves both interpretability and cross-domain robustness, the derivative-based calibration proposed in this work yields a β that remains consistent across perturbations, enabling more reproducible downstream analyses. In settings where the data are high-dimensional, correlated, and noisy, such stability directly translates into both scientific trust-worthiness and operational reliability.

Stable Representation Learning in Industrial Fault Detection. The work of Michau and Fink [33] provides a particularly relevant example of how stability of internal representations translates into both interpretability and operational reliability in large-scale systems. Their study focuses on condition monitoring of industrial assets, where hundreds of heterogeneous sensors continuously record operating parameters under varying conditions. The challenge addressed is that only healthy-state data are typically available for training, while fault patterns are scarce or entirely absent. To handle this, the authors proposed an unsupervised feature learning and one-class classification architecture based on *Hierarchical Extreme Learning Machines* (HELM), integrating a stacked autoencoder for representation learning with a one-class classifier trained on the latent features.

The HELM framework can be seen as a random-feature architecture with analytical ridge-type weight estimation, conceptually close to the models investigated in this thesis. It learns a mapping from raw sensor space to a latent "health indicator," defined as the distance of the observation to the manifold of healthy data. Importantly, the system is trained with random input weights A, and its output coefficients β are determined in closed form through regularized least-squares (ridge or LASSO). This allows an explicit analysis of how the stability of β affects the consistency of the learned health indicators and, consequently, the reliability of fault detection decisions.

In their extensive experiments, Michau and Fink compared HELM with six alternative approaches—standalone one-class ELM and SVM classifiers, the same classifiers preceded by PCA, and a Deep Belief Network (DBN). On a simulated dataset with D=200 sensors and five injected fault types, HELM achieved the highest average accuracy (95% for n=5 intrinsic variables, $A_{\rm fault}\approx 0.95$) and strong magnification coefficients (Mag ≈ 88 –105), far exceeding the next-best competitors (SVM: $A\approx 0.94$, Mag ≈ 1.3 ; DBN: $A\approx 0.72$, Mag ≈ 1.2). These results quantitatively demonstrate that the feature-learning stage produces representations whose responses to perturbations—both random and structured—are highly consistent, yielding fault indicators that remain stable under varying operating conditions.

The findings are reinforced by the real-case application to a nine-month dataset from a hydrogen-cooled power plant generator with D=310 sensors. HELM reached an accuracy of 95% with zero false positives and a magnification coefficient of 20, while the closest competitor (SVM) achieved 87% accuracy and Mag =1.04. Importantly, HELM was able to detect early-stage degradation (day 169) months before the full short-circuit failure (day 247), providing a temporal margin that can translate directly into cost savings in maintenance and downtime prevention.

Conceptually, the relevance to this thesis lies in how HELM's stability emerges from a regularized random-feature mechanism. Like the derivative-based calibration introduced here, the HELM formulation stabilizes the learned β against small perturbations in the data, ensuring that the "health indicator" or model output evolves smoothly across retrainings. This stability, while mathematically grounded in

ridge and LASSO regularization, is functionally identical to the desired low-variance behavior of β in random-feature regression. In both contexts, the goal is not merely to predict correctly, but to ensure that the internal model structure—the coefficients themselves—remains reproducible and interpretable across runs, enabling engineers to trust the system's diagnosis or response over time and under changing conditions.

6

Conclusion

6.1. Summary and Final Remarks

This thesis introduced and analyzed a novel derivative-based, direct-loss calibration rule for selecting the regularization parameter in high-dimensional random-feature ridge regression. The method was developed as an alternative to the deterministic-equivalent canonical approach proposed by Coullet [6], offering a data-driven and theoretically grounded criterion derived directly from the loss-function dynamics.

The work began by revisiting the theoretical framework of ridge regression under random feature mappings, establishing the connection between the regularization parameter γ and the statistical biasvariance trade-off in high-dimensional settings. Building upon this foundation, a new calibration rule was formulated by enforcing the vanishing of the empirical loss derivative with respect to γ , yielding a direct, sample-based condition for optimal regularization.

The proposed method was validated through both synthetic and real-data experiments. In controlled simulations, the derivative-based rule achieved nearly identical predictive performance to Couillet's method calibration, confirming its theoretical consistency. Real-data evaluation on the Fashion-MNIST dataset further demonstrated that the proposed method matches or slightly surpasses Couillet's approach in terms of test accuracy, while producing a more stable coefficient vector $\boldsymbol{\beta}$ under perturbations of the training data. The perturbation analysis showed that this stability translates into smoother, more reproducible model behavior, which is an important property in high-dimensional problems where interpretability and robustness are as valuable as predictive accuracy.

A broader discussion connected these findings to practical contexts where model stability is a crucial requirement. By drawing parallels to studies on feature-selection stability in bioinformatics [22], neuroimaging [49], and industrial condition monitoring [33], the thesis highlighted that the benefits of coefficient stability extend far beyond synthetic benchmarks. In all such settings, the ability to learn consistent and reproducible representations enables more interpretable, trustworthy, and reusable models.

In summary, this research contributes both a new ridge regression hyperparameter calibration methodology and an expanded understanding of the role of stability in high-dimensional learning. The derivative-based rule provides a simple yet effective alternative to deterministic-equivalent approaches, combining strong empirical performance with enhanced reproducibility and interpretability of the learned coefficients.

6.2. Future Work 53

6.2. Future Work

The experimental and computational framework adopted in this thesis was deliberately kept minimal in order to isolate and examine the behavior of the proposed calibration rule under controlled conditions. While this design allowed for a clear theoretical and empirical comparison with Couillet's deterministic-equivalent method, it also leaves several directions open for future exploration.

First, the dimensional ratios employed in the simulations denoted as $c_1=p/n$ and $c_2=N/n$ were limited to values not exceeding approximately 1.5, primarily due to computational constraints on the DelftBlue supercomputer [1]. As a result, the synthetic experiments operated on moderate data sizes, sufficient for statistical validity but not for fully exploring the asymptotic regime where random matrix predictions become most distinctive. Extending these experiments to higher dimensional ratios and larger sample sizes would help verify whether the observed stability and generalization properties of the derivative-based rule persist under more extreme high-dimensional conditions.

Second, the real-data validation was restricted to a single dataset (Fashion-MNIST) and to a binary classification task involving two visually distinct classes. Future work could expand the empirical validation along several complementary axes. Additional real-world datasets. For instance, those from text, sensor, or biomedical domains - would allow assessing the robustness of the proposed rule across modalities. Similarly, multi-class classification scenarios could reveal how the calibration behaves when the decision boundaries interact in more complex ways.

Another natural extension concerns the activation function used in the random feature mapping. All experiments in this work employed the ReLU nonlinearity, chosen for its analytical simplicity and widespread use. Investigating alternative activations, such as leaky ReLU, tanh, or even randomized polynomial kernels, could uncover how the choice of nonlinearity influences the derivative-based calibration and its resulting stability properties.

Finally, while the present study focused on classification, future research could revisit the same framework in a genuine regression context. Testing the proposed calibration rule on real-valued target problems would help establish whether its stabilizing effect on the learned coefficients β extends beyond classification tasks and remains beneficial in continuous-output ridge regression. Identifying suitable high-dimensional regression datasets would therefore be an important next step toward a more complete empirical characterization of the method.

In summary, future work should aim to broaden both the scale and scope of experimentation - ncreasing sample and feature counts, diversifying nonlinear mappings, and extending to regression settings to consolidate the findings of this thesis and to fully explore the potential of the derivative-based regularization rule in modern high-dimensional learning.

References

- [1] Delft High Performance Computing Centre (DHPC). *DelftBlue Supercomputer (Phase 2)*. https://www.tudelft.nl/dhpc/ark:/44463/DelftBluePhase2. 2024.
- [2] Zhidong Bai and Jack W. Silverstein. *Spectral Analysis of Large Dimensional Random Matrices*. 2nd. Springer, 2010. DOI: https://doi.org/10.1007/978-1-4419-0661-8.
- [3] Mikhail Belkin et al. "Reconciling Modern Machine Learning Practice and the Classical Bias—Variance Trade-off". In: *Proceedings of the National Academy of Sciences* 116.32 (2019), pp. 15849—15854. DOI: https://doi.org/10.1073/pnas.1903070116.
- [4] Youngmin Cho and Lawrence K. Saul. "Kernel Methods for Deep Learning". In: Advances in Neural Information Processing Systems (NeurIPS). Vol. 22. 2009, pp. 342–350. URL: https://proceedings.neurips.cc/paper_files/paper/2009/hash/013a006f03dbc5392effeb8f18fda755-Abstract.html.
- [5] Romain Couillet and Mérouane Debbah. *Random Matrix Methods for Wireless Communications*. Cambridge University Press, 2011. DOI: https://doi.org/10.1017/CB09780511994746.
- [6] Romain Couillet and Zhenyu Liao. "Random Matrix Methods for Machine Learning". In: Cambridge University Press, 2022, pp. 299–306. DOI: https://doi.org/10.1017/9781009128490.
- [7] Romain Couillet and Zhenyu Liao. *Random Matrix Methods for Machine Learning*. Cambridge University Press, 2022. DOI: https://doi.org/10.1017/9781009128490.
- [8] George Cybenko. "Approximation by Superpositions of a Sigmoidal Function". In: Mathematics of Control, Signals, and Systems 2.4 (1989), pp. 303–314. DOI: https://doi.org/10.1007/ BF02551274.
- [9] Jia Deng et al. "ImageNet: A Large-Scale Hierarchical Image Database". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2009), pp. 248–255. DOI: https://doi.org/10.1109/CVPR.2009.5206848.
- [10] Edgar Dobriban and Weijie Su. "High-Dimensional Asymptotics of Prediction: Ridge Regression and Classification". In: *The Annals of Statistics* 46.1 (2018), pp. 247–279. DOI: https://doi.org/10.1214/17-A0S1549.
- [11] Jeffrey L. Elman. "Finding Structure in Time". In: *Cognitive Science* 14.2 (1990), pp. 179–211. DOI: https://doi.org/10.1207/s15516709cog1402_1.
- [12] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. Statistical Learning with Sparsity: The Lasso and Generalizations. Chapman and Hall/CRC, 2019. DOI: https://doi.org/10.1201/9780429447273.
- [13] Trevor Hastie et al. "Surprises in High-Dimensional Ridgeless Least Squares Interpolation". In: *The Annals of Statistics* 50.2 (2022), pp. 949–986. DOI: https://doi.org/10.1214/21-A0S2133.
- [14] Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2016, pp. 770–778. DOI: https://doi.org/10.1109/CVPR.2016.90.
- [15] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. "A Fast Learning Algorithm for Deep Belief Nets". In: *Neural Computation* 18.7 (2006), pp. 1527–1554. DOI: https://doi.org/10.1162/neco.2006.18.7.1527.
- [16] Arthur E. Hoerl and Robert W. Kennard. "Ridge Regression: Biased Estimation for Nonorthogonal Problems". In: *Technometrics* 12.1 (1970), pp. 55–67. DOI: https://doi.org/10.1080/00401706.1970.10488634.

References 55

[17] John J. Hopfield. "Neural Networks and Physical Systems with Emergent Collective Computational Abilities". In: *Proceedings of the National Academy of Sciences* 79.8 (1982), pp. 2554–2558. DOI: https://doi.org/10.1073/pnas.79.8.2554.

- [18] Kurt Hornik. "Approximation Capabilities of Multilayer Feedforward Networks". In: *Neural Networks* 4.2 (1991), pp. 251–257. DOI: https://doi.org/10.1016/0893-6080(91)90009-T.
- [19] Guang-Bin Huang, Qin-Yu Zhu, and Chee-Kheong Siew. "Extreme Learning Machine: Theory and Applications". In: *Neurocomputing* 70.1–3 (2006), pp. 489–501. DOI: https://doi.org/10.1016/j.neucom.2005.12.126.
- [20] Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural Tangent Kernel: Convergence and Generalization in Neural Networks". In: *Advances in Neural Information Processing Systems* (NeurlPS) 31 (2018). URL: https://arxiv.org/abs/1806.07572.
- [21] Herbert Jaeger. "The "Echo State" Approach to Analysing and Training Recurrent Neural Networks". In: *GMD Report* 148 (2001). URL: https://www.researchgate.net/publication/215385037.
- [22] Alexandros Kalousis, Julien Prados, and Melanie Hilario. "Stability of Feature Selection Algorithms: A Study on High-Dimensional Spaces". In: *Knowledge and Information Systems* 12.1 (2007), pp. 95–116. DOI: https://doi.org/10.1007/s10115-006-0040-8. URL: https://link.springer.com/article/10.1007/s10115-006-0040-8.
- [23] Teuvo Kohonen. "Self-Organized Formation of Topologically Correct Feature Maps". In: *Biological Cybernetics* 43.1 (1982), pp. 59–69. DOI: https://doi.org/10.1007/BF00337288.
- [24] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. "ImageNet Classification with Deep Convolutional Neural Networks". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 25. 2012, pp. 1097–1105. URL: https://proceedings.neurips.cc/paper_files/paper/2012/hash/c399862d3b9d6b76c8436e924a68c45b-Abstract.html.
- [25] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. "Deep Learning". In: *Nature* 521.7553 (2015), pp. 436–444. DOI: https://doi.org/10.1038/nature14539.
- [26] Yann LeCun et al. "Backpropagation Applied to Handwritten Zip Code Recognition". In: *Neural Computation* 1.4 (1989), pp. 541–551. DOI: https://doi.org/10.1162/neco.1989.1.4.541.
- [27] Yann LeCun et al. "Gradient-Based Learning Applied to Document Recognition". In: *Proceedings of the IEEE* 86.11 (1998), pp. 2278–2324. DOI: https://doi.org/10.1109/5.726791.
- [28] Cosme Louart, Zhenyu Liao, and Romain Couillet. "A Random Matrix Approach to Neural Networks". In: *The Annals of Applied Probability* 28.2 (2018), pp. 1190–1248. DOI: https://doi.org/10.1214/17-AAP1328.
- [29] Cosme Louart, Zhenyu Liao, and Romain Couillet. "A random matrix approach to neural net-works". In: The Annals of Applied Probability (2018). DOI: https://doi.org/10.1214/17-AAP1328.
- [30] Wolfgang Maass, Thomas Natschläger, and Henry Markram. "Real-Time Computing Without Stable States: A New Framework for Neural Computation Based on Perturbations". In: *Neural Computation* 14.11 (2002), pp. 2531–2560. DOI: https://doi.org/10.1162/089976602760407955.
- [31] Vladimir A. Marchenko and Leonid A. Pastur. "Distribution of Eigenvalues for Some Sets of Random Matrices". In: *Mathematics of the USSR-Sbornik* 1.4 (1967), pp. 457–483. DOI: https://doi.org/10.1070/SM1967v001n04ABEH001994.
- [32] Song Mei and Andrea Montanari. "The Generalization Error of Random Features Regression: Precise Asymptotics and Double Descent Curve". In: *Communications on Pure and Applied Mathematics* 75.4 (2022), pp. 667–766. DOI: https://doi.org/10.1002/cpa.21937.
- [33] Gabriel Michau and Olga Fink. "Feature Learning for Fault Detection in High-Dimensional Condition-Monitoring Signals". In: arXiv preprint arXiv:1810.05550 (2019). URL: https://arxiv.org/abs/1810.05550.
- [34] Marvin Minsky and Seymour Papert. *Perceptrons: An Introduction to Computational Geometry*. Cambridge, MA: MIT Press, 1969.

References 56

[35] Vinod Nair and Geoffrey E. Hinton. "Rectified Linear Units Improve Restricted Boltzmann Machines". In: *Proceedings of the 27th International Conference on Machine Learning (ICML)*. 2010, pp. 807–814. URL: https://icml.cc/Conferences/2010/papers/432.pdf.

- [36] Radford M. Neal. *Bayesian Learning for Neural Networks*. Vol. 118. Lecture Notes in Statistics. Springer, 1996. DOI: https://doi.org/10.1007/978-1-4612-0745-0.
- [37] Ian Osband et al. "Deep Exploration via Bootstrapped DQN". In: Advances in Neural Information Processing Systems (NeurIPS). Vol. 29. 2016, pp. 4026–4034. URL: https://proceedings.neurips.cc/paper_files/paper/2016/hash/8d8818c8ce75b17d62e3e0fbde7d0b5c-Abstract.html.
- [38] Jeffrey Pennington and Pratik Worah. "Nonlinear Random Matrix Theory for Deep Learning". In: Advances in Neural Information Processing Systems (NeurIPS). Vol. 30. 2017, pp. 2637–2646. URL: https://proceedings.neurips.cc/paper_files/paper/2017/hash/6070ef0e056174b4 a3d8d005a36cf3c9-Abstract.html.
- [39] Ali Rahimi and Benjamin Recht. "Random Features for Large-Scale Kernel Machines". In: Advances in Neural Information Processing Systems (NeurIPS). Vol. 20. 2008, pp. 1177–1184. URL: https://proceedings.neurips.cc/paper_files/paper/2007/hash/013a006f03dbc5392effe b8f18fda755-Abstract.html.
- [40] Ali Rahimi and Benjamin Recht. "Weighted Sums of Random Kitchen Sinks: Replacing Minimization with Randomization in Learning". In: Advances in Neural Information Processing Systems (NeurIPS). Vol. 21. 2008, pp. 1313–1320. URL: https://proceedings.neurips.cc/paper_files/paper/2008/hash/0efe32849d230d7f53049ddc4a4b0c60-Abstract.html.
- [41] Frank Rosenblatt. "The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain". In: *Psychological Review* 65.6 (1958), pp. 386–408. DOI: https://doi.org/10.1037/h0042519.
- [42] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. "Learning Representations by Back-Propagating Errors". In: *Nature* 323.6088 (1986), pp. 533–536. DOI: https://doi.org/10.1038/323533a0.
- [43] Andrew Saxe et al. "On Random Weights and Unsupervised Feature Learning". In: *Proceedings of the 28th International Conference on Machine Learning (ICML)*. 2011, pp. 1089–1096. URL: https://icml.cc/Conferences/2011/papers/438_icmlpaper.pdf.
- [44] Karen Simonyan and Andrew Zisserman. "Very Deep Convolutional Networks for Large-Scale Image Recognition". In: *International Conference on Learning Representations (ICLR)* (2015). URL: https://arxiv.org/abs/1409.1556.
- [45] Nitish Srivastava et al. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". In: *Journal of Machine Learning Research* 15.56 (2014), pp. 1929–1958. URL: https://jmlr.org/papers/v15/srivastava14a.html.
- [46] Andrey N. Tikhonov and Vasiliy Y. Arsenin. *Solution of Ill-Posed Problems*. Originally published in Russian, 1963. Winston and Sons, 1977.
- [47] Ashish Vaswani et al. "Attention Is All You Need". In: *Advances in Neural Information Processing Systems (NeurIPS)*. Vol. 30. 2017. URL: https://arxiv.org/abs/1706.03762.
- [48] Pascal Vincent et al. "Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion". In: *Journal of Machine Learning Research* 11 (2010), pp. 3371–3408. URL: https://www.jmlr.org/papers/v11/vincent10a.html.
- [49] Yilun Wang et al. "A Novel Approach for Stable Selection of Informative Redundant Features from High Dimensional fMRI Data". In: *IEEE Transactions on Neural Systems and Rehabilitation Engineering* 24.11 (2015), pp. 1236–1248. DOI: https://doi.org/10.1109/TNSRE.2015. 2474015. URL: https://arxiv.org/abs/1506.08301.
- [50] Eugene P. Wigner. "Characteristic Vectors of Bordered Matrices with Infinite Dimensions". In: *Annals of Mathematics* 62.3 (1955), pp. 548–564. DOI: https://doi.org/10.2307/1970079.



Source Code

A.1. Mean Squared Error against the regularization parameter γ

```
1 import numpy as np
2 import random
3 from scipy import linalg
4 from mpi4py import MPI
6 comm = MPI.COMM_WORLD
7 nprocs = comm.Get_size()
8 myrank = comm.Get_rank()
10
11 def sigma(t):
      Small sigma function of choice
13
14
     :param t: input
      :return: output
15
16
     if activation_function == 'linear':
18
19
     if activation_function == 'ReLu':
          return np.maximum(t, 0)
21
22
     if activation_function == 'sign':
         return np.sign(t)
24
25
27
28 def K(x, y):
29
30
      \label{tensor} \mbox{Kernel function. Depends on the choice of a small sigma function}
      :param x: first 'point' set
:param y: second 'point' set
32
      :return: matrix of measures of 'distances' between points
34
     if activation_function == 'linear':
35
          return x.T@y
37
      if activation_function == 'ReLu':
38
          norm_x = np.linalg.norm(x, axis=0) # Shape (n_x,)
           norm_y = np.linalg.norm(y, axis=0) # Shape (n_y,)
40
41
           xTy = x.T @ y # Shape (n_x, n_y)
42
43
           norm_prod = norm_x[:, np.newaxis] * norm_y[np.newaxis, :] # Shape (n_x, n_y)
44
           norm_prod = np.maximum(norm_prod, 1e-10) # Avoid division by zero
45
```

```
46
           cos_theta = xTy / norm_prod
47
           cos_theta = np.clip(cos_theta, -1 + 1e-10, 1 - 1e-10) # Clamp values to [-1+1e-10,
48
               1-1e-10]
49
           theta = np.arccos(-cos\_theta) # Shape (n_x, n_y)
50
           sin_theta = np.sqrt(1 - cos_theta ** 2)
51
52
           return (norm_prod) / (2 * np.pi) * (cos_theta * theta + sin_theta)
53
54
      if activation_function == 'sign':
55
56
           norm_x = np.linalg.norm(x, axis=0) # Shape (n_x,)
           norm_y = np.linalg.norm(y, axis=0) # Shape (n_y,)
57
58
           xTy = x.T @ y # Shape (n_x, n_y)
59
60
           norm_prod = np.outer(norm_x, norm_y) # Shape (n_x, n_y)
61
62
           norm_prod = np.maximum(norm_prod, 1e-10) # Avoid division by zero
63
          cos_theta = xTy / norm_prod
           cos_theta = np.clip(cos_theta, -1 + 1e-10, 1 - 1e-10) # Clamp values to [-1+1e-10,
65
               1-1e-10]
           return (2 / np.pi) * np.arcsin(cos_theta)
67
68
70 def K_(x, y, delta, N, n_train):
       Function to compute kernel approximation
72
73
       :param x: first 'point' set
74
       :param y: second 'point' set
       :param delta: delta parameter
75
76
       :return: approximation for the kernel matrix
77
78
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
      return k_
79
80
81
82 def Etrain(gamma, data, npN, monte_carlo_loops=30):
83
84
       Actual training error compute (random-based)
       :param gamma: ridge penalty value
85
       :param data: train/test data
86
87
       :param npN: dimensionality of the data and num of neurons
       :param monte_carlo_loops: number of iterations for E value averaging for different
88
           generated W
       :return: Training error
90
       # Unpacking variables
91
92
       X_train, X_test, Y_train, Y_test = data
       n_{train}, n_{test}, p, N = npN
93
      E_train_arr = []
95
       p = X_train.shape[0]
96
      for i in range(monte_carlo_loops):
98
99
           # W = np.random.randn(N, p)
100
           # Sigm = sigma(W@X_train)
101
           # Q_y = np.linalg.inv((1/n_train)*Sigm.T@Sigm + gamma*np.eye(n_train))
102
103
           # E_train = (gamma*gamma/n_train)*Y_train@np.linalg.matrix_power(Q_y, 2)@Y_train.T
104
           # E_train_arr.append(E_train)
105
106
107
           W = np.random.randn(N, p)
           Sigm = sigma(W @ X_train)
Sigm_ = sigma(W @ X_test)
108
109
110
111
           inv_tQ_r = linalg.solve(Sigm.T @ Sigm / n_train + gamma * np.eye(n_train), Y_train)
           beta = Sigm / n_train @ inv_tQ_r
112
```

```
E_train = np.linalg.norm(Y_train-Sigm.T@beta)**2/n_train
114
                          E_train_arr.append(E_train)
115
116
                return np.mean(np.array(E_train_arr))
118
119
def Etest(gamma, data, npN, monte_carlo_loops=30):
121
122
                Actual test error compute (random-based)
123
                :param gamma: ridge penalty value
                :param data: train/test data
124
125
                 :param npN: dimensionality of the data and number of neurons
                :param monte_carlo_loops: number of iterations for E value averaging for different
126
                          generated W
127
                :return: Test error
128
                # Unpacking variables
129
                X_train, X_test, Y_train, Y_test = data
130
                n_{train}, n_{test}, p, N = npN
131
132
133
               E test arr = []
134
                p = X_train.shape[0]
135
               for i in range(monte_carlo_loops):
136
137
                          W = np.random.randn(N, p)
138
                          Sigm = sigma(W@X_train)
Sigm_ = sigma(W@X_test)
139
140
141
                          inv_tQ_r = linalg.solve(Sigm.T@Sigm/n_train + gamma * np.eye(n_train), Y_train)
142
143
                          beta = Sigm/n_train @ inv_tQ_r
144
                          \# Q_y = np.linalg.inv((1 / n_train) * Sigm.T @ Sigm + gamma * np.eye(n_train))
145
146
                         # term1 = (1/n_{test})*(Y_{test}QY_{test}.T)
147
                          # term2 = (2/(n_train*n_test))*(Y_train@Q_y@Sigm.T@Sigm_@Y_test.T)
                           \# \ \text{term3} = (1/(n_{\text{train}}**2*n_{\text{test}}))*(Y_{\text{train}}Q_{\text{y}}GSigm.TGSigm_GSigm_.TGSigm}Q_{\text{y}}QY_{\text{train}}. 
149
                                    T)
151
                          E_test = np.linalg.norm(Y_test-Sigm_.T@beta)**2/n_test
152
                          E_test_arr.append(E_test)
153
154
155
                ans = np.mean(np.array(E_test_arr))
156
157
                return ans
159
def find_delta(gamma, X_train, N, accuracy) -> float:
161
                Helper-function that finds delta parameter for the resolvent {\tt Q} iteratively
162
                :param gamma: ridge penalty value
163
                 :param X_train: training set
164
                :param N: number of neurons in the hidden layer
165
                :param accuracy: accuracy for numerical delta finding
                :return: optimal value for delta parameter for the resolvent {\tt Q}
167
168
                n_train = X_train.shape[1]
169
                delta_prev = 1
170
171
                delta_next = 0
                while abs(delta_prev-delta_next) > accuracy:
172
                          delta_prev = delta_next
173
                          \label{eq:Q_section} Q_{-} = \text{np.linalg.inv}((N/n\_\text{train})*(K(X\_\text{train}, X\_\text{train})/(1+\text{delta\_next})) + \text{gamma*np.eye}((N/n\_\text{train})*(N/n\_\text{train})) + (N/n\_\text{train})*(N/n\_\text{train}) + (N/n\_\text{train}) + (N
174
                                   n train))
175
                          delta_next = (1/n_train)*(np.trace(Q_@K(X_train, X_train)))
176
                return delta_next
177
179
180 def Etrain_(gamma, data, npN):
```

```
Estimated train error compute (expectation based, deterministic)
182
                  :param gamma: ridge penalty value
183
                  :param data: train/test data
184
                  :param \operatorname{npN}: dimensionality of the data and number of Neurons
                  :return: estimated training error
186
187
                  # Unpacking variables:
188
                  X_{train}, X_{test}, Y_{train}, Y_{test} = data
189
                  n_train, n_test, p, N = npN
190
191
192
                  delta = find_delta(gamma, X_train, N, delta_accuracy)
193
                  \label{eq:Q_section} Q_{-} = \text{np.linalg.inv}((\text{N/n\_train})*(\text{K(X\_train, X\_train})/(1+\text{delta})) + \text{gamma*np.eye}(\text{n\_train}))
                  K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
194
195
                  \#E_{train} = ((gamma**2)/n_train)*(Y_train)Q_0((((1/N)*np.matrix.trace(Q_0K_0Q_))/((1 - 1/N)*np.matrix))
196
                             \label{eq:normalized} N)*np.matrix.trace(K_@Q_@K_@Q_)))*K_ + np.eye(n_train))@Q_@Y_train.T)
                   E_{train} = ((gamma**2)/n_train)*(Y_{train}Q_Q(((((1/N)*np.trace(Q_QK_QQ_))/((1 - 1/N)*np.trace(Q_QK_QQ_))/((1 - 1/N)*np.trace(Q_QK_QQ_QK_QQ_))/((1 - 1/N)*np.trace(Q_QK_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_
197
                             \label{eq:trace} trace(K_@Q_@K_@Q_)))*K_ + np.eye(n_train))@Q_@Y_train.T)
198
                  return E_train_
200
201
202 def Etest_(gamma, data, npN):
203
204
                  Estimated test error compute (expectation based, deterministic)
                  :param gamma: ridge penalty value
205
206
                  :param data: train/test
                  :param npN: dimensionality of the data and number of Neurons
207
                  :return: estimated test error
208
209
210
                  # Unpacking variables:
                  X_{train}, X_{test}, Y_{train}, Y_{test} = data
211
212
                  n_{train}, n_{test}, p, N = npN
213
214
                  delta = 0.0
                  delta = find_delta(gamma, X_train, N, delta_accuracy)
                  Q_ = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta)) + gamma*np.eye(n_train))
216
                  K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
217
                  K_xX = (N/n_{train})*(K(X_{train}, X_{test})/(1+delta))
                  K_XX = (N/n_{train})*(K(X_{test}, X_{test})/(1+delta))
219
220
                  #E_test_ = (1/n_test)*np.sum((Y_test.T - K_xX.T@Q_@Y_train.T)**2) + (((1/N)*(
                             \texttt{matrix.trace}(\texttt{K}_{\texttt{X}}\texttt{X}\texttt{X}) - (1/\texttt{n}_{\texttt{test}}) * \texttt{np.matrix.trace}((\texttt{np.eye}(\texttt{n}_{\texttt{train}}) + \texttt{gamma*Q}_{\texttt{Q}}) \texttt{Q}(
                             K_xX@K_xX.T@Q_)))
                  222
                              Y_{trainQQ_0K_0Q_0Y_train.T))/((1-1/N)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0
                              \texttt{K\_XX)} \ - \ (1/\texttt{n\_test}) * \texttt{np.trace}((\texttt{np.eye}(\texttt{n\_train}) \ + \ \texttt{gamma*Q\_}) @(\texttt{K\_xX@K\_xX.T@Q\_}))) 
223
224
                  return E_test_
225
def generate_synthetic_data(n, c1, c2, density=0.35, noise_level=0.1):
                  p = int(c1 * n)
228
                  N = int(c2 * n)
                  n_{train} = int(n * 0.8)
230
                  n_test = n - n_train
231
232
233
                  # Generate random weights and sparse beta
                  W = np.random.randn(N, p)
234
                  b = np.random.choice([0, 1], size=(N, 1), p=[1 - density, density])
235
236
237
                  X_train = np.random.randn(p, n_train)
238
239
                  Y_train = (sigma(W @ X_train).T @ b).flatten() + np.random.normal(0, noise_level, n_train
                  X_test = np.random.randn(p, n_test)
240
                  Y_test = (sigma(W @ X_test).T @ b).flatten() + np.random.normal(0, noise_level, n_test)
241
242
                  return X_train, X_test, Y_train, Y_test, W, b
243
```

```
245
246 def save_results_to_file(gammas, E_test_arr, E_train_arr, E_test_bar_arr, E_train_bar_arr,
                    log_filename, result_filename):
247
                    A function designed for HPC to save the results in proper readable and plotable form.
248
249
                    :param E_test_arr:
250
                    :param E_train_arr:
251
252
                     :param E_test_bar_arr:
253
                    :param E_train_bar_arr:
254
                    :param log_filename:
255
                    :param result_filename:
                    :return: none
256
257
                    with open(log_filename, 'w') as log_file:
258
                                for gamma, E_test, E_train, E_test_, E_train_ in zip(gammas, E_test_arr, E_train_arr,
259
                                               E_test_bar_arr,
                                                                                                                                                                                            E_train_bar_arr):
260
261
                                            log_file.write(
                                                         \texttt{f"for}_{\sqcup} \texttt{gamma}_{\sqcup} = _{\sqcup} \{\texttt{gamma} : .6e\}, _{\sqcup} \\ \texttt{E\_train}_{\sqcup} = _{\sqcup} \{\texttt{E\_train} : .6e\}, _{\sqcup} \\ \texttt{E\_test}_{\sqcup} = _{\sqcup} \{\texttt{E\_test} : .6e\}, _{\sqcup} \\ \texttt{E\_test}_{\sqcup} = _{\sqcup} \{\texttt{E\_test}_{\sqcup} : .6e\}, _{\sqcup} \\ \texttt{
                                                                    E_{\text{train}} = \{E_{\text{train}} : .6e\}, E_{\text{test}} = \{E_{\text{test}} : .6e\} \setminus n"\}
263
                    results = {
                                "gammas": gammas,
265
                                 "E_test_arr": E_test_arr,
266
                                "E_train_arr": E_train_arr,
267
                                "E_test_bar_arr": E_test_bar_arr,
268
                                "E_train_bar_arr": E_train_bar_arr
269
270
                    np.savez(result_filename, **results)
271
272
273
274 activation_function = 'ReLu'
276 delta_accuracy = 1e-4
277 b_vector_density = 0.15 #%
278
279 n = 5000
c1 = 1.0 \# p/n
281 c2 = 0.8 \# N/n
282
p = int(c1 * n)
284 N = int(c2 * n)
n_{train} = int(n * 0.8)
n_{test} = n - n_{train}
287
288 npN = n_train, n_test, p, N
b_vector_density)
290 data = X_train, X_test, Y_train, Y_test
# data = get_data_MNIST_binary(npN)
292 gammas = [10**y for y in np.arange(-6, 6, 0.1)]
294 gamma_chunks = np.array_split(gammas, nprocs)
295 my_gammas = gamma_chunks[myrank]
296
297 E_test_arr = []
298 E_train_arr = []
299 E_test_bar_arr = []
300 E_train_bar_arr = []
301 log_entries = []
302
303 for g in my_gammas:
                    E_test = Etest(g, data, npN)
304
305
                    E_test_arr.append(E_test)
306
                    E_test_ = Etest_(g, data, npN)
307
                    E_test_bar_arr.append(E_test_)
308
309
                    E_train = Etrain(g, data, npN)
310
311     E_train_arr.append(E_train)
```

```
312
                      E_train_ = Etrain_(g, data, npN)
313
314
                      E_train_bar_arr.append(E_train_)
                     log_entries.append(
316
317
                                   f"for \_ gamma \_ = \_ \{g\} \_ the \_ error \_ values \_ are : \_ E\_ train = \{E\_ train\}, \_ E\_ test = \{E\_ test\}, \_ E\_ train \_ \{g\}, \_ E\_ train \_ \{
                                                ={E_train_},_E_test_={E_test_}\n")
318
319 E_test_arr = comm.gather(E_test_arr, root=0)
320 E_train_arr = comm.gather(E_train_arr, root=0)
321 E_test_bar_arr = comm.gather(E_test_bar_arr, root=0)
322 E_train_bar_arr = comm.gather(E_train_bar_arr, root=0)
323 log_entries = comm.gather(log_entries, root=0)
324
325 if myrank == 0:
                      E_test_arr = [item for sublist in E_test_arr for item in sublist]
326
                      E_train_arr = [item for sublist in E_train_arr for item in sublist]
327
328
                      E_test_bar_arr = [item for sublist in E_test_bar_arr for item in sublist]
                      E_train_bar_arr = [item for sublist in E_train_bar_arr for item in sublist]
329
                      log_entries = [item for sublist in log_entries for item in sublist]
331
                      np.savez("results.npz",
332
                                                  gammas=gammas,
333
                                                   E_test_arr=E_test_arr,
334
335
                                                  E_train_arr=E_train_arr,
                                                  E_test_bar_arr=E_test_bar_arr,
336
337
                                                  E_train_bar_arr=E_train_bar_arr,
                                                  {\tt activation\_function=activation\_function}\,,
                                                  density=b_vector_density,
339
                                                  n_train=n_train,
340
                                                  n_test=n_test,
                                                  N=N,
342
343
                                                  p=p)
344
                      with open("log.txt", "w") as f:
345
                                  f.writelines(log_entries)
```

A.2. Optimal γ value and error behaviour with increasing amount of data

```
1 import numpy as np
2 import random
3 from scipy import linalg
4 from mpi4py import MPI
5 import os
6 os.environ["OMP_NUM_THREADS"] = "16"
8 comm = MPI.COMM_WORLD
9 nprocs = comm.Get_size()
10 myrank = comm.Get_rank()
12
13 def sigma(t):
14
15
       Small sigma function of choice
      :param t: input
16
      :return: output
17
18
      return t
19
20
21
22 def K(x, y):
      Kernel function. Depends on the choice of a small sigma function :param x\colon first 'point' set
24
25
     :param y: second 'point' set
     :return: matrix of measures of 'distances' between points
27
28
     return x.T@y
30
32 def K_{x}, y, delta, N, n_{train}:
33
      Function to compute kernel approximation
      :param x: first 'point' set
35
36
      :param y: second 'point' set
      :param delta: delta parameter
37
      :return: approximation for the kernel matrix
38
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
40
41
      return k_
43
44 def Etrain(gamma, data, npN, monte_carlo_loops=30):
45
      Actual training error compute (random-based)
46
47
       :param gamma: ridge penalty value
      :param data: train/test data
48
49
      :param \operatorname{npN}: dimensionality of the data and \operatorname{num} of neurons
      :param monte_carlo_loops: number of iterations for E value averaging for different
          generated W
51
      :return: Training error
52
      # Unpacking variables
53
      X_train, X_test, Y_train, Y_test = data
      n_train, n_test, p, N = npN
55
56
     E_train_arr = []
58
     p = X_{train.shape[0]}
59
      for i in range(monte_carlo_loops):
61
62
           # W = np.random.randn(N, p)
           # Sigm = sigma(W@X_train)
63
           # Q_y = np.linalg.inv((1/n_train)*Sigm.T@Sigm + gamma*np.eye(n_train))
64
           # E_train = (gamma*gamma/n_train)*Y_train@np.linalg.matrix_power(Q_y, 2)@Y_train.T
```

```
# E_train_arr.append(E_train)
67
68
69
           W = np.random.randn(N, p)
           Sigm = sigma(W @ X_train)
           Sigm_ = sigma(W @ X_test)
71
72
           inv_tQ_r = linalg.solve(Sigm.T @ Sigm / n_train + gamma * np.eye(n_train), Y_train)
73
           beta = Sigm / n_train @ inv_tQ_r
74
75
76
           E_train = np.linalg.norm(Y_train-Sigm.T@beta)**2/n_train
77
           E_train_arr.append(E_train)
78
       return np.mean(np.array(E_train_arr))
79
80
81
82 def Etest(gamma, data, npN, monte_carlo_loops=30):
83
84
       Actual test error compute (random-based)
85
       :param gamma: ridge penalty value
       :param data: train/test data
       :param \operatorname{np} \mathbb{N}\colon \operatorname{dimensionality} of the data and number of neurons
87
       : \verb"param monte_carlo_loops: number of iterations for E value averaging for different
88
           generated W
       :return: Test error
89
90
       # Unpacking variables
91
92
       X_{train}, X_{test}, Y_{train}, Y_{test} = data
       n_{train}, n_{test}, p, N = npN
93
94
       E_test_arr = []
95
       p = X_train.shape[0]
97
98
       for i in range(monte_carlo_loops):
99
100
           W = np.random.randn(N, p)
           Sigm = sigma(W@X_train)
           Sigm_ = sigma(W@X_test)
102
103
           inv_tQ_r = linalg.solve(Sigm.T@Sigm/n_train + gamma * np.eye(n_train), Y_train)
           beta = Sigm/n_train @ inv_tQ_r
105
106
           # Q_y = np.linalg.inv((1 / n_train) * Sigm.T @ Sigm + gamma * np.eye(n_train))
107
108
109
           # term1 = (1/n_test)*(Y_test@Y_test.T)
           # term2 = (2/(n_train*n_test))*(Y_train@Q_y@Sigm.T@Sigm_@Y_test.T)
110
            \texttt{\# term3} = (1/(n\_train**2*n\_test))*(Y\_train@Q\_y@Sigm.T@Sigm\_@Sigm\_.T@SigmQQ\_y@Y\_train.
111
112
113
114
           E_test = np.linalg.norm(Y_test-Sigm_.T@beta)**2/n_test
           E_test_arr.append(E_test)
115
       ans = np.mean(np.array(E_test_arr))
117
118
       return ans
119
120
121
def find_delta(gamma, X_train, N, accuracy) -> float:
123
124
       Helper-function that finds delta parameter for the resolvent Q iteratively
       :param gamma: ridge penalty value
125
       :param X_train: training set
126
       :param \mathbb{N}: number of neurons in the hidden layer
127
       :param accuracy: accuracy for numerical delta finding
128
129
       :return: optimal value for delta parameter for the resolvent {\tt Q}
130
       n_train = X_train.shape[1]
131
       delta_prev = 1
132
133
       delta_next = 0
       while abs(delta_prev-delta_next) > accuracy:
134
          delta_prev = delta_next
```

```
Q_ = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta_next)) + gamma*np.eye(
136
               n_train))
           delta_next = (1/n_train)*(np.trace(Q_@K(X_train, X_train)))
137
       return delta_next
138
139
140
141
142 def Etrain (gamma, data, npN):
143
      Estimated train error compute (expectation based, deterministic)
144
145
      :param gamma: ridge penalty value
146
       :param data: train/test data
       :param npN: dimensionality of the data and number of Neurons
147
148
      :return: estimated training error
149
      # Unpacking variables:
150
      X_{train}, X_{test}, Y_{train}, Y_{test} = data
151
      n_{train}, n_{test}, p, N = npN
152
153
      delta = find_delta(gamma, X_train, N, delta_accuracy)
      Q_{-} = \text{np.linalg.inv}((\text{N/n\_train})*(\text{K(X\_train, X\_train})/(1+\text{delta})) + \text{gamma*np.eye}(\text{n\_train}))
155
      K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
156
157
      \#E\_train\_ = ((gamma**2)/n\_train)*(Y\_trainQQ\_Q(((((1/N)*np.matrix.trace(Q_QK_QQ_)))/((1 - 1/N))))
158
           N)*np.matrix.trace(K_@Q_@K_@Q_)))*K_ + np.eye(n_train))@Q_@Y_train.T)
      159
           \label{eq:trace} trace(\texttt{K}_@Q_@K_@Q_)))*\texttt{K}_+ np.eye(\texttt{n}_train))@Q_@Y_train.T)
      return E train
161
162
163
164 def Etest_(gamma, data, npN):
165
      Estimated test error compute (expectation based, deterministic)
166
167
      :param gamma: ridge penalty value
       :param data: train/test
      :param npN: dimensionality of the data and number of Neurons
169
      :return: estimated test error
170
171
      # Unpacking variables:
172
      X_train, X_test, Y_train, Y_test = data
173
      n_{train}, n_{test}, p, N = npN
174
175
176
      delta = 0.0
      delta = find_delta(gamma, X_train, N, delta_accuracy)
177
      \label{eq:Q_section} Q_{-} = \text{np.linalg.inv}((\text{N/n\_train})*(\text{K(X\_train, X\_train})/(1+\text{delta})) + \text{gamma*np.eye}(\text{n\_train}))
178
      K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
179
      K_xX = (N/n_{train})*(K(X_{train}, X_{test})/(1+delta))
180
      K_XX = (N/n_{train})*(K(X_{test}, X_{test})/(1+delta))
181
182
      \#E\_test\_ = (1/n\_test)*np.sum((Y\_test.T - K\_xX.T@Q\_@Y\_train.T)**2) + (((1/N)*(
183
           matrix.trace(K_XX) - (1/n_test)*np.matrix.trace((np.eye(n_train) + gamma*Q_)@(
           K xX@K xX.T@Q )))
       (1/n\_test)*np.trace((np.eye(n\_train) + gamma*Q_)@(K_xX@K_xX.T@Q_))) 
185
186
      return E test
187
188
def find_optimal_gamma(gammas, E_func, data, npN):
       errors = [E_func(g, data, npN) for g in gammas]
      return gammas[np.argmin(errors)], errors[np.argmin(errors)]
191
192
193
194 def generate_synthetic_data_regression(npN, density, d=1, noise_level=0.1):
195
196
      Generating X randomly; Y according to the law Y = \Sigma^T * b + eps
       :param \operatorname{np}{\tt N}\colon\operatorname{dimensionality} of the data and number of neurons
197
      :param d: dimensionality of the output vector
```

```
:param noise_level: level of noisiness in Y
199
               :return: training and test data for both X and Y
200
201
               # Unpacking values:
203
               n_train, n_test, p, N = npN
204
205
               X_train = np.random.randn(p, n_train)
206
207
               X_test = np.random.randn(p, n_test)
208
               # Generating Y according to the law: Y = \Sigma^T * b + eps
209
210
               dens = density
               b = np.array([np.random.choice([0, 1], size=d, p=[1-dens/100, dens/100]) for _ in range(N
211
                      )])
212
               W = np.random.randn(N, p)
               Y_generator = lambda X: (sigma(W@X).T@b + np.random.normal(0, noise_level, (X.shape[1], d
213
                       ))).flatten()
214
               Y_train = Y_generator(X_train)
215
               Y_test = Y_generator(X_test)
217
               # Shuffling the dataset
218
               shuffle_train = np.random.permutation(n_train)
219
               shuffle_test = np.random.permutation(n_test)
220
221
              X_train = X_train[:, shuffle_train]
222
223
              Y_train = Y_train[shuffle_train]
               X_test = X_test[:, shuffle_test]
224
               Y_test = Y_test[shuffle_test]
225
226
227
               return X_train, X_test, Y_train, Y_test, W, b
228
229
 \texttt{230} \texttt{ \# def log\_results(log\_file, n, p, N, optimal\_gamma\_Etest, optimal\_gamma\_Etest\_bar, E\_test, optimal\_gamma\_Etest\_bar, E\_test, optimal\_gamma\_Etest\_bar, E\_test, optimal\_gamma\_Etest\_bar, E\_test, optimal\_gamma\_Etest\_bar, E\_test, optimal\_gamma\_Etest\_bar, optimal\_ga
               E_test_):
231 #
                   with open(log_file, 'a') as f:
232 #
                          f.write(f"n = {n}, p = {p}, N = {N}\n")
                           f.write(f"Optimal \ gamma \ for \ Etest: \ \{optimal\_gamma\_Etest\} \ ")
233 #
                          f.write(f"Optimal gamma for Etest_bar: {optimal_gamma_Etest_bar}\n")
                           f.write(f"Etest: {E_test}\n")
235 #
236 #
                           f.write(f"Etest\_bar: \{E\_test\_\}\n\n")
238
239 comm = MPI.COMM_WORLD
240 size = comm.Get_size()
241 rank = comm.Get_rank()
243 delta_accuracy = 1e-3
244 b_vector_density = 50 #%
c1 = 1.0 \# p/n
c2 = 0.8 \# N/n
249 gammas = [10**y for y in np.arange(-7, 1, 0.1)]
250 n_values = np.arange(100, 3000, 100)
251
252 # Split the n_values among available ranks
253 n_split = np.array_split(n_values, size)
254 local_n_values = n_split[rank]
256 optimal_gammas_Etest = []
257 optimal_gammas_Etest_bar = []
259 Etest_values = []
260 Etest_bar_values = []
262
263 for i, n in enumerate(local_n_values):
264
              p = int(c1 * n) # Calculate p based on the constant ratio c1
               N = int(c2 * n) # Calculate N based on the constant ratio c2
265
npN = (int(n * 0.8), int(n * 0.2), p, N) # Recalculate n_train and n_test
```

```
267
        print(f'Rank_\(\text{rank}\)_\-\(\text{Iteration}\(\text{i+1}\)/\(\text{len(local_n_values)}\)')
268
         \begin{array}{l} \textbf{print} (\texttt{f'Rank}_{\sqcup} \{\texttt{rank}\}_{\sqcup} - {}_{\sqcup} \texttt{Trying}_{\sqcup} \texttt{for}_{\sqcup} n_{\sqcup} = {}_{\sqcup} \{\texttt{n}\}, {}_{\sqcup} p_{\sqcup} = {}_{\sqcup} \{\texttt{p}\}, {}_{\sqcup} N_{\sqcup} = {}_{\sqcup} \{\texttt{N}\}_{\sqcup} \dots \text{'}) \end{array} 
269
        # Generate synthetic data
271
        X_train, X_test, Y_train, Y_test, W, b = generate_synthetic_data_regression(npN,
272
             b_vector_density)
        data = X_train, X_test, Y_train, Y_test
273
274
275
        # Find the optimal gammas that minimize Etest and Etest_bar
        optimal_gamma_Etest, E_test = find_optimal_gamma(gammas, Etest, data, npN)
276
277
        optimal_gamma_Etest_bar, E_test_ = find_optimal_gamma(gammas, Etest_, data, npN)
278
        279
        print(f'Ranku{rank}u-uGammauthatuminimizesuEtest_bar:u{optimal_gamma_Etest_bar}')
280
281
282
        print(f'Ranku{rank}u-uEtest:u{E_test}')
        print(f'Ranku{rank}u-uEtest_bar:u{E_test_}')
283
284
        optimal_gammas_Etest.append(optimal_gamma_Etest)
        optimal_gammas_Etest_bar.append(optimal_gamma_Etest_bar)
286
287
        Etest_values.append(E_test)
288
        Etest_bar_values.append(E_test_)
289
290
        # log_results(log_file, n, p, N, optimal_gamma_Etest, optimal_gamma_Etest_bar, E_test,
291
             E_test_)
293 # Gather results from all ranks
optimal_gammas_Etest = comm.gather(optimal_gammas_Etest, root=0)
295 optimal_gammas_Etest_bar = comm.gather(optimal_gammas_Etest_bar, root=0)
296 Etest_values = comm.gather(Etest_values, root=0)
297 Etest_bar_values = comm.gather(Etest_bar_values, root=0)
299 if rank == 0:
        optimal_gammas_Etest = np.concatenate(optimal_gammas_Etest)
300
        optimal_gammas_Etest_bar = np.concatenate(optimal_gammas_Etest_bar)
301
        Etest_values = np.concatenate(Etest_values)
302
        Etest_bar_values = np.concatenate(Etest_bar_values)
303
304
        np.savez('results_parallel.npz',
305
                   n_values=n_values,
306
                   {\tt optimal\_gammas\_Etest=optimal\_gammas\_Etest}\,,
307
308
                   {\tt optimal\_gammas\_Etest\_bar=optimal\_gammas\_Etest\_bar}\,,
                   Etest_values=Etest_values,
309
                  Etest_bar_values=Etest_bar_values,
310
                   c1=c1,
311
                  c2=c2,
312
                  b_vector_density=b_vector_density)
313
```

A.3. Normalized ridge error against regularization parameter γ

```
1 import numpy as np
2 from scipy import linalg
3 import matplotlib.pyplot as plt
6 def sigma(t):
      return t # Linear activation function
def generate_synthetic_data(n, c1, c2, density=0.9, noise_level=0.1):
      p = int(c1 * n)
11
      N = int(c2 * n)
     n_{train} = int(n * 0.8)
13
      n_{test} = n - n_{train}
14
      # Generate random weights and sparse beta
16
      W = np.random.randn(N, p)
      b = np.random.choice([0, 1], size=(N, 1), p=[1 - density, density])
18
19
20
      # Generate data
      X_train = np.random.randn(p, n_train)
21
      Y_train = (sigma(W @ X_train).T @ b).flatten() + np.random.normal(0, noise_level, n_train
      X_test = np.random.randn(p, n_test)
23
      Y_test = (sigma(W @ X_test).T @ b).flatten() + np.random.normal(0, noise_level, n_test)
25
      return X_train, X_test, Y_train, Y_test, W, b
26
29 def compute_ridge_errors(gammas, X_train, Y_train, W, b):
      n_train = X_train.shape[1]
30
      Sigm = sigma(W @ X_train)
31
      SigmT_Sigm = Sigm.T @ Sigm
32
      ridge_errors = []
33
34
      for gamma in gammas:
36
           try:
37
               # Compute ridge solution
               beta_ridge = (Sigm / n_train) @ linalg.solve(
38
                   SigmT_Sigm / n_train + gamma * np.eye(n_train),
39
                   Y_{train}
41
42
               # Compute normalized error
               error = np.linalg.norm(beta_ridge - b, 'fro') ** 2 / np.linalg.norm(b, 'fro') **
44
               ridge_errors.append(error)
45
46
47
          except np.linalg.LinAlgError:
              ridge_errors.append(np.nan)
48
49
      return np.array(ridge_errors)
53 # Parameters (adjust these for faster testing)
n = 512 # Reduced for local testing
55 c1 = 1.0 \# p/n
56 c2 = 0.7 \# N/n
57 b_density = 0.95
58 gammas = [10**y \text{ for } y \text{ in } np.arange(-7, 7, 0.1)]
60 # Generate synthetic data
61 X_train, X_test, Y_train, Y_test, W, b = generate_synthetic_data(
      n=n, c1=c1, c2=c2, density=b_density
62
63 )
65 # Compute ridge errors
66 ridge_errors = compute_ridge_errors(gammas, X_train, Y_train, W, b)
```

A.4. Error Metrics and γ value over number of data points

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy import linalg
6 ################################
7 # 1) Basic definitions
8 ##############################
10 def sigma(t):
      # For this example, just use linear activation
11
      return t
13
14
15 def K(x, y):
      # Kernel function for the linear activation
16
17
      return x.T @ y
18
19
20 def find_delta(gamma, X_train, N, accuracy=1e-3):
21
22
      Numerically solves for the delta parameter used in the deterministic equivalents.
23
     n_train = X_train.shape[1]
24
     delta_prev = 1.0
      delta next = 0.0
26
      while abs(delta_prev - delta_next) > accuracy:
27
          delta_prev = delta_next
          Q_{-} = np.linalg.inv((N / n_train) * K(X_train, X_train) / (1 + delta_next) + gamma *
29
              np.eye(n_train))
          delta_next = (1 / n_train) * np.trace(Q_ @ K(X_train, X_train))
30
      return delta_next
31
34 ##############################
35 # 2) Error metrics
36 ###############################
37
38 def Etest(gamma, X_train, Y_train, X_test, Y_test, N_loops=10):
39
40
      Monte Carlo (random) test MSE:
        Randomly draw W, compute Sigm = sigma(W X_train)Solve ridge, apply to test set, average over N_loops
41
42
      n_train = X_train.shape[1]
44
      n_test = X_test.shape[1]
45
      p = X_train.shape[0]
46
47
48
      E_{test_vals} = []
      for _ in range(N_loops):
49
           W = np.random.randn(N, p)
50
           Sigm = sigma(W @ X_train) # shape (N, n_train)
          Sigm_test = sigma(W @ X_test) # shape (N, n_test)
52
53
           # Solve for ridge coefficients: (1/n_train)*Sigm * inv( Sigm^T Sigm/n_train + gamma I
54
               ) * Y_{train}
55
          A = Sigm.T @ Sigm / n_train + gamma * np.eye(n_train)
           invA_Y = linalg.solve(A, Y_train, assume_a='pos') # shape (n_train, )
56
          beta = (Sigm / n_train) @ invA_Y
57
59
          # Evaluate test error
           pred_test = Sigm_test.T @ beta # shape (n_test,)
60
           err = np.mean((Y_test - pred_test) ** 2)
          E_test_vals.append(err)
62
63
      return np.mean(E_test_vals)
64
65
67 def Etest_det(gamma, X_train, Y_train, X_test, Y_test):
```

```
0.00
68
       Deterministic equivalent test MSE (approximation).
69
70
       n_train = X_train.shape[1]
       n_test = X_test.shape[1]
72
       N = int(c2 * (n_train + n_test)) # or pass as a parameter
73
       delta = find_delta(gamma, X_train, N)
75
       Q_ = np.linalg.inv((N / n_train) * K(X_train, X_train) / (1 + delta) + gamma * np.eye(
76
           n_train))
        \begin{array}{l} \text{K}\_ = \text{(N / n\_train)} * \text{K(X\_train, X\_train)} / \text{(1 + delta)} \\ \text{K}\_\text{xX} = \text{(N / n\_train)} * \text{K(X\_train, X\_test)} / \text{(1 + delta)} \\ \end{array} 
77
78
       K_XX = (N / n_{train}) * K(X_{test}, X_{test}) / (1 + delta)
79
80
       # first "direct" MSE part: (1/n_test)* ||Y_test - K_xX^T Q_ Y_train||^2
81
       residual = Y_test - (K_xX.T @ (Q_ @ Y_train))
82
83
       partA = np.mean(residual ** 2)
84
       # second "correction" part from random matrix theory
85
       num = (1 / N) * (Y_train.T @ Q_ @ K_ @ Q_ @ Y_train)
       den = (1 - (1 / N) * np.trace(K_ @ Q_ @ K_ @ Q_))
87
88
       partB = num / den * ((1 / n_test) * np.trace(K_XX)
                               - (1 / n_test) * np.trace((np.eye(n_train) + gamma * Q_) @ (K_xX @
90
                                   K_xX.T @ Q_)))
92
       return partA + partB
95 def Lridge(gamma, X_train, Y_train, b):
96
       Ridge 'coefficient error': || beta_hat - b ||^2 / ||b||^2
97
98
       with beta_hat found by ridge. We *don't* average over W here,
       since 'b' depends on a *specific* W used to generate data.
99
100
       n_train = X_train.shape[1]
       Sigm = sigma(W_gen @ X_train) # use the same W that generated 'b'
102
103
       A = Sigm.T @ Sigm / n_train + gamma * np.eye(n_train)
       invA_Y = linalg.solve(A, Y_train, assume_a='pos')
105
106
       beta = (Sigm / n_train) @ invA_Y
       return np.linalg.norm(beta - b, 'fro') ** 2 / np.linalg.norm(b, 'fro') ** 2
108
109
110
112 # 3) Local runner
113 ##################################
114
115 # Smaller test parameters
c1 = 1.0 \# p/n
117 c2 = 0.7 \# N/n
118 gammas = [10 ** x for x in np.arange(-7, 5, 0.5)] # e.g. gamma from 1e-3 to 1e3
120 n_values = [500, 600, 700] # smaller range
121 results_gamma_etest = []
122 results_gamma_etest_ = []
results_gamma_lridge = []
125 results_etest_vals = []
126 results_etest_vals_ = []
127 results_lridge_vals = []
129 for n in n values:
130
       # Dimensions
       n_{train} = int(n * 0.8)
       n_test = n - n_train
132
       p = int(c1 * n)
133
       N = int(c2 * n)
134
135
# Generate synthetic data for this n
```

```
# NOTE: store W_gen, b so we can use the same underlying W in Lridge
137
       density = 50 # e.g. 50%
138
       noise\_level = 0.1
139
       \# Create a random W_{\underline{}}gen for generating the data
141
142
       W_gen = np.random.randn(N, p)
143
144
       def Y_generator(X):
145
           # We create a random 'b' with certain density:
146
           b\_local = np.random.choice([0, 1], size=(N, 1),
147
148
                                        p=[1 - density / 100, density / 100])
           # We'll keep track of b_local so we can measure Lridge
149
           # but for clarity, let's store it outside in a closure
150
           return sigma(W_gen @ X).T @ b_local, b_local
151
152
153
       # Make X and Y
154
       X_tr = np.random.randn(p, n_train)
155
       Y_tr_gen, b_local = Y_generator(X_tr)
       # Add noise
157
       Y_tr = Y_tr_gen.ravel() + noise_level * np.random.randn(n_train)
158
159
       X_te = np.random.randn(p, n_test)
160
161
       Y_te_gen, _ = Y_generator(X_te)
       Y_te = Y_te_gen.ravel() + noise_level * np.random.randn(n_test)
162
163
       # Shuffle
164
       perm_tr = np.random.permutation(n_train)
165
       X_tr = X_tr[:, perm_tr]
166
167
       Y_tr = Y_tr[perm_tr]
168
169
       perm_te = np.random.permutation(n_test)
170
       X_te = X_te[:, perm_te]
       Y_te = Y_te[perm_te]
171
173
       # We define a small function that can compute each metric for a given gamma
174
       def measure_all_metrics(g):
           # test error (random-based)
176
           e_test = Etest(g, X_tr, Y_tr, X_te, Y_te, N_loops=5)
177
           # test error (det-based)
178
           e_test_approx = Etest_det(g, X_tr, Y_tr, X_te, Y_te)
179
180
           # ridge error in param space
           # but we need W_gen and b_local from above
181
182
           # put them in global scope or closure
           return e_test, e_test_approx
183
184
185
186
       # Now search for optimal gamma for each metric
       best_e_test = np.inf
187
       best_e_test_g = None
188
189
       best_e_test_approx = np.inf
190
       best_e_test_approx_g = None
192
193
       best_lridge = np.inf
       best_lridge_g = None
194
195
       for g in gammas:
196
            # Etest
197
           val_test = Etest(g, X_tr, Y_tr, X_te, Y_te, N_loops=5)
198
           if val_test < best_e_test:</pre>
199
               best_e_test = val_test
200
201
                best_e_test_g = g
202
           # Etest det
203
           val_test_det = Etest_det(g, X_tr, Y_tr, X_te, Y_te)
204
205
           if val_test_det < best_e_test_approx:</pre>
                best_e_test_approx = val_test_det
206
                best_e_test_approx_g = g
```

```
208
            # Lridge
209
            # We need references to W_gen, b_local in scope
210
            # so let's define them as global, or do it inline:
            Sigm_gen = sigma(W_gen @ X_tr)
212
213
            A = Sigm_gen.T @ Sigm_gen / n_train + g * np.eye(n_train)
            invA_Y = linalg.solve(A, Y_tr, assume_a='pos')
214
            beta_est = (Sigm_gen / n_train) @ invA_Y
215
            this_lridge = (np.linalg.norm(beta_est - b_local, 'fro') ** 2
216
                            / np.linalg.norm(b_local, 'fro') ** 2)
217
            if this_lridge < best_lridge:</pre>
218
219
                best_lridge = this_lridge
                best_lridge_g = g
220
221
       # store results
       results_gamma_etest.append(best_e_test_g)
223
224
       results_etest_vals.append(best_e_test)
225
       results_gamma_etest_.append(best_e_test_approx_g)
226
       results_etest_vals_.append(best_e_test_approx)
228
       results_gamma_lridge.append(best_lridge_g)
229
       results_lridge_vals.append(best_lridge)
231
232 ###############################
233 # 4) Plot results
234 ###############################
235
236 fig, axs = plt.subplots(1, 2, figsize=(12, 5))
237
238 # (a) Plot the optimal gamma for each metric vs n
axs[0].plot(n_values, results_gamma_etest, 'bo-', label='Gamma_{\sqcup}for_{\sqcup}Etest')
240 axs[0].plot(n_values, results_gamma_etest_, 'mo-', label='Gamma_for_Etest__approx')
241 axs[0].plot(n_values, results_gamma_lridge, 'ro-', label='Gamma_for_Lridge')
242 axs[0].set_xscale('log')
243 axs[0].set_yscale('log')
244 axs[0].set_xlabel('n')
245 axs[0].set_ylabel('Optimal_gamma_l(log_scale)')
246 axs[0].legend()
247 axs[0].set_title('Optimal_gamma_vs.un')
248
249 # (b) Plot the *minimum* error achieved vs n (for each metric)
axs[1].plot(n_values, results_etest_vals, 'b*-', label='Min_Etest')
axs[1].plot(n_values, results_etest_vals_, 'm*-', label='Min_Etest_\sqcupapprox')
252 axs[1].plot(n_values, results_lridge_vals, 'r*-', label='Min_Lridge')
253 axs[1].set_xscale('log')
254 axs[1].set_yscale('log')
255 axs[1].set_xlabel('n')
axs[1].set_ylabel('Error_{\sqcup}(log_{\sqcup}scale)')
257 axs[1].legend()
258 axs[1].set\_title('Minimum_achieved_error_vs._n')
260 plt.tight_layout()
261 plt.show()
```

A.5. Derivative-based loss functions against γ parameter (different noise levels)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
4 def sigma(t):
      Small sigma function of choice
6
      :param t: input
      :return: output
8
9
      return t
11
12
13 def K(x, y):
14
15
      Kernel function. Depends on the choice of a small sigma function
      :param x: first 'point' set
16
      :param y: second 'point' set
17
      :return: matrix of measures of 'distances' between points
18
19
20
     return x.T@y
21
22
23 def K_(x, y, delta, N, n_train):
24
25
      Function to compute kernel approximation
      :param x: first 'point' set
      :param y: second 'point' set
27
28
      :param delta: delta parameter
      :return: approximation for the kernel matrix
29
30
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
      return k
32
33
35 def find_delta(gamma, X_train, N, accuracy) -> float:
36
37
      Helper-function that finds delta parameter for the resolvent {\tt Q} iteratively
      :param gamma: ridge penalty value
38
      :param X_train: training set
      :param N: number of neurons in the hidden layer
40
41
      :param accuracy: accuracy for numerical delta finding
      :return: optimal value for delta parameter for the resolvent {\tt Q}
43
44
      n_train = X_train.shape[1]
      delta_prev = 1
45
      delta_next = 0
46
47
      while abs(delta_prev-delta_next) > accuracy:
          delta_prev = delta_next
48
           Q_{-} = \text{np.linalg.inv}((\text{N/n\_train})*(\text{K(X\_train, X\_train})/(1+\text{delta\_next})) + \text{gamma*np.eye}(
49
               n train))
          delta_next = (1/n_train)*(np.trace(Q_@K(X_train, X_train)))
50
51
      return delta_next
52
53
54 def Etrain_(gamma, data, npN):
55
      Estimated train error compute (expectation based, deterministic)
56
      :param gamma: ridge penalty value
58
      :param data: train/test data
      :param \operatorname{np} \mathbb{N}\colon dimensionality of the data and number of Neurons
59
      :return: estimated training error
61
      # Unpacking variables:
62
      X_train, X_test, Y_train, Y_test = data
63
      n_train, n_test, p, N = npN
64
delta = find_delta(gamma, X_train, N, delta_accuracy)
```

```
Q_ = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta)) + gamma*np.eye(n_train))
  67
                         K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
  68
  69
                         \label{eq:proposed_property} \texttt{\#E\_train} = ((\texttt{gamma**2})/\texttt{n\_train}) * (\texttt{Y\_trainQQ\_Q}((((1/N)*\texttt{np.matrix.trace}(\texttt{Q\_QK\_QQ_})))/((1 - 1/N))) + ((\texttt{pamma**2})/\texttt{n\_train}) * (\texttt{pamma**2})/\texttt{n\_train}) * (\texttt{pamma**2})/\texttt{n\_train} * (\texttt{pamma**2})/\texttt{n\_train} * (\texttt{pamma**2})/\texttt{n\_train}) * (\texttt{pamma**2})/\texttt{n\_train} * (\texttt{pamma**2})/\texttt{n\_
                                        N)*np.matrix.trace(K_@Q_@K_@Q_)))*K_ + np.eye(n_train))@Q_@Y_train.T)
                         71
                                         trace(K_0Q_0K_0Q_)))*K_+ np.eye(n_train))@Q_0Y_train.T)
  72
                         return E_train_
  73
  75
  76 def Etest_(gamma, data, npN):
  77
                         Estimated test error compute (expectation based, deterministic)
  78
  79
                         :param gamma: ridge penalty value
                         :param data: train/test
  80
                         :param npN: dimensionality of the data and number of Neurons
  81
  82
                         :return: estimated test error
  83
                         # Unpacking variables:
                         X_train, X_test, Y_train, Y_test = data
  85
                         n_{train}, n_{test}, p, N = npN
  86
                         delta = 0.0
  88
                         delta = find_delta(gamma, X_train, N, delta_accuracy)
  89
                         Q_{-} = \text{np.linalg.inv}((N/n_{-}\text{train})*(K(X_{-}\text{train}, X_{-}\text{train})/(1+\text{delta})) + \text{gamma*np.eye}(n_{-}\text{train}))
  90
  91
                         K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
                         K_xX = (N/n_{train})*(K(X_{train}, X_{test})/(1+delta))
  92
                         K_XX = (N/n_{train})*(K(X_{test}, X_{test})/(1+delta))
  93
  94
                         #E_test_ = (1/n_test)*np.sum((Y_test.T - K_xX.T@Q_@Y_train.T)**2) + (((1/N)*(
                                         Y_{\text{trainQQ_QK_QQ_QY_train.T})} / ((1-1/N)*np.matrix.trace(K_QQ_QK_QQ_)))*((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_))) * ((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_))) * ((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQ
                                         \verb|matrix.trace(K_XX)| - (1/n_test)*np.matrix.trace((np.eye(n_train) + gamma*Q_)@(
                                         K_xX@K_xX.T@Q_)))
                         96
                                          Y_{trainQQ_0K_0Q_0Y_{train.T}))/((1-1/N)*np.trace(K_0Q_0K_0Q_)))*((1/n_{test})*np.trace(K_0Q_0K_0Q_)))*((1/n_{test})*np.trace(K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0
                                         K_XX) - (1/n_test)*np.trace((np.eye(n_train) + gamma*Q_)@(K_xX@K_xX.T@Q_)))
  97
                         return E_test_
 99
100
def estimate_variance_RSS(Y, Sigm, n, gamma):
                         QyT = np.linalg.solve(Sigm.T @ Sigm / n + gamma * np.eye(n), Y.T)
102
103
                         beta_hat = Sigm / n @ QyT
104
105
                         v est = beta hat.T @ Sigm
                         # variance already squared
107
                         variance = ((Y - y_est) @ (Y - y_est).T)/n # sigm^2
108
109
                         return variance
110
def dL_dy_expected(gamma, X, Y, W, variance):
113
                         "Oracle". Conditional expectation w.r.t. epsilon, taken
 114
                         from an unfolded derivative of Loss function ||beta-b||**2
115
116
                         by gamma (set equal to 0). Some 'oracle' prediction for loss.
117
                         :param gamma: gamma parameter (1x1)
                         :param X: input data (pxn)
118
                          :param Y: output data (dxn)
119
                         :param W: weights matrix (Nxp)
120
                         :param eps: noise vector (dxn)
121
                          :return: calculated value of the derivative
122
123
                         Sigm = sigma(W @ X) # Size Nxn
124
125
                         n = X.shape[1]
126
                         Q = (1/n * (Sigm.T @ Sigm) + gamma * np.eye(n))
127
128
                         Q_inv = np.linalg.inv(Q)
                         Q_inv2 = Q_inv @ Q_inv
129
                         Q_{inv3} = Q_{inv2} @ Q_{inv}
```

```
131
       dL_dy_exp = gamma*(Y @ Q_inv3 @ Y.T) - variance * np.trace(Q_inv2)
132
       return dL_dy_exp
133
135
def dL_dy_estimated(gamma, X, Y, W, variance=1):
137
       Calculation of the unfolded derivative of Loss function ||beta-b||**2
138
       under expectation w.r.t. epsilon, with estimated variance instead of real 'oracle' one.
139
       :param gamma: gamma parameter (1x1)
140
       :param X: input data (pxn)
141
142
       :param Y: output data (dxn)
       :param W: weights matrix (Nxp)
143
144
       :param eps: noise vector (dxn)
       :return: calculated value of the derivative
145
146
       Sigm = sigma(W @ X) # Size Nxn
147
148
       n = X.shape[1]
149
       Q = (1 / n * (Sigm.T @ Sigm) + gamma * np.eye(n))
       Q_inv = np.linalg.inv(Q)
Q_inv2 = Q_inv @ Q_inv
151
152
       Q_inv3 = Q_inv2 @ Q_inv
153
154
       estimated_var = estimate_variance_RSS(Y, Sigm, n, gamma)
155
156
       \label{eq:dl_dy_exp} dL_dy_exp = gamma * (Y @ Q_inv3 @ Y.T) - estimated_var * np.trace(Q_inv2)
157
       return dL_dy_exp
158
159
160
161 def generate_synthetic_data_regression(npN, density, d=1, noise_level=0.1):
162
163
       Generates data for the model:
           Y = b^T * Sigma + eps
164
       with Sigma = W @ X and b has 'density'% of 1's.
165
       Shapes:
167
         X_train: (p, n_train)
168
         X_test: (p, n_test)
         W:
                   (N, p)
170
                             -- if d=1 => shape (N,1)
171
         b:
                   (N, d)
172
         Y_train, e_train: (d, n_train)
         Y_test, e_test: (d, n_test)
173
174
       :param npN: (n_train, n_test, p, N)
175
176
       :param density: percentage for b's 1's
       :param d: dimension of output (1 => scalar outputs)
177
       :param noise_level: std-dev of Gaussian noise
178
       :return: (X_train, X_test, Y_train, Y_test, W, b, e_train, e_test)
179
180
       n_train, n_test, p, N = npN
181
182
       # 1) Generate X
183
       X_train = np.random.randn(p, n_train) # (p, n_train)
184
       X_test = np.random.randn(p, n_test) # (p, n_test)
186
187
       # 2) Generate b (N x d), if d=1 \Rightarrow (N,1)
       b = np.random.choice([0, 1],
188
                              size=(N, d),
189
                              p=[1 - density / 100, density / 100])
190
191
       # 3) Generate W (N x p)
192
       W = np.random.randn(N, p)
193
194
195
       def make_data(X):
            Given X with shape (p, n), returns:
197
             Y: (d, n)
e: (d, n)
198
199
              Sigm: (N, n) [might be useful if needed]
200
```

```
# Sigma = W @ X \Rightarrow (N, n)
202
            Sigm = sigma(W @ X)
203
204
            # Noise eps \Rightarrow (d, n)
            e = np.random.normal(loc=0, scale=noise_level, size=(d, X.shape[1]))
206
207
            # b^T shape => (d, N), so b^T @ Sigm => (d, n)
208
           Y = (b.T @ Sigm) + e
209
210
211
           return Y, e, Sigm
212
213
       # 4) Make training data (Y_train_full, e_train_full)
       Y_train_full, e_train_full, Sigm_train = make_data(X_train)
214
215
       # 5) Make test data (Y_test_full, e_test_full)
216
       Y_test_full, e_test_full, Sigm_test = make_data(X_test)
217
218
       # 6) Shuffle columns (the "sample" axis) in X, Y, e
219
            X is (p, n), Y,e are (d, n). We shuffle axis=1 for each.
220
       idx_train = np.random.permutation(n_train)
       idx_test = np.random.permutation(n_test)
222
223
       X_train = X_train[:, idx_train]
       Y_train_full = Y_train_full[:, idx_train]
225
       e_train_full = e_train_full[:, idx_train]
226
227
228
       X_test = X_test[:, idx_test]
       Y_test_full = Y_test_full[:, idx_test]
229
       e_test_full = e_test_full[:, idx_test]
230
231
       # 7) Return them all in consistent shapes
       \# No flattening is required; Y and e remain (d,n).
233
234
       return X_train, X_test, Y_train_full, Y_test_full, W, b, e_train_full, e_test_full
235
236
237 def L_ridge(gamma, X, Y, W, b):
       n = X.shape[1]
238
       Sigm = sigma(W @ X)
239
       QyT = np.linalg.solve(Sigm.T@Sigm/n + gamma*np.eye(n), Y.T)
241
       beta = Sigm/n @ QyT
242
243
       numerator = np.linalg.norm(beta - b,'fro')**2
denominator = np.linalg.norm(b, 'fro')**2
244
245
246
247
       return numerator/denominator
249
250 def plot_estimated_vs_oracle_different_xranges():
251
252
       Left subplot:
         dL/dgamma (estimated RSS) over gamma in [1e-5, 1e3] (step 0.005 in log10).
253
       Right subplot:
254
         dL/dgamma (oracle) over gamma in [1e-3, 1e2] (step 0.005 in log10).
255
       We loop over noise_levels=[1,2,3,4,5].
257
258
       # Basic scenario
259
260
       noise\_levels = [1,2,3,4,5]
261
       \# c1, c2 => p/n=1.0, N/n=0.8
262
       n = 200
263
       c1 = 1.0
       c2 = 0.8
265
266
       p = int(c1*n)
267
       N = int(c2*n)
268
       n_{train} = int(0.8*n)
269
270
       n_{test} = n - n_{train}
       npN = (n_train, n_test, p, N)
271
```

```
273
        # Different gamma ranges
        gammas_est = [10**y \text{ for } y \text{ in } np.arange(-4, 3, 0.005)]
                                                                       # ~1e-5 to ~1e3
274
        gammas_oracle = [10**y for y in np.arange(-4, 2, 0.005)]# ~1e-3 to ~1e2
275
276
        fig, (ax_est, ax_oracle) = plt.subplots(1, 2, figsize=(12,5))
277
278
        # Loop over noise
279
        for nl in noise_levels:
280
281
            # Generate data for each noise
282
            X_train, X_test, Y_train, Y_test, W, b, e_train, e_test = \
                 \tt generate\_synthetic\_data\_regression(npN, density=50, d=1, noise\_level=nl)
283
284
            # Evaluate dL/dgamma (est) over gammas_est
285
286
            dL_est_vals = []
287
            for g in gammas_est:
                 val_est = dL_dy_estimated(g, X_train, Y_train, W, variance=1.0)
288
                 dL_est_vals.append(val_est.item() if hasattr(val_est, 'item') else float(val_est)
289
290
            # Evaluate dL/dgamma (oracle) over gammas_oracle
            dL_oracle_vals = []
292
            # We pass in the real variance = nl^2
293
            for g in gammas_oracle:
                 val_oracle = dL_dy_expected(g, X_train, Y_train, W, variance=n1**2)
dL_oracle_vals.append(val_oracle.item() if hasattr(val_oracle, 'item') else float
295
296
                      (val_oracle))
297
            # Plot them
            ax_est.loglog(gammas_est, dL_est_vals, label=f"noise={nl}")
299
            ax_oracle.loglog(gammas_oracle, dL_oracle_vals, label=f"noise={nl}")
300
        # Label/Legend subplots
302
303
        ax\_est.set\_title("Estimated_{\sqcup}Derivative_{\sqcup}(RSS)_{\sqcup}-_{\sqcup}Gamma_{\sqcup}in_{\sqcup}[1e-5,1e3]")
        ax_est.set_xlabel("Gamma_{\sqcup}(log_{\sqcup}scale)")
304
        ax_est.set_ylabel("dL/dGamma_(EstVar)")
305
        ax_est.grid(True, which='both', ls='--', alpha=0.7)
        ax_est.legend()
307
308
        ax\_oracle.set\_title("Oracle\_Derivative\_-\_Gamma\_in\_[1e-3,1e2]")
        ax\_oracle.set\_xlabel("Gamma_{\sqcup}(log_{\sqcup}scale)")
310
        ax_oracle.set_ylabel("dL/dGamma_(Oracle)")
311
        ax_oracle.grid(True, which='both', ls='--', alpha=0.7)
312
        ax_oracle.legend()
313
314
       fig.tight_layout()
315
316
        plt.show()
317
318
319 if __name__ == "__main__":
plot_estimated_vs_oracle_different_xranges()
```

A.6. Direct loss with γ derived from minimizing different loss functions (including the derivative-based)

```
1 import numpy as np
2 import matplotlib.pyplot as plt
4 def sigma(t):
       Small sigma function of choice
6
      :param t: input
      :return: output
8
9
      return t
11
12
13 def K(x, y):
14
      Kernel function. Depends on the choice of a small sigma function
      :param x: first 'point' set
16
      :param y: second 'point' set
17
      :return: matrix of measures of 'distances' between points
18
19
     return x.T@y
20
21
23 def K_(x, y, delta, N, n_train):
24
25
      Function to compute kernel approximation
      :param x: first 'point' set
      :param y: second 'point' set
27
28
       :param delta: delta parameter
      :return: approximation for the kernel matrix
29
30
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
      return k
32
33
35 def find_delta(gamma, X_train, N, accuracy) -> float:
36
37
      Helper-function that finds delta parameter for the resolvent Q iteratively
      :param gamma: ridge penalty value
38
      :param X_train: training set
      :param N: number of neurons in the hidden layer
40
41
      :param accuracy: accuracy for numerical delta finding
      :return: optimal value for delta parameter for the resolvent {\tt Q}
43
44
      n_train = X_train.shape[1]
      delta_prev = 1
45
      delta_next = 0
46
47
      while abs(delta_prev-delta_next) > accuracy:
          delta_prev = delta_next
48
           \label{eq:Q_section} Q_= = \text{np.linalg.inv}((N/n\_\text{train})*(K(X\_\text{train}, X\_\text{train})/(1+\text{delta\_next})) + \text{gamma*np.eye}(
49
               n train))
          delta_next = (1/n_train)*(np.trace(Q_@K(X_train, X_train)))
50
51
      return delta_next
52
53
54 def Etrain_(gamma, data, npN):
55
       Estimated train error compute (expectation based, deterministic)
56
      :param gamma: ridge penalty value
58
      :param data: train/test data
       :param \operatorname{np} \mathbb{N}\colon dimensionality of the data and number of Neurons
59
      :return: estimated training error
61
      # Unpacking variables:
      X_train, X_test, Y_train, Y_test = data
63
      n_{train}, n_{test}, p, N = npN
64
delta = find_delta(gamma, X_train, N, delta_accuracy)
```

```
Q_ = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta)) + gamma*np.eye(n_train))
 67
                 K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
 68
 69
                 \label{eq:normalized} \mbox{N)*np.matrix.trace($K_QQ_QK_QQ_)))*K_ + \mbox{np.eye(n_train)} \mbox{$QQ_QY_train.T)$}
                 71
                           trace(K_@Q_@K_@Q_)))*K_ + np.eye(n_train))@Q_@Y_train.T)
 72
                 return E_train_
 73
 75
 76 def Etest_(gamma, data, npN):
 77
                 Estimated test error compute (expectation based, deterministic)
 78
 79
                 :param gamma: ridge penalty value
                 :param data: train/test
 80
                 :param npN: dimensionality of the data and number of Neurons
 81
                 :return: estimated test error
 82
 83
                 # Unpacking variables:
                 X_train, X_test, Y_train, Y_test = data
 85
                 n_{train}, n_{test}, p, N = npN
 86
 87
                 delta = 0.0
 88
                 delta = find_delta(gamma, X_train, N, delta_accuracy)
 89
                 Q_{-} = \text{np.linalg.inv}((N/n_{-}\text{train})*(K(X_{-}\text{train}, X_{-}\text{train})/(1+\text{delta})) + \text{gamma*np.eye}(n_{-}\text{train}))
 90
 91
                 K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
                 K_xX = (N/n_{train})*(K(X_{train}, X_{test})/(1+delta))
 92
                 K_XX = (N/n_{train})*(K(X_{test}, X_{test})/(1+delta))
 93
 94
                 #E_test_ = (1/n_test)*np.sum((Y_test.T - K_xX.T@Q_@Y_train.T)**2) + (((1/N)*(
                            Y_{\text{trainQQ_QK_QQ_QY_train.T})} / ((1-1/N)*np.matrix.trace(K_QQ_QK_QQ_)))*((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_))) * ((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_))) * ((1/n_{\text{test}})*np.matrix.trace(K_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQQ_QK_QQ
                           \verb|matrix.trace(K_XX)| - (1/n_test)*np.matrix.trace((np.eye(n_train) + gamma*Q_)@(
                           K_xX@K_xX.T@Q_)))
                 96
                            Y_{trainQQ_0K_0Q_0Y_{train.T}))/((1-1/N)*np.trace(K_0Q_0K_0Q_)))*((1/n_{test})*np.trace(K_0Q_0K_0Q_)))*((1/n_{test})*np.trace(K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0
                           K_XX) - (1/n_{test})*np.trace((np.eye(n_train) + gamma*Q_)@(K_xX@K_xX.T@Q_)))
 97
                 return E_test_
 99
100
def estimate_variance_RSS(Y, Sigm, n, gamma):
                 102
103
                 beta_hat = Sigm / n @ QyT
104
                 y_est = beta_hat.T @ Sigm
105
                 # variance already squared
107
                 variance = ((Y - y_est) @ (Y - y_est).T)/n # sigm^2
108
109
                 return variance
110
def dL_dy_expected(gamma, X, Y, W, variance):
113
                 "Oracle". Conditional expectation w.r.t. epsilon, taken
114
                 from an unfolded derivative of Loss function ||beta-b||**2
115
116
                 by gamma (set equal to 0). Some 'oracle' prediction for loss.
                 :param gamma: gamma parameter (1x1)
117
                 :param X: input data (pxn)
118
                 :param Y: output data (dxn)
119
                 :param W: weights matrix (Nxp)
120
                 :param eps: noise vector (dxn)
121
                 :return: calculated value of the derivative
122
123
124
                 Sigm = sigma(W @ X) # Size Nxn
125
                 n = X.shape[1]
126
                 Q = (1/n * (Sigm.T @ Sigm) + gamma * np.eye(n))
127
128
                 Q_inv = np.linalg.inv(Q)
                 Q_inv2 = Q_inv @ Q_inv
129
                 Q_{inv3} = Q_{inv2} @ Q_{inv}
```

```
131
       dL_dy_exp = gamma*(Y @ Q_inv3 @ Y.T) - variance * np.trace(Q_inv2)
132
       return dL_dy_exp
133
135
def dL_dy_estimated(gamma, X, Y, W, variance=1):
137
       Calculation of the unfolded derivative of Loss function ||beta-b||**2
138
139
       under expectation w.r.t. epsilon, with estimated variance instead of real 'oracle' one.
       :param gamma: gamma parameter (1x1)
140
       :param X: input data (pxn)
141
142
       :param Y: output data (dxn)
       :param W: weights matrix (Nxp)
143
144
       :param eps: noise vector (dxn)
       :return: calculated value of the derivative
145
146
       Sigm = sigma(W @ X) # Size Nxn
147
       n = X.shape[1]
148
149
       Q = (1 / n * (Sigm.T @ Sigm) + gamma * np.eye(n))
       Q_inv = np.linalg.inv(Q)
Q_inv2 = Q_inv @ Q_inv
151
152
       Q_inv3 = Q_inv2 @ Q_inv
153
154
       estimated_var = estimate_variance_RSS(Y, Sigm, n, gamma)
155
156
       \label{eq:dl_dy_exp} dL_dy_exp = gamma * (Y @ Q_inv3 @ Y.T) - estimated_var * np.trace(Q_inv2)
157
       return dL_dy_exp
158
159
160
161 def generate_synthetic_data_regression(npN, density, d=1, noise_level=0.1):
162
163
       Generates data for the model:
           Y = b^T * Sigma + eps
164
       with Sigma = W @ X and b has 'density'% of 1's.
165
       Shapes:
167
         X_train: (p, n_train)
168
         X_test: (p, n_test)
         W:
                   (N, p)
170
                             -- if d=1 => shape (N,1)
171
         b:
                   (N, d)
         Y_train, e_train: (d, n_train)
172
         Y_test, e_test: (d, n_test)
173
174
       :param npN: (n_train, n_test, p, N)
175
176
       :param density: percentage for b's 1's
       :param d: dimension of output (1 => scalar outputs)
177
       :param noise_level: std-dev of Gaussian noise
178
       :return: (X_train, X_test, Y_train, Y_test, W, b, e_train, e_test)
179
180
       n_{train}, n_{test}, p, N = npN
181
182
       # 1) Generate X
183
       X_train = np.random.randn(p, n_train) # (p, n_train)
184
       X_test = np.random.randn(p, n_test) # (p, n_test)
186
187
       # 2) Generate b (N x d), if d=1 \Rightarrow (N,1)
       b = np.random.choice([0, 1],
188
                              size=(N, d),
189
                              p=[1 - density / 100, density / 100])
190
191
       # 3) Generate W (N x p)
192
       W = np.random.randn(N, p)
193
194
195
       def make_data(X):
196
            Given X with shape (p, n), returns:
197
             Y: (d, n)
e: (d, n)
198
199
             Sigm: (N, n) [might be useful if needed]
200
```

```
# Sigma = W @ X => (N, n)
202
            Sigm = sigma(W @ X)
203
204
           # Noise eps => (d, n)
           e = np.random.normal(loc=0, scale=noise_level, size=(d, X.shape[1]))
206
207
            # b^T shape => (d, N), so b^T @ Sigm => (d, n)
208
           Y = (b.T @ Sigm) + e
209
210
           return Y, e, Sigm
211
212
213
       # 4) Make training data (Y_train_full, e_train_full)
       Y_train_full, e_train_full, Sigm_train = make_data(X_train)
214
215
       # 5) Make test data (Y_test_full, e_test_full)
216
       Y_test_full, e_test_full, Sigm_test = make_data(X_test)
217
218
       # 6) Shuffle columns (the "sample" axis) in X, Y, e
219
          X is (p, n), Y,e are (d, n). We shuffle axis=1 for each.
220
       idx_train = np.random.permutation(n_train)
       idx_test = np.random.permutation(n_test)
222
223
       X_train = X_train[:, idx_train]
224
       Y_train_full = Y_train_full[:, idx_train]
225
       e_train_full = e_train_full[:, idx_train]
226
227
228
       X_test = X_test[:, idx_test]
       Y_test_full = Y_test_full[:, idx_test]
       e_test_full = e_test_full[:, idx_test]
230
231
       # 7) Return them all in consistent shapes
       \# No flattening is required; Y and e remain (d,n).
233
234
       return X_train, X_test, Y_train_full, Y_test_full, W, b, e_train_full, e_test_full
235
236
237 def L_ridge(gamma, X, Y, W, b):
       n = X.shape[1]
238
       Sigm = sigma(W @ X)
239
240
       QyT = np.linalg.solve(Sigm.T@Sigm/n + gamma*np.eye(n), Y.T)
241
       beta = Sigm/n @ QyT
242
243
       numerator = np.linalg.norm(beta - b,'fro')**2
denominator = np.linalg.norm(b, 'fro')**2
244
245
246
       return numerator/denominator
247
249
250 def plot_estimated_vs_oracle_different_xranges():
251
252
       Left subplot:
         dL/dgamma (estimated RSS) over gamma in [1e-5, 1e3] (step 0.005 in log10).
       Right subplot:
254
        dL/dgamma (oracle) over gamma in [1e-3, 1e2] (step 0.005 in log10).
255
       We loop over noise_levels=[1,2,3,4,5].
257
258
       # Basic scenario
259
260
       noise\_levels = [1,2,3,4,5]
261
       \# c1, c2 => p/n=1.0, N/n=0.8
262
       n = 200
263
       c1 = 1.0
       c2 = 0.8
265
266
267
       p = int(c1*n)
       N = int(c2*n)
268
       n_{train} = int(0.8*n)
269
270
       n_{test} = n - n_{train}
       npN = (n_{train}, n_{test}, p, N)
271
```

```
273
        # Different gamma ranges
        gammas_est = [10**y for y in np.arange(-4, 3, 0.005)]
                                                                      # ~1e-5 to ~1e3
274
        gammas_oracle = [10**y for y in np.arange(-4, 2, 0.005)]# ~1e-3 to ~1e2
275
       fig, (ax_est, ax_oracle) = plt.subplots(1, 2, figsize=(12,5))
277
278
        # Loop over noise
279
        for nl in noise_levels:
280
281
            # Generate data for each noise
282
            X_train, X_test, Y_train, Y_test, W, b, e_train, e_test = \
                 \tt generate\_synthetic\_data\_regression(npN, density=50, d=1, noise\_level=nl)
283
284
            # Evaluate dL/dgamma (est) over gammas_est
285
286
            dL_est_vals = []
287
            for g in gammas_est:
                 val_est = dL_dy_estimated(g, X_train, Y_train, W, variance=1.0)
288
289
                 dL_est_vals.append(val_est.item() if hasattr(val_est, 'item') else float(val_est)
290
            # Evaluate dL/dgamma (oracle) over gammas_oracle
            dL_oracle_vals = []
292
            # We pass in the real variance = nl^2
293
            for g in gammas_oracle:
                 val_oracle = dL_dy_expected(g, X_train, Y_train, W, variance=n1**2)
dL_oracle_vals.append(val_oracle.item() if hasattr(val_oracle, 'item') else float
295
296
                     (val_oracle))
297
            # Plot them
            ax_est.loglog(gammas_est, dL_est_vals, label=f"noise={nl}")
299
            ax_oracle.loglog(gammas_oracle, dL_oracle_vals, label=f"noise={nl}")
300
301
        # Label/Legend subplots
302
303
        ax\_est.set\_title("Estimated_{\sqcup}Derivative_{\sqcup}(RSS)_{\sqcup}-_{\sqcup}Gamma_{\sqcup}in_{\sqcup}[1e-5,1e3]")
        ax_est.set_xlabel("Gamma_{\sqcup}(log_{\sqcup}scale)")
304
        ax_est.set_ylabel("dL/dGamma_(EstVar)")
305
        ax_est.grid(True, which='both', ls='--', alpha=0.7)
        ax_est.legend()
307
308
        ax\_oracle.set\_title("Oracle\_Derivative\_-\_Gamma\_in\_[1e-3,1e2]")
        ax\_oracle.set\_xlabel("Gamma_{\sqcup}(log_{\sqcup}scale)")
310
        ax_oracle.set_ylabel("dL/dGamma_(Oracle)")
311
        ax_oracle.grid(True, which='both', ls='--', alpha=0.7)
312
        ax_oracle.legend()
313
314
       fig.tight_layout()
315
316
        plt.show()
318
319 if __name__ == "__main__":
plot_estimated_vs_oracle_different_xranges()
```

A.7. Convergence Study

```
1 import numpy as np
2 from scipy import linalg
4 def sigma(t):
       Small sigma function of choice
       :param t: input
7
      :return: output
8
      0.00
      return t
10
11
13 def K(x, y):
       Kernel function. Depends on the choice of a small sigma function
15
      :param x: first 'point' set
16
17
       :param y: second 'point' set
      :return: matrix of measures of 'distances' between points
18
19
20
      return x.T@y
21
23 def K_(x, y, delta, N, n_train):
24
       Function to compute kernel approximation
      :param x: first 'point' set
:param y: second 'point' set
26
27
      :param delta: delta parameter
      :return: approximation for the kernel matrix
29
30
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
31
32
      return k_
35 def Etrain(gamma, data, npN, monte_carlo_loops=30):
       Actual training error compute (random-based)
37
38
       :param gamma: ridge penalty value
       :param data: train/test data
39
       :param \ensuremath{\mathtt{npN}}\xspace dimensionality of the data and \ensuremath{\mathtt{num}}\xspace of neurons
40
       :param monte_carlo_loops: number of iterations for E value averaging for different
           generated W
      :return: Training error
42
       # Unpacking variables
44
       X_train, X_test, Y_train, Y_test = data
45
      n_train, n_test, p, N = npN
46
47
48
      E_train_arr = []
      p = X_train.shape[0]
49
50
      for i in range(monte_carlo_loops):
52
           # W = np.random.randn(N, p)
53
           # Sigm = sigma(W@X_train)
54
           # Q_y = np.linalg.inv((1/n_train)*Sigm.T@Sigm + gamma*np.eye(n_train))
55
           # E_train = (gamma*gamma/n_train)*Y_train@np.linalg.matrix_power(Q_y, 2)@Y_train.T
57
           # E_train_arr.append(E_train)
58
60
           W = np.random.randn(N, p)
           Sigm = sigma(W @ X_train)
61
           Sigm_ = sigma(W @ X_test)
63
64
           inv_tQ_r = linalg.solve(Sigm.T @ Sigm / n_train + gamma * np.eye(n_train), Y_train)
           beta = Sigm / n_train @ inv_tQ_r
65
66
           E_train = np.linalg.norm(Y_train-Sigm.T@beta)**2/n_train
           E_train_arr.append(E_train)
```

```
69
70
       return np.mean(np.array(E_train_arr))
71
73 def Etest(gamma, data, npN, monte_carlo_loops=20):
74
       Actual test error compute (random-based)
75
76
       :param gamma: ridge penalty value
77
       :param data: train/test data
       :param npN: dimensionality of the data and number of neurons
78
       :param monte_carlo_loops: number of iterations for E value averaging for different
79
            generated W
       :return: Test error
80
81
       # Unpacking variables
82
       X_{train}, X_{test}, Y_{train}, Y_{test} = data
83
84
       n_train, n_test, p, N = npN
85
       E_test_arr = []
86
       p = X_train.shape[0]
88
       for i in range(monte_carlo_loops):
89
           W = np.random.randn(N, p)
91
92
           Sigm = sigma(W@X_train)
           Sigm_ = sigma(W@X_test)
93
94
           inv_tQ_r = linalg.solve(Sigm.T@Sigm/n_train + gamma * np.eye(n_train), Y_train)
95
           beta = Sigm/n_train @ inv_tQ_r
96
97
           # Q_y = np.linalg.inv((1 / n_train) * Sigm.T @ Sigm + gamma * np.eye(n_train))
99
100
           \# term1 = (1/n_test)*(Y_test@Y_test.T)
           # term2 = (2/(n_train*n_test))*(Y_train@Q_y@Sigm.T@Sigm_@Y_test.T)
101
            \texttt{\# term3} = (1/(n\_train**2*n\_test))*(Y\_train@Q\_y@Sigm.T@Sigm\_@Sigm\_.T@Sigm@Q\_y@Y\_train.
102
                T)
103
104
           E_test = np.linalg.norm(Y_test-Sigm_.T@beta)**2/n_test
           E_test_arr.append(E_test)
106
107
       ans = np.mean(np.array(E_test_arr))
108
109
110
       return ans
111
112
113 def find_delta(gamma, X_train, N, accuracy) -> float:
114
        \label{thm:equation:equation:equation} \mbox{Helper-function that finds delta parameter for the resolvent } \mbox{\em Q} \mbox{ iteratively} 
115
116
       :param gamma: ridge penalty value
       :param X_train: training set
117
       :param N: number of neurons in the hidden layer
118
       :param accuracy: accuracy for numerical delta finding
119
       :return: optimal value for delta parameter for the resolvent Q
120
       n_train = X_train.shape[1]
122
123
       delta_prev = 1
       delta_next = 0
124
125
       while abs(delta_prev-delta_next) > accuracy:
            delta_prev = delta_next
126
           Q_ = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta_next)) + gamma*np.eye(
127
                n_train))
           delta_next = (1/n_train)*(np.trace(Q_@K(X_train, X_train)))
       return delta next
129
130
132
133 def Etrain_(gamma, data, npN):
134
       Estimated train error compute (expectation based, deterministic)
135
:param gamma: ridge penalty value
```

```
:param data: train/test data
137
                        :param npN: dimensionality of the data and number of Neurons
138
139
                        :return: estimated training error
                       # Unpacking variables:
141
                       X_{train}, X_{test}, Y_{train}, Y_{test} = data
142
                       n_train, n_test, p, N = npN
143
144
                       delta = find_delta(gamma, X_train, N, delta_accuracy)
145
                       Q = np.linalg.inv((N/n_train)*(K(X_train, X_train)/(1+delta)) + gamma*np.eye(n_train))
146
                       K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
147
 148
                       \#E_{train} = ((gamma**2)/n_{train})*(Y_{train}@Q_0((((1/N)*np.matrix.trace(Q_0K_0Q_)))/((1 - 1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N))*((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N)*np.matrix.trace(Q_0K_0Q_))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q_0K_0Q_0))/((1/N)*np.matrix.trace(Q
149
                                      \label{eq:normatrix.trace(K_QQ_0K_QQ_))} $$ N)*np.matrix.trace(K_QQ_0K_QQ_))*K_ + np.eye(n_train))@Q_0Y_train.T) $$
                       150
                                     \label{eq:trace} trace(\texttt{K}_@Q_@K_@Q_)))*\texttt{K}_+ np.eye(\texttt{n}_train))@Q_@Y_train.T)
151
                       return E_train_
152
153
155 def Etest_(gamma, data, npN):
156
                       Estimated test error compute (expectation based, deterministic)
157
                       :param gamma: ridge penalty value
158
159
                        :param data: train/test
                       :param npN: dimensionality of the data and number of Neurons
160
161
                       :return: estimated test error
162
                       # Unpacking variables:
163
                       X_{train}, X_{test}, Y_{train}, Y_{test} = data
164
165
                       n_train, n_test, p, N = npN
166
167
                       delta = 0.0
                       delta = find_delta(gamma, X_train, N, delta_accuracy)
168
169
                       Q = \text{np.linalg.inv}((N/n_{\text{train}})*(K(X_{\text{train}}, X_{\text{train}})/(1+\text{delta})) + \text{gamma*np.eye}(n_{\text{train}}))
                       K_{-} = (N/n_{train})*(K(X_{train}, X_{train})/(1+delta))
                       K_xX = (N/n_train)*(K(X_train, X_test)/(1+delta))

K_xX = (N/n_train)*(K(X_test, X_test)/(1+delta))
171
172
                       \#E_{test} = (1/n_{test})*np.sum((Y_{test}.T - K_xX.TQQ_QY_{train}.T)**2) + (((1/N)*(1/N)*(1/N)*1) + ((1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)*(1/N)
174
                                      matrix.trace(K_XX) - (1/n_test)*np.matrix.trace((np.eye(n_train) + gamma*Q_)@(
                                     K \times X \otimes K \times X \cdot T \otimes Q )))
                       Y_{trainQQ_0K_0Q_0Y_train.T))/((1-1/N)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_)))*((1/n_test)*np.trace(K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q_0K_0Q
                                       \texttt{K_XX}) \ - \ (1/n\_test)*np.trace((np.eye(n\_train) \ + \ gamma*Q\_)@(\texttt{K_xX@K\_xX.T@Q\_}))) 
                       return E_test_
177
178
179
def estimate_variance_RSS(Y, Sigm, n, gamma):
                       QyT = np.linalg.solve(Sigm.T @ Sigm / n + gamma * np.eye(n), Y.T)
181
                       beta_hat = Sigm / n @ QyT
182
183
                       y_est = beta_hat.T @ Sigm
185
186
                       # variance already squared
                       variance = ((Y - y_est) @ (Y - y_est).T)/n # sigm^2
187
188
                       return variance
189
190
def dL_dy_real(gamma, X, Y, W, eps):
192
                       Real unfolded derivative of Loss function ||beta-b||**2 by gamma (set equal to 0)
193
194
                       :param gamma: gamma parameter (1x1)
195
                        :param X: input data (pxn)
                       :param Y: output data (dxn)
196
                       :param W: weights matrix (Nxp)
197
198
                       :param eps: noise vector (dxn)
                       :return: calculated value of the derivative
199
```

```
Sigm = sigma(W @ X) # Size Nxn
201
       n = X.shape[1]
202
203
       Q = (1/n) * (Sigm.T @ Sigm) + gamma * np.eye(n)
       Q_inv = np.linalg.inv(Q)
205
       Q_{inv2} = Q_{inv} @ Q_{inv}
206
       Q_inv3 = Q_inv2 @ Q_inv
207
208
       dL_dy = gamma*(Y @ Q_inv3 @ Y.T) - (Y @ Q_inv2 @ eps.T)
209
       return dL_dy
210
211
212
213 def dL_dy_expected(gamma, X, Y, W, variance):
214
215
       "Oracle". Conditional expectation w.r.t. epsilon, taken
       from an unfolded derivative of Loss function ||beta-b||**2
216
217
       by gamma (set equal to 0). Some 'oracle' prediction for loss.
       :param gamma: gamma parameter (1x1)
218
       :param X: input data (pxn)
219
       :param Y: output data (dxn)
       :param W: weights matrix (Nxp)
221
       :param eps: noise vector (dxn)
222
       :return: calculated value of the derivative
224
       Sigm = sigma(W @ X) # Size Nxn
225
       n = X.shape[1]
226
227
       Q = (1/n * (Sigm.T @ Sigm) + gamma * np.eye(n))
228
       Q_inv = np.linalg.inv(Q)
229
       Q_inv2 = Q_inv @ Q_inv
230
       Q_inv3 = Q_inv2 @ Q_inv
232
233
       dL_dy_exp = gamma*(Y @ Q_inv3 @ Y.T) - variance * np.trace(Q_inv2)
234
       return dL_dy_exp
235
237 def dL_dy_estimated(gamma, X, Y, W, variance=1):
238
       Calculation of the unfolded derivative of Loss function ||beta-b||**2
       under expectation w.r.t. epsilon, with estimated variance instead of real 'oracle' one.
240
241
       :param gamma: gamma parameter (1x1)
       :param X: input data (pxn)
242
243
       :param Y: output data (dxn)
244
       :param W: weights matrix (Nxp)
       :param eps: noise vector (dxn)
245
246
       :return: calculated value of the derivative
247
       Sigm = sigma(W @ X) # Size Nxn
248
249
       n = X.shape[1]
250
       Q = (1 / n * (Sigm.T @ Sigm) + gamma * np.eye(n))
251
       Q_inv = np.linalg.inv(Q)
252
       Q_inv2 = Q_inv @ Q_inv
253
       Q inv3 = Q_inv2 @ Q_inv
254
       estimated_var = estimate_variance_RSS(Y, Sigm, n, gamma)
256
257
       dL_dy_exp = gamma * (Y @ Q_inv3 @ Y.T) - estimated_var * np.trace(Q_inv2)
258
259
       return dL_dy_exp
260
261
def dL_dy_estimated_fixed_rss(gamma, X, Y, W, variance, gamma_rss):
263
       "Estimated" derivative: dL/dgamma, but the RSS-based variance is computed
264
265
       at 'gamma_rss'. Then the derivative uses the current 'gamma' for the {\tt Q}
       in the standard formula:
           derivative = gamma*(Y Q_inv^3 Y^T) - est_var * trace(Q_inv^2)
267
       where est_var is from 'estimate_variance_RSS(Y, Sigm, n, gamma_rss)'.
268
269
       n = X.shape[1]
270
       Sigm = W @ X # shape (N, n)
```

```
272
273
       # Build Q for *this gamma*:
       Q = (1.0/n) * (Sigm.T @ Sigm) + gamma*np.eye(n)
274
       Q_inv = np.linalg.inv(Q)
       Q_{inv2} = Q_{inv} @ Q_{inv}
276
277
       Q_{inv3} = Q_{inv2} @ Q_{inv}
       # We do the RSS-based variance at 'gamma_rss'
279
       # Must call 'estimate_variance_RSS':
280
281
       est_var = estimate_variance_RSS(Y, Sigm, n, gamma_rss)
282
283
       # derivative:
       term1 = gamma * (Y @ Q_inv3 @ Y.T)
                                              # shape (1,1)
284
       term2 = est_var * np.trace(Q_inv2)
285
       dL = term1 - term2
287
288
       # Convert to float
       if hasattr(dL, "item"):
289
           return float(dL.item())
290
       return float(dL)
292
293
294 def generate_synthetic_data_regression(npN, density, d=1, noise_level=0.1):
295
296
       Generates data for the model:
           Y = b^T * Sigma + eps
297
       with Sigma = W @ X and b has 'density'% of 1's.
298
299
       Shapes:
300
         X_train: (p, n_train)
301
302
         X_test: (p, n_test)
                   (N, p)
         W:
303
                             -- if d=1 => shape (N,1)
304
         b:
                   (N, d)
         Y_train, e_train: (d, n_train)
305
         Y_{test}, e_{test}: (d, n_{test})
306
307
       :param npN: (n_train, n_test, p, N)
308
       :param density: percentage for b's 1's
309
       :param d: dimension of output (1 => scalar outputs)
       :param noise_level: std-dev of Gaussian noise
311
       :return: (X_train, X_test, Y_train, Y_test, W, b, e_train, e_test)
312
313
       n_train, n_test, p, N = npN
314
315
       # 1) Generate X
316
       X_train = np.random.randn(p, n_train) # (p, n_train)
317
       X_test = np.random.randn(p, n_test) # (p, n_test)
318
319
       # 2) Generate b (N x d), if d=1 \Rightarrow (N,1)
320
       b = np.random.choice([0, 1],
                              size=(N, d),
322
                              p=[1 - density / 100, density / 100])
323
324
       # 3) Generate W (N x p)
325
       W = np.random.randn(N, p)
327
328
       def make_data(X):
329
           Given X with shape (p, n), returns:
330
             Y: (d, n) e: (d, n)
331
332
              Sigm: (N, n) [might be useful if needed]
333
334
           # Sigma = W @ X \Rightarrow (N, n)
335
336
           Sigm = sigma(W @ X)
           # Noise eps \Rightarrow (d, n)
338
           e = np.random.normal(loc=0, scale=noise_level, size=(d, X.shape[1]))
339
340
           # b^T shape => (d, N), so b^T @ Sigm => (d, n)
341
           Y = (b.T @ Sigm) + e
```

```
343
           return Y, e, Sigm
344
345
       # 4) Make training data (Y_train_full, e_train_full)
       Y_train_full, e_train_full, Sigm_train = make_data(X_train)
347
348
       # 5) Make test data (Y_test_full, e_test_full)
349
       Y_test_full, e_test_full, Sigm_test = make_data(X_test)
350
351
       # 6) Shuffle columns (the "sample" axis) in X, Y, e
352
            X is (p, n), Y, e are (d, n). We shuffle axis=1 for each.
353
354
       idx_train = np.random.permutation(n_train)
       idx_test = np.random.permutation(n_test)
355
356
357
       X_train = X_train[:, idx_train]
       Y_train_full = Y_train_full[:, idx_train]
358
       e_train_full = e_train_full[:, idx_train]
359
360
       X_test = X_test[:, idx_test]
361
       Y_test_full = Y_test_full[:, idx_test]
       e_test_full = e_test_full[:, idx_test]
363
364
       # 7) Return them all in consistent shapes
       \# No flattening is required; Y and e remain (d,n).
366
       return X_train, X_test, Y_train_full, Y_test_full, W, b, e_train_full, e_test_full
367
368
369
370 def find_optimal_gamma_grid(gammas, f, data, npN):
371
       Minimizes f(gamma, data, npN) by scanning over 'gammas'.
372
373
       Returns (best_gamma, best_value).
374
375
       vals = [f(g, data, npN) for g in gammas]
376
       i_min = np.argmin(vals)
377
       return float(gammas[i_min]), float(vals[i_min])
379
380 def find_zero_derivative_grid(gammas, deriv_func):
       0.000
       Finds gamma that yields derivative closest to zero in absolute value,
382
383
       scanning 'gammas'. Returns (best_gamma, best_deriv_value).
       The caller is expected to provide a function deriv_func(g).
384
385
386
       vals = [deriv_func(g) for g in gammas]
       abs_vals = np.abs(vals)
387
388
       i_min = np.argmin(abs_vals)
       return float(gammas[i_min]), float(vals[i_min])
389
390
391
392 def L_ridge(gamma, X, Y, W, b):
       n = X.shape[1]
393
       Sigm = sigma(W @ X)
395
       QyT = np.linalg.solve(Sigm.T@Sigm/n + gamma*np.eye(n), Y.T)
396
       beta = Sigm/n @ QyT
398
       numerator = np.linalg.norm(beta - b,'fro')**2
399
       denominator = np.linalg.norm(b, 'fro')**2
400
401
402
       return numerator/denominator
403
def find_optimal_gamma_grid(gammas, f, data, npN):
405
       Minimizes f(gamma, data, npN) by scanning over 'gammas'.
406
407
       Returns (best_gamma, best_value).
408
       best_gamma is a float, best_value is float.
409
       vals = [f(g, data, npN) for g in gammas]
410
411
       i_min = np.argmin(vals)
       return float(gammas[i_min]), float(vals[i_min])
412
```

```
414
415 if __name__ == "__main__":
416
       # ----- fixed experiment parameters -----
                           = 1e-3
       delta accuracy
418
       b_vector_density
                              = 50
                                              # (% of ones in true b)
419
                             = 8
       noise_level
                                              # of additive noise
420
                             = noise_level**2
       real_variance
421
422
       c1, c2
                             = 1.0, 0.8 # p/n and N/n
                             = [100, 200, 300]
423
       n_values
                              = 10
424
       num_repeats
                                       # averages per n
425
       gammas_grid
                             = [10**x for x in np.arange(-6, 3, 0.01)]
426
427
428
       # containers
                             = []
429
       n arr
       gEtrain_arr, gEtest_arr = [], []
430
431
       gOracle_arr, gEstVar_arr = [], []
432
       lrEtrain_arr, lrEtest_arr = [], []
lrOracle_arr, lrEstVar_arr = [], []
434
435
       # ----- main loop on sample size -----
436
       for n in n_values:
437
438
            # running sums to form averages across repeats
439
440
            sum\_gEtrain = sum\_gEtest = sum\_gOracle = sum\_gEstVar = 0.0
441
            sum_lrEtrain = sum_lrEtest = sum_lrOracle = sum_lrEstVar = 0.0
442
            for _ in range(num_repeats):
443
444
                # dimensions for this n
445
                p = int(c1 * n)
N = int(c2 * n)
446
447
                n_{train} = int(0.8 * n)
448
                n_{test} = n - n_{train}
449
                npN
                        = (n_train, n_test, p, N)
450
451
                # synthetic data
                X_tr, X_te, Y_tr, Y_te, W, b, _, _ = generate_synthetic_data_regression(
    npN, density=b_vector_density, d=1, noise_level=noise_level
453
454
455
                data = (X_tr, X_te, Y_tr, Y_te)
456
457
                # ----- four 's (grid search on the same grid) -----
458
                # 1) from \hat{E}_{t}
459
                g_Etrain, _ = find_optimal_gamma_grid(gammas_grid, Etrain_, data, npN)
460
461
462
                # 2) from \hat{E}_{t}
463
                g_Etest, _ = find_optimal_gamma_grid(gammas_grid, Etest_, data, npN)
464
                \# 3) from oracle derivative dL/d_{expected} = 0
465
                def d_oracle(g):    return dL_dy_expected(g, X_tr, Y_tr, W, real_variance)
466
                g_Oracle, _ = find_zero_derivative_grid(gammas_grid, d_oracle)
467
468
                # 4) from -datadriven derivative dL/d estimated = 0
469
470
                g_EstVar, _ = find_zero_derivative_grid(gammas_grid, d_est)
472
473
                # ridge errors
474
                lr_Etrain = L_ridge(g_Etrain, X_tr, Y_tr, W, b)
lr_Etest = L_ridge(g_Etest, X_tr, Y_tr, W, b)
lr_Oracle = L_ridge(g_Oracle, X_tr, Y_tr, W, b)
475
476
477
478
                lr_EstVar = L_ridge(g_EstVar, X_tr, Y_tr, W, b)
479
                # accumulate
480
                sum_gEtrain += g_Etrain; sum_lrEtrain += lr_Etrain
481
                sum_gEtest += g_Etest; sum_lrEtest += lr_Etest
sum_gOracle += g_Oracle; sum_lrOracle += lr_Oracle
482
483
                sum_gEstVar += g_EstVar; sum_lrEstVar += lr_EstVar
```

```
485
486
            # averages over repeats
487
            n_arr.append(n)
488
            gEtrain_arr.append(sum_gEtrain / num_repeats)
            gEtest_arr.append( sum_gEtest / num_repeats)
gOracle_arr.append(sum_gOracle / num_repeats)
489
490
            gEstVar_arr.append(sum_gEstVar / num_repeats)
491
492
            lrEtrain_arr.append(sum_lrEtrain / num_repeats)
493
            lrEtest_arr.append( sum_lrEtest / num_repeats)
494
            lrOracle_arr.append(sum_lrOracle / num_repeats)
495
496
            lrEstVar_arr.append(sum_lrEstVar / num_repeats)
497
            print(f"n={n:4d}_{\sqcup\sqcup}done.")
498
499
       # ----- save for later plotting -----
500
       np.savez("conv_analysis_n100_200_300.npz",
501
502
                 n_values
                                   = np.array(n_arr),
503
                 gamma_Etrain
                                    = np.array(gEtrain_arr),
                                 = np.array(gEtest_arr),
= np.array(gOracle_arr),
                 gamma_Etest
505
                 gamma_Oracle
506
                                   = np.array(gEstVar_arr),
                 gamma_EstVar
508
                 lr_Etrain
                                    = np.array(lrEtrain_arr),
509
                 lr_Etest
                                   = np.array(lrEtest_arr),
510
                 lr_Oracle
                                    = np.array(lrOracle_arr),
511
512
                 lr_EstVar
                                    = np.array(lrEstVar_arr),
513
                 noise_level
                                    = noise_level,
514
515
                 repeats
                                   = num_repeats,
                 b_vector_density = b_vector_density,
516
517
                 c1 = c1, c2 = c2
```

A.8. Fashion-MNIST validation

```
1 import numpy as np
2 from scipy import linalg
3 from tensorflow.keras.datasets import mnist,fashion_mnist
5 import time
7 SEED = 42
                                 # <-- one knob for reproducibility
8 np.random.seed(SEED)
9 random.seed(SEED)
10
11 def sigma(t):
      Small sigma function of choice
13
14
      :param t: input
      :return: output
15
16
17
      if activation_function == 'linear':
          return t
18
19
      if activation_function == 'ReLu':
20
          return np.maximum(t, 0)
21
22
      if activation_function == 'sign':
23
          return np.sign(t)
24
26
27 def K(x, y):
      Kernel function. Depends on the choice of a small sigma function :param x\colon first 'point' set
29
30
      :param y: second 'point' set
31
      :return: matrix of measures of 'distances' between points
32
33
      if activation_function == 'linear':
34
35
          return x.T@y
      if activation_function == 'ReLu':
37
38
          norm_x = np.linalg.norm(x, axis=0) # Shape (n_x,)
          norm_y = np.linalg.norm(y, axis=0) # Shape (n_y,)
39
40
41
          xTy = x.T @ y # Shape (n_x, n_y)
42
          norm_prod = norm_x[:, np.newaxis] * norm_y[np.newaxis, :] # Shape (n_x, n_y)
43
          norm_prod = np.maximum(norm_prod, 1e-9) # Avoid division by zero
45
          cos_theta = xTy / norm_prod
46
          cos_theta = np.clip(cos_theta, -1 + 1e-9, 1 - 1e-9) # Clamp values to [-1+1e-10, 1-1
47
               e-10]
48
          theta = np.arccos(-cos_theta) # Shape (n_x, n_y)
49
          sin_theta = np.sqrt(1 - cos_theta ** 2)
50
          return (norm_prod) / (2 * np.pi) * (cos_theta * theta + sin_theta)
52
53
      if activation_function == 'sign':
54
          norm_x = np.linalg.norm(x, axis=0) # Shape (n_x,)
55
          norm_y = np.linalg.norm(y, axis=0) # Shape (n_y,)
57
          xTy = x.T @ y # Shape (n_x, n_y)
58
          norm_prod = np.outer(norm_x, norm_y) # Shape (n_x, n_y)
60
          norm_prod = np.maximum(norm_prod, 1e-10) # Avoid division by zero
61
          cos_theta = xTy / norm_prod
63
          cos_{theta} = np.clip(cos_{theta}, -1 + 1e-10, 1 - 1e-10) # Clamp values to [-1+1e-10,
64
65
          return (2 / np.pi) * np.arcsin(cos_theta)
```

```
69 def K_(x, y, delta, N, n_train):
70
       Function to compute kernel approximation
       :param x: first 'point' set
72
73
       :param y: second 'point' set
       :param delta: delta parameter
       :return: approximation for the kernel matrix
75
76
      k_{-} = (N/n_{train})*((K(x, y))/(1+delta))
77
78
       return k_
79
80
81 def find_delta(gamma, K_train, N, n_train, accuracy) -> float:
       Fixed-point solver for delta without forming any matrix inverse.
83
84
       M(delta) = (N/n)*(K_train/(1+delta)) + gamma*I
85
       delta_next = (1/n) * trace( M(delta)^{-1}) @ K_train )
86
       delta_prev = 1.0
       delta_next = 0.0
88
89
       I = np.eye(n_train)
       while abs(delta_prev - delta_next) > accuracy:
91
92
           delta_prev = delta_next
           a = (N / n_train) / (1.0 + delta_prev)
                                                                  # scaling factor
93
94
           M = a * K_train + gamma * I
                                                                  # (n x n), SPD
95
           cf = linalg.cho_factor(M, lower=True, check_finite=False)
           # U = M^{-1} K_train (solve for all columns at once)
96
           U = linalg.cho_solve(cf, K_train, check_finite=False)
97
98
           delta_next = (1.0 / n_train) * np.trace(U)
99
100
       return float(delta_next)
102
103 def find_delta_from_eigs(gamma, eigvals, N, n_train, accuracy) -> float:
104
       Fixed-point iteration for delta using only eigenvalues of K_train.
105
       a = (N/n) / (1+delta)
       delta_next = (1/n) * sum_i [ lambda_i / (a*lambda_i + gamma) ]
107
108
       delta_prev = 1.0
109
       delta_next = 0.0
110
111
       while abs(delta_prev - delta_next) > accuracy:
           delta_prev = delta_next
112
           a = (N / n_{train}) / (1.0 + delta_{prev})
113
           delta_next = (1.0 / n_train) * float(np.sum(eigvals / (a * eigvals + gamma)))
       return float(delta_next)
115
116
117
118 def Etest_(gamma, data, npN):
119
       Eigen/SVD-based E_test (no explicit inverses).
120
       Expects 'data' to be a dict:
121
           {
             "K_train": (n,n),
123
             "K_xX":
124
                        (n,n_test),
             "K_XX":
125
                         (n_test, n_test),
             "Y_train": (n,),
126
             "Y_test": (n,),
127
128
           }
129
       npN = (n_train, n_test, p, N) # N here is ignored; we take N from 'data'
130
131
132
       n_{train}, n_{test}, p, _{n} = npN
133
       K_train = data["K_train"]
134
135
       K_xX
               = data["K_xX"]
               = data["K_XX"]
136
       K_XX
               = data["Y_train"]
137
       y_tr
   y_te = data["Y_test"]
```

```
N = int(data["N"])
139
140
       # 1) Eigendecomposition of K_{train} (SPD \rightarrow eigh is ideal)
141
       lam, V = np.linalg.eigh(K_train) # lam: (n,), V: (n,n)
143
144
       # 2) Solve for delta using only eigenvalues
       delta = find_delta_from_eigs(gamma, lam, N, n_train, delta_accuracy)
145
146
147
       # Common scaling
       a = (N / n_{train}) / (1.0 + delta)
148
                                                 # scalar
       d = a * lam + gamma
149
                                                 # (n,)
150
       # 3) Prediction term
151
       y_{til} = V.T @ y_{tr}
                                                 # (n.)
152
       u = V @ (y_til / d)
153
                                                 # (n.)
       K_xX_scaled = a * K_xX
                                                 # (n, n_test)
154
       resid = y_te - K_xX_scaled.T @ u
                                                 # (n_test,)
155
       term_pred = (1.0 / n_test) * float(resid @ resid)
156
157
       # 4) Ratio term: num / denom
       \# \text{ num} = (1/N) * \text{ sum_i} (a*lam_i / d_i^2) * (y_til_i)^2
159
       num = (1.0 / N) * float(np.sum((a * lam / (d * d)) * (y_til * y_til)))
160
       # denom = (1 - 1/N) * sum_i (a^2 * lam_i^2 / d_i^2)
       denom = (1.0 - 1.0 / N) * float(np.sum((a * a) * (lam * lam) / (d * d)))
162
163
       # 5) Bracket term
164
165
       # First piece: tr(a*K_XX)
       tr_aKXX = a * float(np.trace(K_XX))
       # Second piece: tr(Q S) + gamma tr(Q^2 S), with S = (a K_xX)(a K_xX)^T
167
       # If Z = V^T (a K_xX), then diag(V^T S V) = rowwise sum of Z^2
168
       Z = V.T @ K_xX_scaled
                                                 # (n, n_test)
       diag_Sprime = np.sum(Z * Z, axis=1)
                                                 # (n,)
170
171
       tr_QS = float(np.sum(diag_Sprime / d))
172
       tr_Q2S = float(np.sum(diag_Sprime / (d * d)))
       bracket = (tr_aKXX - (tr_QS + gamma * tr_Q2S)) / n_test
173
       # 6) Combine
175
       E_test_ = term_pred + (num / denom) * bracket
176
       return float(E_test_)
178
179
def estimate_variance_RSS(Y, Sigm, n, gamma):
181
182
       RSS-based noise variance estimate using SVD of A = Sigm / sqrt(n).
       Equivalent to solving (Sigm^T Sigm / n + gamma I) z = Y, but faster and
183
       numerically stable.
184
       Returns a Python float.
186
187
       A = Sigm / np.sqrt(n)
188
                                                        # (N, n)
       U, S, Vt = np.linalg.svd(A, full_matrices=False) # S: (r,), Vt: (r, n), r = rank
189
       y = np.asarray(Y).reshape(-1)
                                                        # (n,)
190
       y_v = Vt @ y
                                                        # coords in span(A)
191
       r = S.shape[0]
192
       denom = S**2 + gamma
                                                        # (r,)
194
       # Residual in span(A): y_v - (S^2/(S^2+gamma)) y_v = (gamma/(S^2+gamma)) y_v
195
       rss_span = np.sum((gamma * y_v / denom)**2)
196
197
       # Residual in null(A): unchanged (no fit there)
198
       if r < y.size:</pre>
199
           y_perp = y - Vt.T @ y_v
200
           rss_null = float(y_perp @ y_perp)
201
       else:
202
203
           rss_null = 0.0
       variance = (rss_span + rss_null) / n
205
       return float(variance)
206
207
208
209 def dL_dy_estimated(gamma, X, Y, W, variance=1):
```

```
0.00
210
       Optimized derivative:
211
          dL/dgamma = gamma * (Y Q^{-3} Y^T) - est_var * tr(Q^{-2}),
212
       with Q = (Sigm^T Sigm)/n + gamma I, Sigm = sigma(W @ X).
213
214
215
       Uses SVD of A = Sigm / sqrt(n) to avoid explicit inverses.
216
       # Design in hidden layer
217
218
       Sigm = sigma(W @ X)
                                          # (N, n)
       n = X.shape[1]
219
220
       # SVD of A = Sigm / sqrt(n) -> Q = V diag(S^2 + gamma) V^T
221
       A = Sigm / np.sqrt(n)
222
        \label{eq:continuous} \mbox{U, S, Vt = np.linalg.svd(A, full_matrices=False)} \quad \mbox{\# U:(N,r), S:(r,), Vt:(r,n)} 
223
       r = S.shape[0]
       denom = S**2 + gamma
225
226
227
       \# Coordinates of y in the V basis
       y = np.asarray(Y).reshape(-1) # (n,)
228
                                          # (r,)
       y_v = Vt @ y
       if r < n:
230
           y_perp = y - Vt.T @ y_v
231
           y_perp_norm2 = float(y_perp @ y_perp)
       else:
233
234
           y_perp_norm2 = 0.0
235
       # First term: gamma * y^T Q^{-3} y
236
       term1 = gamma * ( np.sum((y_v**2) / (denom**3)) + (y_perp_norm2 / (gamma**3)) )
237
238
       # Noise variance estimate (same as before but SVD-based and exact)
239
240
       est_var = estimate_variance_RSS(y, Sigm, n, gamma) # returns float
241
242
       # Second term: est_var * tr(Q^{-2})
       tr_Qinv2 = np.sum(1.0 / (denom**2)) + (n - r) * (1.0 / (gamma**2))
243
244
       term2 = est_var * tr_Qinv2
245
       return float(term1 - term2)
246
247
249 def dL_dy_estimated_fixed_rss(gamma, X, Y, W, variance, gamma_rss):
250
251
       "Estimated" derivative: dL/dgamma, but the RSS-based variance is computed
       at 'gamma_rss'. Then the derivative uses the current 'gamma' for the {\tt Q}
252
253
       in the standard formula:
           derivative = gamma*(Y Q_inv^3 Y^T) - est_var * trace(Q_inv^2)
254
       where est_var is from 'estimate_variance_RSS(Y, Sigm, n, gamma_rss)'.
255
256
       n = X.shape[1]
257
       Sigm = sigma(W @ X) # shape (N, n)
258
259
       # Build Q for *this gamma*:
260
       Q = (1.0/n) * (Sigm.T @ Sigm) + gamma*np.eye(n)
261
       Q_inv = np.linalg.inv(Q)
262
       Q_inv2 = Q_inv @ Q_inv
263
       Q_inv3 = Q_inv2 @ Q_inv
264
265
266
       # We do the RSS-based variance at 'gamma_rss'
       # Must call 'estimate_variance_RSS':
267
268
       est_var = estimate_variance_RSS(Y, Sigm, n, gamma_rss)
269
       # derivative:
270
       term1 = gamma * (Y @ Q_inv3 @ Y.T)
                                               # shape (1,1)
271
       term2 = est_var * np.trace(Q_inv2)
272
       dL = term1 - term2
273
274
       # Convert to float
       if hasattr(dL, "item"):
276
           return float(dL.item())
277
278
       return float(dL)
279
```

```
281 def make_two_class_ridge_from_arrays(
                                        # e.g. (60000, 28, 28) uint8
# e.g. (60000,) int
       init_data: np.ndarray,
282
       init_labels: np.ndarray,
283
       selected_labels: list[int],
                                         # exactly two labels, e.g. [1, 2]
284
       c1: float,
                                         # target p/n ratio (features-to-sample-size); n round
285
           (p / c1)
       train_fraction: float = 0.8,
                                        # desired train share of (train + test)
       cs: tuple[float, float] = (0.5, 0.5), # per-class proportions (sum 1.0)
287
                                         # RNG seed (None -> non-deterministic like authors)
288
       seed: int | None = None
289 ):
290
291
       Rebuild (X, y, X_{test}, y_{test}) for a 2-class MNIST/Fashion setup, mirroring the authors:
        1) Sort by label.
292
293
         2) Flatten to (p, init_n), rescale to [0,1].
         3) Global mean-center and scale so average ||x||^2 = p.
         4) Select the two classes, pool them, mean-center & scale AGAIN within the pooled
295
             subset so
            pooled average ||x||^2 = p.
296
         5) For each class: shuffle WITHIN class; TRAIN = first int(cs[i]*n);
297
            TEST = columns [n : n + int(cs[i]*n_test)] (i.e., test slice starts at absolute
                index n).
         6) Labels: class 0 \rightarrow -1, class 1 \rightarrow +1; class blocks contiguous.
299
300
       Differences from the raw authors' script:
301
         - You control n via c1 (p/n), and we derive n_test from `train_fraction`:
302
              n_test = round( n * (1 - train_fraction) / train_fraction )
303
304
          This preserves their splitting scheme while giving you a clean 80/20 (or any) split.
305
       # -----
306
       # A) Sort by label (authors do this)
307
       idx_sorted = np.argsort(np.array(init_labels))
309
310
       labels_sorted = np.array(init_labels)[idx_sorted]
311
       imgs_sorted = np.array(init_data)[idx_sorted] # (init_n, H, W)
312
       # -----
       # B) Flatten to (p, init_n), cast to float
314
315
       H, W = imgs_sorted.shape[1], imgs_sorted.shape[2]
      p = H * W
317
       init_n = imgs_sorted.shape[0]
318
       data = imgs_sorted.reshape(init_n, p).T.astype(np.float64) # (p, init_n); columns =
319
          samples
320
       # -----
321
       # C) Global rescale to [0, 1] (authors: data = data / data.max())
322
       max_val = data.max()
324
325
       if max_val > 0:
326
          data /= max_val
327
       # -----
328
       # D) Global mean-center & renormalize so avg ||x||^2 = p
329
330
       mean_data = np.mean(data, axis=1, keepdims=True)
                                                                # (p,1)
       centered_global = data - mean_data
                                                                # (p, init_n)
332
       norm2_data = np.mean(np.sum(centered_global**2, axis=0)) # scalar
333
       scale_global = np.sqrt(p) / np.sqrt(norm2_data) if norm2_data > 0 else 1.0
334
       data_std = centered_global * scale_global
335
                                                                # (p, init_n)
336
       # -----
337
       # E) Extract the two classes and pool them
338
339
       if len(selected_labels) != 2 or len(cs) != 2:
340
341
          raise ValueError("Provide_exactly_two_labels_and_two_class_proportions_cs=(c0,_c1).")
342
       selected data = []
343
       for lab in selected_labels:
344
          cols = (labels_sorted == lab)
345
           selected_data.append(data_std[:, cols])
                                                                # each: (p, n_class)
346
```

```
# Pool both classes (fix authors' missing assignment in np.concatenate)
348
       cascade_selected = np.concatenate(selected_data, axis=1) # (p, n_pool)
349
350
       # F) Recenter & renormalize AGAIN within the pooled two-class subset
352
353
       mean_pool = np.mean(cascade_selected, axis=1, keepdims=True) # (p,1)
354
       centered_pool = cascade_selected - mean_pool
355
       norm2_pool = np.mean(np.sum(centered_pool**2, axis=0))
356
       scale_pool = np.sqrt(p) / np.sqrt(norm2_pool) if norm2_pool > 0 else 1.0
357
358
359
       for j in range(2):
            selected_data[j] = (selected_data[j] - mean_pool) * scale_pool
360
361
362
       # G) Choose n from c1 = p/n -> n round(p / c1)
363
364
            Then compute n_{test} from the desired train fraction:
365
              n / (n + n_test) = train_fraction => n_test = n * (1 - tf) / tf
366
       if c1 <= 0:</pre>
            raise ValueError("c1_{\square}must_{\square}be_{\square}positive_{\square}(c1_{\square}=_{\square}p/n).")
368
       if not (0 < train_fraction < 1):</pre>
369
            raise ValueError("train_fraction_must_be_in_(0, 1).")
371
       n = int(np.round(p / c1))
372
       n = \max(1, n)
373
       n_test = int(np.round(n * (1.0 - train_fraction) / train_fraction))
374
       n_{test} = max(1, n_{test}) # ensure at least one test column overall
375
       c1_actual = p / n
376
       train\_share\_actual = n / (n + n\_test)
377
379
380
       \# H) Allocate X, X_test and compute per-class block sizes as authors do
381
382
       rng = np.random.default_rng(seed)
       X = np.zeros((p, n), dtype=np.float64)
       X_test = np.zeros((p, n_test), dtype=np.float64)
384
385
       cs_arr = np.array(cs, dtype=float)
386
       \label{eq:train_bounds} \mbox{ = (np.cumsum(np.concatenate([[0.0], cs_arr])) * n).astype(int)}
387
       test_bounds = (np.cumsum(np.concatenate([[0.0], cs_arr])) * n_test).astype(int)
388
389
       train_counts = np.diff(train_bounds)
                                                    # [int(cs[0]*n), int(cs[1]*n)]
390
391
       test_counts = np.diff(test_bounds)
                                                     # [int(cs[0]*n_test), int(cs[1]*n_test)]
392
       # Each class must have >= n + test_counts[i] samples for the authors' indexing
393
       for i in range(2):
            need = n + test_counts[i]
                                                               # because test slice is data_i[:, n : n
395
                + test_counts[i]]
396
            have = selected_data[i].shape[1]
            if have < need:</pre>
397
                raise ValueError(
398
                    f"Class_{\{i\}_{\sqcup}}(label_{\sqcup}\{selected\_labels[i]\})_{\sqcup}has_{\sqcup}\{have\}_{\sqcup}samples;_{\sqcup}need_{\sqcup}at_{\sqcup}least_{\sqcup}\{
399
                         need}..."
                     f"to_{\sqcup}reproduce_{\sqcup}the_{\sqcup}authors'_{\sqcup}exact_{\sqcup}slicing_{\sqcup}(test_{\sqcup}starts_{\sqcup}at_{\sqcup}index_{\sqcup}n)."
400
401
402
403
       # ------
       \mbox{\tt\#} I) Fill X and X_test: shuffle WITHIN each class, then slice as authors do
404
405
       for i in range(2):
406
            perm = rng.permutation(selected_data[i].shape[1])
407
            data_i = selected_data[i][:, perm]
408
409
            # Train block for class i
410
            X[:, train_bounds[i]:train_bounds[i+1]] = data_i[:, :train_counts[i]]
411
412
            # Test block for class i (note: slice starts at absolute index n)
413
414
            X_test[:, test_bounds[i]:test_bounds[i+1]] = data_i[:, n : n + test_counts[i]]
415
       # -----
```

```
# J) Labels: class 0 -> -1, class 1 -> +1; contiguous class blocks
417
418
        y = np.concatenate([
419
            -np.ones(train_counts[0], dtype=np.int8),
420
            +np.ones(train_counts[1], dtype=np.int8),
421
        1)
422
        y_test = np.concatenate([
423
            -np.ones(test_counts[0], dtype=np.int8),
424
            +np.ones(test_counts[1], dtype=np.int8),
425
426
427
428
        info = dict(
            p=p, n=n, n_test=n_test,
429
            selected_labels=tuple(selected_labels),
430
            class_sizes=tuple(d.shape[1] for d in selected_data),
431
            train_counts=tuple(train_counts.tolist()),
432
433
            test_counts=tuple(test_counts.tolist()),
            c1_target=c1,
434
435
            c1_actual=c1_actual,
            {\tt train\_fraction\_target=train\_fraction}\,,
            train_fraction_actual=round(train_share_actual, 6),
437
438
            global_scale=scale_global,
            pooled_scale=scale_pool,
            \verb"note="Exact_{\sqcup} \verb"authors" \verb||| preprocessing; \verb||| test_{\sqcup} \verb|slice_{\sqcup} starts_{\sqcup} \verb|at_{\sqcup} index_{\sqcup} \verb|n_{\sqcup} within_{\sqcup} each_{\sqcup} class."
440
441
442
        return X, y, X_test, y_test, info
443
444
def find_optimal_gamma_grid(gammas, f, data, npN):
446
447
        Minimizes f(gamma, data, npN) by scanning over 'gammas'.
        Returns (best_gamma, best_value).
448
449
450
        vals = [f(g, data, npN) for g in gammas]
        i_min = np.argmin(vals)
451
        return float(gammas[i_min]), float(vals[i_min])
452
453
454
455 def find_zero_derivative_grid(gammas, deriv_func):
456
457
        Finds gamma that yields derivative closest to zero in absolute value,
458
        scanning 'gammas'. Returns (best_gamma, best_deriv_value).
        The caller is expected to provide a function deriv\_func(g).
459
460
        vals = [deriv_func(g) for g in gammas]
461
462
        abs_vals = np.abs(vals)
        i_min = np.argmin(abs_vals)
463
        return float(gammas[i_min]), float(vals[i_min])
464
465
466
def get_beta(W, X_train, Y_train, gamma):
468
        SVD version (optional):
469
          Solve z = (\widetilde{Sigm}^T \operatorname{Sigm} / n + \operatorname{gamma} I)^{-1} Y via SVD of A = Sigm / \operatorname{sqrt}(n).
470
         Then beta = (Sigm / n) @ z.
471
472
        Sigm = sigma(W @ X_train)
473
                                                        # (N, n)
        n = X_train.shape[1]
474
        A = Sigm / np.sqrt(n)
475
                                                        # (N. n)
476
        # economy SVD
477
        U, S, Wt = np.linalg.svd(A, full_matrices=False) # U:(N,r), S:(r,), Vt:(r,n), r=min(N,n)
478
479
        y_v = Vt @ Y_train
                                                          # (r,)
480
        z_{span} = (Vt.T * (1.0 / (S**2 + gamma))) @ y_v # V diag(1/(S^2+)) V^T Y in span(A)
481
        if Vt.shape[0] < n:</pre>
                                                          # nullspace component (if n > r)
483
            y_perp = Y_train - Vt.T @ y_v
484
485
            z = z_{span} + (1.0 / gamma) * y_{perp}
486
        else:
          z = z_{span}
```

```
488
                                                   # (N,)
489
       beta = (Sigm / n) @ z
       return beta
490
491
492
493
494 def build_kernels(X_train, X_test):
495
       Compute unscaled kernels once. Scaling by 'a' happens inside Etest_.
496
497
       K_train = K(X_train, X_train)
498
       K_xX = K(X_train, X_test)
K_XX = K(X_test, X_test)
499
500
       return K_train, K_xX, K_XX
501
503
504
505 if __name__ == "__main__":
       import os
506
      t0 = time.perf counter()
508
509
       # ----- Experiment knobs -----
       511
                               = 1e-3
512
       delta_accuracy
                               = 1.0
513
514
                               = 1.0
                               = [10**x for x in np.arange(-6, 4, 0.01)]
       gammas_grid_full
515
                               = 'fashion' # 'fashion' | 'MNIST'
       testcase
516
                               = 0.8
      train_fraction
517
518
      RNG_SCENARIO_SEED
                               = SEED + 123
519
520
      N_PIX_SCEN
                               = 10
                                              # how many single-pixel scenarios
                               = 5
                                              # how many 5x5 patch scenarios
521
       N_PATCH_SCEN
       PATCH_SIZE
                               = 5
522
       DELTA_STD
                               = 0.5
                                             # additive perturbation in standardized space
       LOCAL_LOG_WIDTH
                               = 1.0
                                              # +- decades around baseline
524
                               = 0.05
       LOCAL LOG STEP
                                              # grid resolution in log10
525
       RESULTS_DIR
                               = "results"
       os.makedirs(RESULTS_DIR, exist_ok=True)
527
528
529
       # --- Load data (two classes) ---
if testcase == 'MNIST':
530
531
          selected_labels = [7, 9]
532
533
           (init_data, init_labels), _ = mnist.load_data()
534
           selected_labels = [1, 2]
535
           (init_data, init_labels), _ = fashion_mnist.load_data()
536
537
       # --- Authors' preprocessing / split (deterministic with SEED) ---
538
       X_train, y_train, X_test, y_test, info = make_two_class_ridge_from_arrays(
539
           init_data=init_data,
540
           init labels=init labels.
541
           selected_labels=selected_labels,
           c1=c1,
543
544
           train_fraction=train_fraction,
          seed=SEED
545
546
547
       # --- Sizes & model shapes ---
548
       p = info['p']
549
       n_train = info['n']
550
       n_test = info['n_test']
551
               = int(round(c2 * n_train))
552
                = np.random.randn(N, p)
                                                  # fixed W for all scenarios
553
554
       # Sanity
555
556
       assert X_train.shape == (p, n_train)
       assert X_test.shape == (p, n_test)
557
assert y_train.shape == (n_train,)
```

```
assert y_test.shape == (n_test,)
559
560
       # --- Precompute baseline kernels once ---
561
       K_train_base, K_xX_base, K_XX_base = build_kernels(X_train, X_test)
563
       # --- Baseline: gammas (full grid once), betas, and caches ---
564
       data_base = {
565
           "K_train": K_train_base,
566
           "K_xX":
567
                      K_xX_base,
           "K_XX":
568
                     K_XX_base,
           "Y_train": y_train,
569
           "Y_test": y_test,
570
           "N":
571
572
573
       npN = (n_train, n_test, p, N)
574
575
       gamma_Etest_base, _ = find_optimal_gamma_grid(gammas_grid_full, Etest_, data_base, npN)
       def d_est(g): return dL_dy_estimated(g, X_train, y_train, W)
576
       gamma_direct_base, _ = find_zero_derivative_grid(gammas_grid_full, d_est)
577
       beta_Etest_base = get_beta(W, X_train, y_train, gamma_Etest_base)
579
       beta_direct_base = get_beta(W, X_train, y_train, gamma_direct_base)
580
       # Helpers: local gamma grids around baseline (± LOCAL_LOG_WIDTH decades)
582
       def local_grid(g0, width=LOCAL_LOG_WIDTH, step=LOCAL_LOG_STEP):
583
           g0 = float(g0)
584
585
           g0_log = np.log10(g0)
           return 10.0**np.arange(g0_log - width, g0_log + width + 1e-12, step)
586
587
       # Random generator for scenario design
588
589
       rng = np.random.default_rng(RNG_SCENARIO_SEED)
590
591
       # --- Utility: run a single scenario given a perturbed X_train ---
592
       def run_scenario(X_train_pert, scenario_meta, tag):
           # Kernels for this scenario
593
           K_tr, K_xX, K_XX = build_kernels(X_train_pert, X_test)
           data_s = {
595
               "K_train": K_tr,
596
               "K_xX": K_xX,
               "K_XX":
                          K_XX,
598
               "Y_train": y_train,
599
               "Y_test": y_test,
600
               "N":
601
602
           }
603
604
           # Local gamma search around baseline gammas (fast, robust)
           gE_grid = local_grid(gamma_Etest_base)
           gD_grid = local_grid(gamma_direct_base)
606
607
608
                _ = find_optimal_gamma_grid(gE_grid, Etest_, data_s, npN)
           def d_est_local(g): return dL_dy_estimated(g, X_train_pert, y_train, W)
609
           gD, _ = find_zero_derivative_grid(gD_grid, d_est_local)
610
611
           # Betas
612
           beta_E = get_beta(W, X_train_pert, y_train, gE)
           beta_D = get_beta(W, X_train_pert, y_train, gD)
614
615
616
           # Save all artifacts
617
           out = {
               "testcase": testcase,
618
               "activation_function": activation_function,
619
               "selected_labels": np.array(selected_labels, dtype=int),
620
               "SEED": SEED,
621
               "rng_scenario_seed": RNG_SCENARIO_SEED,
622
623
               "p": p, "n_train": n_train, "n_test": n_test, "N": N,
               "c1": c1, "c2": c2, "train_fraction": train_fraction,
625
               "gamma_Etest_base": float(gamma_Etest_base),
626
627
               "gamma_direct_base": float(gamma_direct_base),
               "gamma_Etest_scen": float(gE),
628
               "gamma_direct_scen": float(gD),
```

```
630
               "beta_Etest_base": beta_Etest_base,
631
               "beta_direct_base": beta_direct_base,
632
               "beta_Etest_scen": beta_E,
               "beta_direct_scen": beta_D,
634
635
               "scenario_meta": scenario_meta
636
           }
637
638
           fname = os.path.join(RESULTS_DIR, f"{tag}.npz")
           np.savez(fname, **out)
639
           print(f"[saved]_|{fname}")
640
641
       # --- 10 single-pixel scenarios (std-space) ---
642
643
       # choose unique pixels
       pix_ids = rng.choice(p, size=N_PIX_SCEN, replace=False)
644
       # random signs per scenario (+1 or -1)
645
       pix_signs = rng.choice([-1.0, 1.0], size=N_PIX_SCEN)
646
647
       for idx, (pix, sgn) in enumerate(zip(pix_ids, pix_signs), start=1):
648
           Xp = X_train.copy()
           Xp[pix, :] += sgn * DELTA_STD
650
651
           meta = {
               "type": "pixel",
652
               "feature_indices": np.array([int(pix)], dtype=int),
653
654
               "delta_std": float(DELTA_STD),
               "sign": float(sgn),
655
               "applied_space": "standardized_training_only"
656
           }
657
           tag = f"pert_pix_{idx:02d}_{testcase}_af-{activation_function}_pix{int(pix)}_d{
658
                DELTA_STD:+.2f}sgn{int(sgn)}_seed{SEED}"
           run_scenario(Xp, meta, tag)
660
661
       # --- 5 patch (5x5) scenarios (std-space) ---
       # infer H, W (MNIST/Fashion are 28x28; derive from p if possible)
662
663
       H_px = W_px = int(round(np.sqrt(p)))
       if H_px * W_px != p:
           H_px = W_px = 28
665
666
       max_r = H_px - PATCH_SIZE
       max_c = W_px - PATCH_SIZE
668
669
       rc_pairs = set()
       while len(rc_pairs) < N_PATCH_SCEN:</pre>
670
           rc_pairs.add((int(rng.integers(0, max_r + 1)), int(rng.integers(0, max_c + 1))))
671
672
       rc_pairs = list(rc_pairs)
       patch_signs = rng.choice([-1.0, 1.0], size=N_PATCH_SCEN)
673
674
       for idx, ((r0, c0), sgn) in enumerate(zip(rc_pairs, patch_signs), start=1):
675
           # build flat indices of the patch
676
677
           feats = []
678
           for dr in range(PATCH_SIZE):
               rr = r0 + dr
679
               cc0 = c0
680
               row_start = rr * W_px
681
               feats.extend(range(row_start + cc0, row_start + cc0 + PATCH_SIZE))
682
           feats = np.array(feats, dtype=int)
684
           Xp = X_train.copy()
685
           Xp[feats, :] += sgn * DELTA_STD
686
           meta = {
687
               "type": "patch",
688
               "patch_size": int(PATCH_SIZE),
689
               "top_left_rc": (int(r0), int(c0)),
690
               "feature_indices": feats,
               "delta_std": float(DELTA_STD),
692
693
               "sign": float(sgn),
                "applied_space": "standardized_training_only"
694
           }
695
           tag = (f"pert_patch_{idx:02d}_{testcase}_af-{activation_function}_r{r0}c{c0}"
696
697
                  f"
                     _sz{PATCH_SIZE}_d{DELTA_STD:+.2f}sgn{int(sgn)}_seed{SEED}")
           run_scenario(Xp, meta, tag)
698
```

```
total_sec = time.perf_counter() - t0

print(f"[TOTAL]_\{total_sec:.3f}_\s")

print(f"Baseline_gammas:_Etest={gamma_Etest_base:.4g},_\direct={gamma_direct_base:.4g}")
```