

RISK PHASE TRANSITIONS IN SPIKED REGRESSION: ALIGNMENT DRIVEN BENIGN AND CATASTROPHIC OVERFITTING

Jiping Li

Department of Mathematics
University of California, Los Angeles
jipingli0324@g.ucla.edu

Rishi Sonthalia

Department of Mathematics
Boston College
rishi.sonthalia@bc.edu

ABSTRACT

This paper analyzes the generalization error of minimum-norm interpolating solutions in linear regression using spiked covariance data models. The paper characterizes how varying spike strengths and target-spike alignments can affect risk, especially in overparameterized settings. The study presents an exact expression for the generalization error, leading to a comprehensive classification of benign, tempered, and catastrophic overfitting regimes based on spike strength, the aspect ratio $c = d/n$ (particularly as $c \rightarrow \infty$), and target alignment. Notably, in well-specified aligned problems, increasing spike strength can surprisingly induce catastrophic overfitting before achieving benign overfitting. The paper also reveals that target-spike alignment is not always advantageous, identifying specific, sometimes counterintuitive, conditions for its benefit or detriment. Alignment with the spike being detrimental is empirically demonstrated to persist in nonlinear models.

1 INTRODUCTION

Understanding the generalization error of overparameterized models is a central challenge in modern machine learning. Phenomena such as double descent (Belkin et al., 2019; Hastie et al., 2022) and benign overfitting (Bartlett et al., 2020; Mallinar et al., 2022; Tsigler and Bartlett, 2023) have spurred research underscoring the critical role of the data’s spectral structure (Bartlett et al., 2020; Dobriban and Wager, 2018; Hastie et al., 2022; Kausik et al., 2024; Mei et al., 2022; Sonthalia and Nadakuditi, 2023; Tsigler and Bartlett, 2023; Wang et al., 2024a). The spiked covariance model is one commonly considered spectral structure Couillet and Liao (2022). In this model, the data matrix $\mathbf{X} = \mathbf{Z} + \mathbf{A} \in \mathbb{R}^{d \times n}$, comprising n data points in \mathbb{R}^d , is decomposed into a rank-one signal component (“spike”) \mathbf{Z} and an isotropic noise component (“bulk”) \mathbf{A} . Spiked covariance models emerge naturally in practice, for instance, in the features learned by neural networks during training Sonthalia et al. (2025); Ba et al. (2022; 2023); Damian et al. (2022); Dandi et al. (2024); Martin and Mahoney (2021); Moniri et al. (2023); Wang et al. (2024b). While recent studies have examined benign overfitting in spiked models (Ba et al., 2023; Kausik et al., 2024), they lack a systematic taxonomy spanning spike strength, target–spike alignment, model misspecification, and train–test covariate shift. This paper closes the gap for linear regression.

This work explores how general spike sizes and target alignments affect generalization error in least squares linear regression. We consider targets y generated by:

$$y = \alpha_Z \beta_*^\top \mathbf{z} + \alpha_A \beta_*^\top \mathbf{a} + \varepsilon$$

Here, $\mathbf{z} \in \mathbb{R}^d$ represents the signal (spike) component, $\mathbf{a} \in \mathbb{R}^d$ corresponds to the noise (bulk) component, ε is observation noise, and $\beta_* \in \mathbb{R}^d$. The coefficients α_Z and α_A model the target’s dependence on the spike and bulk components. Notably, if $\alpha_A \neq \alpha_Z$, the targets are non-linear functions of $\mathbf{x} = \mathbf{z} + \mathbf{a}$, introducing model mis-specification. We address two fundamental questions:

- **Q1:** For a fixed aspect ratio $c = d/n$, in the asymptotically proportional regime, under what conditions does alignment of the target signal with the data spike improve or impair generalization?

- **Q2:** In the high-dimensional limit where $c \rightarrow \infty$, when do we observe benign, tempered, or catastrophic overfitting regimes?

Contributions. We present precise characterization of the generalization performance of minimum-norm interpolating solutions in linear regression. Our exact risk decomposition pinpoints conditions for transitions between benign and catastrophic overfitting. This reveals alignment-dependent phenomena obscured by isotropic theories, clarifying how signal structure, data scaling, and overparameterization shape generalization. Our primary contributions are as follows:

- **Precise Risk Characterization:** We derive an exact generalization error decomposition (Theorem 5) into interpretable bias, variance, data noise, and alignment terms.
- **Comprehensive Categorization of Overfitting Regimes:** We precisely classify benign, tempered, or catastrophic overfitting regimes based on spike strength, overparameterization ($c = d/n$), and target alignment (Table 1). Surprisingly, for well-specified aligned problems, increasing spike strength can induce catastrophic overfitting before achieving benign overfitting. Misspecified problems show distinct transitions, often precluding benign overfitting.
- **Conditions for Beneficial Alignment:** Challenging conventional wisdom, we show spike alignment is not always beneficial and depends on spike strength meeting critical thresholds (Table 2). For misspecified problems, beneficial alignment requires α_Z/α_A in a specific, non-trivial range. Counterintuitively, very strong spike dependence (α_Z/α_A) can render alignment detrimental.
- **Empirical Validation:**¹ Empirical validation confirms our theoretical phenomena, including surprising negative alignment impacts, persist in nonlinear models, underscoring broader relevance.

Benign Overfitting in Linear Regression. Significant research has explored benign overfitting in linear regression (Bartlett et al., 2020; Cao et al., 2021; Chatterji and Long, 2021; Karhadkar et al., 2024; Koehler et al., 2021; Liang and Rakhlin, 2020; Mallinar et al., 2022; Muthukumar et al., 2020; Shamir, 2022; Tsigler and Bartlett, 2023; Wu and Xu, 2020). Many studies assume a uniformly bounded largest covariance eigenvalue or lack precise characterizations of its interplay with target alignment and generalization. *Our work allows this eigenvalue to grow, offering precise performance characterizations based on this growth and alignment.* While Kausik et al. (2024) considers spiked models, their focus is on noiseless, well-specified scenarios with specific spike scaling. *Our analysis is broader, encompassing observation noise, misspecification, and general spike scaling.*

Many prior works (Karhadkar et al., 2024; Shamir, 2022; Tsigler and Bartlett, 2023) on benign overfitting with low-rank signals plus isotropic noise require near-orthogonality between signal and noise, sometimes imposing strong conditions like $d = \Omega(n^2 \log n)$. *We instead consider the proportional regime $d/n \rightarrow c = \Theta(1)$, subsequently examining $c \rightarrow \infty$.* This setting is morally similar to allowing $d = \omega(n)$ and aligns with approaches like (Karhadkar et al., 2024) which, for classification, shows misclassification probability can be upper bounded by $Ce^{-d/n}$, vanishing as $d/n \rightarrow \infty$.

Generalization Error with Spiked Covariance. While recovering spike properties (Sonthalia and Nadakuditi, 2023; Kausik et al., 2024; Nadakuditi, 2014; Benaych-Georges and Nadakuditi, 2011; 2012) and analyzing generalization error in spiked models (Kobak et al., 2020; Ba et al., 2022; 2023; Mousavi-Hosseini et al., 2023; Moniri et al., 2023; Freeman, 2025) are active research areas, existing analyses often characterize generalization implicitly (e.g., via fixed-point equations) or focus on specific spike strengths/alignments. *In contrast, we provide explicit, generic formulae for generalization error, enabling precise categorization of overfitting regimes and conditions for beneficial spike alignment.*

We further discuss connections to selected prior works, direct theoretical extensions to other settings, and proof techniques in Appendix C.

Notation The subscript on $o, O, \omega, \Omega, \Theta$ will denote which quantity is being sent to infinity.

¹Our code is available at the anonymous GitHub repository: link

Table 1: Asymptotic Generalization Regimes. This table summarizes conditions for when overfitting is benign, tempered, or catastrophic in the limit where $d/n \rightarrow c$ and subsequently $c \rightarrow \infty$. The behavior depends on the spike scaling relative to the bulk, target alignment (β_* relative to spike direction \mathbf{u}), and target specifications α_A, α_Z (train) and $\tilde{\alpha}_A, \tilde{\alpha}_Z$ (test). Here, θ^2 quantifies the scaled spike strength and τ^2 the scaled bulk variance; the two primary scaling regimes are operator norm based ($\theta^2 = \gamma\tau^2$) and Frobenius norm based ($\theta^2 = d\tau^2$). The ω, o, O, Θ are all as we send $c \rightarrow \infty$.

Scaling	Benign	Tempered	Catastrophic
Well-Specified, No Covariate Shift: $\alpha_A = \tilde{\alpha}_A = \alpha_Z = \tilde{\alpha}_Z = \alpha > 0$			
$\theta^2 = \gamma\tau^2$	$\gamma = \omega_c(c^2), \beta_* \parallel \mathbf{u}$	All other cases	$o_c(c^2) \geq \gamma \geq \omega_c(1), \beta_* \not\parallel \mathbf{u}$
$\theta^2 = d\tau^2$	$\beta_* \parallel \mathbf{u}$	$\beta_* \not\parallel \mathbf{u}$	Never
Misspecified, No Covariate Shift: $\alpha_A = \tilde{\alpha}_A, \alpha_Z = \tilde{\alpha}_Z, \alpha_A \neq \alpha_Z$			
$\theta^2 = \gamma\tau^2$	Never	All other cases	$o_c(c^2) \geq \gamma \geq \omega_c(1), \beta_* \not\parallel \mathbf{u}$
$\theta^2 = d\tau^2$	Never	Always	Never
Misspecified with Covariate Shift: $\alpha_A \neq \tilde{\alpha}_A$ or $\alpha_Z \neq \tilde{\alpha}_Z$			
$\theta^2 = \gamma\tau^2$	Never	All other cases	$\alpha_Z \neq \tilde{\alpha}_Z, \beta_* \not\parallel \mathbf{u}, \gamma = \omega_c(1)$ or $\alpha_Z = \tilde{\alpha}_Z, \beta_* \not\parallel \mathbf{u}, \omega_c(1) \leq \gamma \leq o_c(c^2)$
$\theta^2 = d\tau^2$	$\alpha_Z = \tilde{\alpha}_Z = \tilde{\alpha}_A,$ $\beta_* \parallel \mathbf{u}$	All other cases	$\alpha_Z \neq \tilde{\alpha}_Z$ and $\beta_* \not\parallel \mathbf{u}$
Spike Recovery: $\alpha_A = \tilde{\alpha}_A = 0, \alpha_Z = \tilde{\alpha}_Z$			(Appendix D)
$\theta^2 = \gamma\tau^2$	$\gamma\tau^2 = o_c(1)$	$\gamma\tau^2 = \Theta_c(1)$	$\gamma\tau^2 = \omega_c(1)$
$\theta^2 = d\tau^2$	$\tau^2 = o_c(1)$	$\tau^2 = \Theta_c(1)$	Never

2 PROBLEM SETTING

We study the generalization of minimum-norm interpolators in high-dimensional linear regression. Using a spiked covariance data model, we quantify how spike strength and alignment influence generalization and the emergence of benign, tempered, or catastrophic overfitting.

Data Model. We consider a data matrix $\mathbf{X} = \mathbf{Z} + \mathbf{A} \in \mathbb{R}^{d \times n}$ with *signal component* \mathbf{Z} and *isotropic noise component* \mathbf{A} that satisfy the following assumptions. Specifically, we shall that the population feature covariance is $\Sigma = \theta^2 \mathbf{u}\mathbf{u}^\top + \tau^2 \mathbf{I}_d$, modeling a rank-one perturbation of isotropic noise.

Assumption 1 (Signal). Let $\mathbf{u} \in \mathbb{R}^d$ be a fixed unit vector representing the spike direction. Then

$$\mathbf{Z} = \theta \mathbf{u}\mathbf{v}^\top, \quad (1)$$

where $\theta > 0$ controls the spike strength, and the vector $\mathbf{v} \in \mathbb{R}^n$ has i.i.d. standard normal entries.

Assumption 2 (Noise). The entries of \mathbf{A} have zero mean and variance τ^2 . The matrix \mathbf{A} satisfies:

- Its entries are uncorrelated and possess finite fourth moments.
- Its distribution is invariant under left and right orthogonal transformations.
- The empirical spectral distribution of $\frac{1}{\tau^2 d} \mathbf{A}\mathbf{A}^\top$ converges to the Marchenko–Pastur law as $n, d \rightarrow \infty$ with $d/n \rightarrow c \in (0, \infty)$.

Spike Strength Normalizations. We consider two key scaling regimes for the spike strength relative to the bulk noise. These lead to distinct generalization behaviors.

Table 2: Conditions for Beneficial Spike Alignment at Finite Aspect Ratios ($c = d/n$). This table outlines the specific regions where alignment of the target signal with the data’s principal spike direction improves generalization. Conditions depend on the problem setting (well-specified vs. mis-specified), the spike scaling regime (operator or frobenius norm based), the overparameterization level $c = d/n$, and the relative dependence of the targets y on the spike versus the bulk α_Z/α_A .

Setting	Alignment Beneficial Region
Well-Specified, Operator Norm	$\gamma > c(c-2)$
Well-Specified, Frobenius Norm	$c > 1$
Misspecified, No Covariate Shift, Operator Norm	$\frac{1}{c} \leq \frac{\alpha_Z}{\alpha_A} \leq \frac{1}{c} \left(\frac{3c^2 - \gamma + 2c\gamma - 2c}{c^2 + \gamma} \right)$
Misspecified, No Covariate Shift, Frobenius Norm	$\frac{1}{c} < \frac{\alpha_Z}{\alpha_A} < 2 - \frac{1}{c}$

- Operator Norm Scaling** ($\theta^2 = \gamma\tau^2$): Here γ tunes the spike strength θ^2 relative to the noise variance τ^2 . When $\gamma = (1 + \sqrt{c})^2$, the spectral norm of the signal component \mathbf{Z} is comparable to that of the noise component \mathbf{A} . If $\gamma > (1 + \sqrt{c})^2$, the spike emerges as an isolated eigenvalue beyond the bulk spectrum established by \mathbf{A} , a phenomenon known as the Baik–Ben Arous–Péché (BBP) transition (Baik et al., 2005). This scaling reflects spikes in learned neural network features (Ba et al., 2022; Moniri et al., 2023).
- Frobenius Norm Scaling** ($\theta^2 = d\tau^2$): Here $\theta^2 = d\tau^2$ matches expected signal and noise Frobenius norms ($\mathbb{E}[\|\mathbf{Z}\|_F^2] = \mathbb{E}[\|\mathbf{A}\|_F^2]$) and the spike has macroscopic proportion of the energy. Such strong signals can lead to improved sample complexity, potentially overcoming limitations observed in purely isotropic models (Ba et al., 2023; Mei et al., 2022).

Target Model. Given $x_i = z_i + a_i$, the targets y are obtained as follows:

$$y_i = \alpha_Z z_i^\top \beta_* + \alpha_A a_i^\top \beta_* + \varepsilon_i, \quad (2)$$

where $\beta_* \in \mathbb{R}^d$ is uniformly distributed in the subspace $\{\beta \in \mathbb{S}^{d-1} : \beta^\top \mathbf{u} = \text{fixed constant}\}$ is the true underlying parameter vector. The terms z_i and a_i are the i -th columns of \mathbf{Z} and \mathbf{A} respectively. The observation noise ε_i are i.i.d. with $\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \tau_\varepsilon^2$. The coefficients $\alpha_Z, \alpha_A \in \mathbb{R}$ control the target’s dependence on the signal and noise components. If $\alpha_Z \neq \alpha_A$, the true data generating process for y differentially weights components of x_i , causing model misspecification.

Generalization Risk. We study the minimum-norm interpolating ordinary least squares estimator:

$$\beta_{int} = \mathbf{X}^\dagger \mathbf{y}, \quad \text{with} \quad \hat{\mathbf{y}} = (\tilde{\mathbf{z}} + \tilde{\mathbf{a}}) \beta_{int} \quad (3)$$

where \mathbf{X}^\dagger denotes the pseudoinverse. Given a new test data point (\tilde{x}, \tilde{y}) , where $\tilde{x} = \tilde{z} + \tilde{a}$ and targets $\tilde{y} = \tilde{\alpha}_Z \tilde{z}^\top \beta_* + \tilde{\alpha}_A \tilde{a}^\top \beta_* + \tilde{\varepsilon}$ with potentially with different coefficients $\tilde{\alpha}_Z, \tilde{\alpha}_A$ and model parameters $\tilde{\tau}, \tilde{\tau}_\varepsilon$, the generalization risk is defined as the expected squared prediction error:

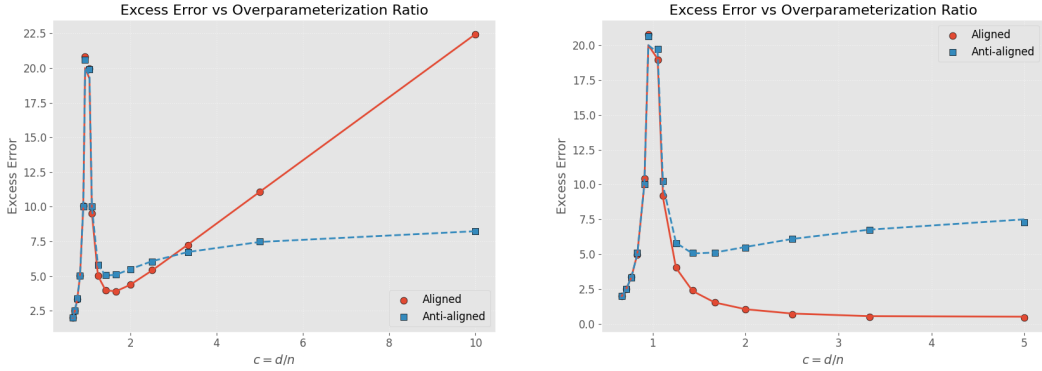
$$\mathcal{R}(\beta_{int}) = \mathbb{E}_{\mathbf{X}, \varepsilon, \{\tilde{x}, \tilde{\varepsilon}\}} [(\tilde{y} - \hat{\mathbf{y}})^2] = \mathbb{E}_{\mathbf{X}, \varepsilon, \tilde{x}, \tilde{\varepsilon}} [(\tilde{y} - \tilde{x}^\top \beta_{int})^2]. \quad (4)$$

The expectation is over the training data $(\mathbf{X}, \varepsilon)$ and the test data realization $(\tilde{x}, \tilde{\varepsilon})$. We shall denote the asymptotic excess risk in the proportional regime as follows:

$$\mathcal{R}_c = \lim_{n, d \rightarrow \infty, d/n \rightarrow c} \mathcal{R}(\beta_{int}) - \tilde{\tau}_\varepsilon^2.$$

Remark 1 (Generalizing Prior Work). *This problem formulation encompasses several existing models as special cases. For instance, isotropic regression settings studied in Hastie et al. (2022) are recovered by setting $\theta = 0$ (no spike) and $\alpha_Z = 0$. Spike recovery models, such as in Sonthalia and Nadakuditi (2023), correspond to specific choices like $\tau^2 = 1/d$, $\tau_\varepsilon^2 = 0$, and $\alpha_A = 0$. Our generalized setup allows for a nuanced investigation of the interplay between signal structure, target alignment, and overparameterization.*

Remark 2. *Although stylized for theoretical simplicity, our setting naturally arises as a low-order Hermite approximation of nonlinear multi-index models, captured by our (α_A, α_Z) -parameterization.*



(a) Operator norm scaling ($\theta^2 = c\tau^2$). Alignment initially improves generalization, but have catastrophic risk as $c \rightarrow \infty$. Anti-alignment yields tempered risk. (b) Equal Frobenius norm scaling ($\theta^2 = d\tau^2$). Alignment leads to benign overfitting, while anti-alignment results in tempered risk.

Figure 1: Excess error vs. overparameterization ratio $c = d/n$ in the well-specified case. Each plot shows the risk for aligned and anti-aligned targets under different spike scaling regimes. **The scatter plots are empirically obtained and the lines are theory.**

This essentially makes our model a tractable surrogate for a broad class of nonlinear targets, while still being simple enough to allow a complete, closed-form risk analysis. Additionally, recent one-step feature-learning analyses like Moniri et al. (2024) show that learned features are **polynomials** of the inputs and the learned spikes. Motivated by this, we study arguably the **simplest non-trivial polynomial target** $y = \alpha_A a + \alpha_Z z$

Quantifying the Benefit of Alignment. A key aspect of our investigation is to determine when the alignment of the true parameter vector β_* with the data’s principal spike direction \mathbf{u} is beneficial for generalization. We define alignment as *beneficial* if the generalization risk $\mathcal{R}(\beta_{int})$ (or \mathcal{R}_c), is monotonically decreasing as a function of $(\beta_*^\top \mathbf{u})^2 \in [0, 1]$. Conversely, alignment is *detrimental* if the risk is a monotonically increasing function of $(\beta_*^\top \mathbf{u})^2$.

Characterizing Overfitting Regimes. Following Bartlett et al. (2020); Mallinar et al. (2022), we classify the asymptotic behavior of the excess risk, \mathcal{R}_c as $c \rightarrow \infty$ as benign, tempered or catastrophic. We say the overfitting is **benign** if $\lim_{c \rightarrow \infty} \mathcal{R}_c$ is zero, **tempered** if this limit is positive and finite, **catastrophic** if this limit is infinite.

3 THEORETICAL RESULTS

Our core theoretical contribution is a precise analytical formula for excess risk in the spiked covariance model. This result relies on Assumption 3, which encompasses both the operator norm scaling ($\theta^2 = \gamma\tau^2$) and Frobenius norm scaling ($\theta^2 = d\tau^2$) regimes. We develop our general risk theorem by analyzing progressively complex scenarios. Specifically, our forthcoming theorems provide specific conditions for benign, tempered, or catastrophic overfitting (as $c \rightarrow \infty$), and determine when, for finite c , alignment of β_* with spike \mathbf{u} is beneficial or detrimental.

Assumption 3 (Scaling). As $n, d \rightarrow \infty$ with $d/n \rightarrow c \in (0, \infty)$, we assume that θ^2 and τ^2 satisfy $\Omega(\tau^2) \leq \theta^2 \leq O(d\tau^2)$ and $\tau^2 = \Theta(1)$.

3.1 WELL SPECIFIED PROBLEM

We begin by analyzing the well-specified case, where the target \mathbf{y} is a direct linear function of the observed covariates $\mathbf{X} = \mathbf{Z} + \mathbf{A}$. This scenario is realized by setting:

$$\alpha_Z = \alpha_A = \tilde{\alpha}_Z = \tilde{\alpha}_A = \alpha > 0.$$

Consequently, $y_i = \alpha \mathbf{x}_i^\top \beta_* + \varepsilon_i$, and the model is properly specified.

Theorem 1 (Well-Specified Risk). *Given data (\mathbf{X}, \mathbf{y}) and $(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})$ generated according to Assumptions 1 (Signal), 2 (Noise), Equation 2 (Target Model), and Assumption 3 (Scaling). If the well-specification condition $\alpha_Z = \alpha_A = \tilde{\alpha}_Z = \tilde{\alpha}_A = \alpha > 0$ holds, the asymptotic excess risk \mathcal{R}_c is:*

$$\mathcal{R}_c = \begin{cases} \tau_\varepsilon^2 \frac{c}{1-c} & \text{if } c < 1 \\ \tau_\varepsilon^2 \frac{1}{c-1} + \alpha^2 \tau^2 \left(1 - \frac{1}{c}\right) \left[\|\beta_*\|^2 + (\beta_*^\top \mathbf{u})^2 \frac{\theta^2 \tau^2 c^2 - 2\theta^2 \tau^2 c - \theta^4}{(\theta^2 + \tau^2 c)^2} \right] & \text{if } c > 1 \end{cases}$$

where \mathbf{u} is the unit vector defining the spike direction.

Remark 3. *If $\theta^2 = \gamma \tau^2$ with $\gamma = o(1)$ (a regime not allowed by Assumption 3 but useful for sanity checks), the coefficient of $(\beta_*^\top \mathbf{u})^2$ vanishes, the risk expression aligns with that of isotropic models, such as in (Hastie et al., 2022, Theorem 1).*

Operator Norm Scaling ($\theta^2 = \gamma \tau^2$). In this regime, the excess risk for $c > 1$ becomes:

$$\mathcal{R}_c = \alpha^2 \tau^2 \left(1 - \frac{1}{c}\right) \left(\|\beta_*\|^2 + \frac{\gamma c^2 - 2\gamma c - \gamma^2}{(\gamma + c)^2} (\beta_*^\top \mathbf{u})^2 \right) + \tau_\varepsilon^2 \frac{1}{c-1}.$$

The formula shows that alignment with the spike direction \mathbf{u} is beneficial if and only if the coefficient of $(\beta_*^\top \mathbf{u})^2$ is negative, which occurs when $\gamma > c(c-2)$. We consider different scalings for γ .

Case 1: $\gamma = \Theta_c(1)$ (constant with respect to c). The condition for beneficial alignment, $\gamma > c(c-2)$, interacts intricately with the BBP phase transition condition, $\gamma > (1 + \sqrt{c})^2$. Let $c_* \approx 4.212$ be the unique solution to $c(c-2) = (1 + \sqrt{c})^2$ for $c > 1$.

- For $1 < c < c_*$: Here, $c(c-2) < (1 + \sqrt{c})^2$. If $c(c-2) < \gamma < (1 + \sqrt{c})^2$, alignment is beneficial even though the BBP transition has *not* occurred (the spike is not resolved from the bulk).
- For $c > c_*$: Here, $c(c-2) > (1 + \sqrt{c})^2$. For alignment to be beneficial ($\gamma > c(c-2)$), the BBP transition must have occurred (as $\gamma > c(c-2) \implies \gamma > (1 + \sqrt{c})^2$). However, the BBP transition occurring is not sufficient for beneficial alignment. If $(1 + \sqrt{c})^2 < \gamma < c(c-2)$, the BBP transition occurs, yet alignment is detrimental.

Regarding the type of overfitting as $c \rightarrow \infty$ (while γ remains constant):

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = \alpha^2 \tau^2 (\|\beta_*\|^2 + \gamma (\beta_*^\top \mathbf{u})^2).$$

Since this limit is a positive constant, we consistently observe *tempered overfitting* when $\gamma = \Theta_c(1)$.

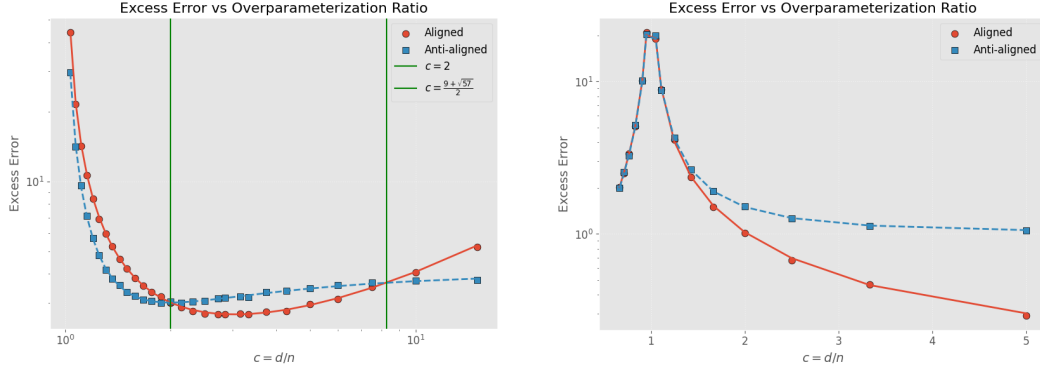
Case 2: $\gamma = \omega_c(1)$ (γ grows with c). The behavior depends on the growth rate of γ relative to c . The limit of the excess risk for $\beta_*^\top \mathbf{u} \neq 0$ as $c \rightarrow \infty$ is:

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = \alpha^2 \tau^2 \cdot \begin{cases} \infty & \text{if } \omega_c(1) \leq \gamma \leq o_c(c^2) \\ \|\beta_*\|^2 + (\frac{1}{\phi} - 1)(\beta_*^\top \mathbf{u})^2 & \text{if } \gamma = \phi c^2 \text{ for const. } \phi > 0 \\ \|\beta_*\|^2 - (\beta_*^\top \mathbf{u})^2 & \text{if } \gamma = \omega_c(c^2) \end{cases}$$

Surprisingly, while $\gamma = \Theta_c(1)$ gives tempered overfitting, increasing spike strength to $\omega_c(1) \leq \gamma \leq o_c(c^2)$ results in *catastrophic overfitting*, even though morally, this version of the problem has less noise. Additionally, we see that this catastrophic overfitting is not present in the anti-aligned $(\beta_*^\top \mathbf{u})$ case. More, aligned with intuition, we see that further increasing the size of the spike improves the generalization performance. Specifically, we get *tempered overfitting* if $\gamma = \phi c^2$ and *benign overfitting* if $\gamma = \omega_c(c^2)$, $\beta_* \parallel \mathbf{u}$ and $\|\beta_*\| = 1$.

For $\gamma = c$, the $(\beta_*^\top \mathbf{u})^2$ coefficient is $(c-3)/4$. Thus, for $1 < c < 3$, alignment is beneficial and for $c > 3$, alignment becomes detrimental. As $c \rightarrow \infty$, if $\beta_* \parallel \mathbf{u}$, the excess risk grows approximately as $\alpha^2 \tau^2 \frac{c}{4} (\beta_*^\top \mathbf{u})^2$, indicating *catastrophic overfitting*. In contrast, if $\beta_* \perp \mathbf{u}$, the excess risk grows like $\alpha^2 \tau^2 (1 - 1/c) \|\beta_*\|^2$, leading to *tempered overfitting*. This transition is illustrated in Figure 1a.

Intuition: This result is a special case of a more general result Theorem 5. At a high level, the complicated phase transitions are a direct consequence of competitions among the four terms seen in Theorem 5. In particular among (1) **Bias**, (2) **Variance**, (3) **Target Alignment**, where (1) and (2) are non-negative (harmful) and (3) is negative (beneficial). For the well specified case, alignment completely cancels the variance, but not the bias, which then drives catastrophic overfitting.



(a) Under operator norm scaling ($\theta^2 = c\tau^2$) with $\alpha_Z = 1$, $\alpha_A = 2$, alignment initially improves generalization for small c , but becomes harmful beyond a critical point, leading to catastrophic overfitting.

(b) Under Frobenius norm scaling ($\theta = \sqrt{d}\tau$) with $\alpha_A = 1$ and $\alpha_Z = 1.1$, alignment remains better than anti-alignment across all c , but benign overfitting is not achieved unless $\alpha_Z = \alpha_A$.

Figure 2: Transition from beneficial to harmful alignment under mild misspecification. The scatter plots are empirically obtained and the lines are theory.

Frobenius Norm Scaling ($\theta^2 = d\tau^2$). The excess risk for $c > 1$ simplifies to:

$$\mathcal{R}_{c>1} = \alpha^2 \tau^2 \left(1 - \frac{1}{c}\right) (\|\beta_*\|^2 - (\beta_*^\top \mathbf{u})^2) + \tau_\varepsilon^2 \frac{1}{c-1}.$$

We have a few observations. First, if $\beta_* \parallel \mathbf{u}$ and $\|\beta_*\| = 1$, the excess risk \mathcal{R}_c tends to 0 as $c \rightarrow \infty$ (*benign overfitting*). Second, if β_* is not perfectly aligned with \mathbf{u} , $\mathcal{R}_c \rightarrow \alpha^2 \tau^2 (\|\beta_*\|^2 - (\beta_*^\top \mathbf{u})^2) > 0$ as $c \rightarrow \infty$ (*tempered overfitting*). Finally, the coefficient of $(\beta_*^\top \mathbf{u})^2$ in the risk formula is negative. Hence, in contrast with the operator norm regime, *alignment is always beneficial* in this regime for $c > 1$, and we visualize these behaviors in Figure 1b.

Takeaways for the Well-Specified Case. As a result, spike scaling profoundly impacts overfitting, especially with target alignment, which is not always beneficial. For aligned targets, increasing spike strength can drive transitions from tempered \rightarrow catastrophic \rightarrow tempered \rightarrow benign overfitting, while anti-alignment ($\beta_* \perp \mathbf{u}$) can mitigate catastrophic overfitting.

3.2 MISSPECIFIED CASE AND NO COVARIATE SHIFT

We next consider misspecified targets \mathbf{y} with differing dependence on spike \mathbf{Z} and noise \mathbf{A} feature components. Specifically, we assume $\alpha_Z \neq \alpha_A$ but introduce no covariate shift between training and test distributions, i.e., $\tilde{\alpha}_Z = \alpha_Z$ and $\tilde{\alpha}_A = \alpha_A$. This scenario models situations where intrinsic feature properties lead to differential correlations with the target, a common occurrence in practice. For notational convenience, we define $\Delta_c := \alpha_Z - \frac{\alpha_A}{c}$ with $\Delta_1 := \alpha_Z - \alpha_A$.

Theorem 2 (Misspecified). *Let $\mathbf{Z}, \tilde{\mathbf{Z}}$ satisfy Assumption 1, $\mathbf{A}, \tilde{\mathbf{A}}$ satisfy Assumption 2 and $\mathbf{y}, \tilde{\mathbf{y}}$ according to Equation (2). If Assumption 3 holds with $\alpha_Z = \tilde{\alpha}_Z$, $\alpha_A = \tilde{\alpha}_A$, then*

$$\mathcal{R}_c = \begin{cases} \tau_\varepsilon^2 \frac{c}{1-c} + \tau^2 (\beta_*^\top \mathbf{u})^2 \frac{\Delta_1^2}{1-c} \frac{\theta^2}{\theta^2 + \tau^2} & c < 1 \\ \tau_\varepsilon^2 \frac{1}{c-1} + \alpha_A^2 \tau^2 \|\beta_*\|^2 \left(1 - \frac{1}{c}\right) + \tau^2 (\beta_*^\top \mathbf{u})^2 \Delta_c^2 \frac{\theta^2}{\theta^2 + \tau^2 c} \left[\frac{c}{c-1} \frac{\theta^2 + \tau^2 c^2}{\theta^2 + \tau^2 c} - 2 \frac{\alpha_A}{\Delta_c}\right] & c > 1 \end{cases}$$

A key observation is that misspecification ($\alpha_Z \neq \alpha_A$) can itself induce double descent, even if $\tau_\varepsilon^2 = 0$. This contrasts with the well-specified case where, if $\tau_\varepsilon^2 = 0$, double descent is absent. However, in the misspecified case, we do not observe double descent if there is no alignment $\beta_*^\top \mathbf{u} = 0$.

Equal Operator Norm Case. For $\theta^2 = \gamma\tau^2$, the excess risk is

$$\mathcal{R} = \begin{cases} \tau^2(\beta_*^\top \mathbf{u})^2 \frac{\Delta_c^2}{1-c} \frac{\gamma}{\gamma+1} + \tau_\varepsilon^2 \frac{c}{1-c} & c < 1 \\ \tau^2 \frac{\gamma}{\gamma+c} (\beta_*^\top \mathbf{u})^2 \Delta_c^2 \left[\left(\frac{c^2+\gamma}{\gamma+c} \frac{c}{c-1} \right) - 2\frac{\alpha_A}{\Delta_c} \right] + \alpha_A^2 \tau^2 \|\beta_*\|^2 \left(1 - \frac{1}{c}\right) + \tau_\varepsilon^2 \frac{1}{c-1} & c > 1 \end{cases}$$

For $c < 1$, the spike is *detrimental*. For $c > 1$, the behavior depends on α_Z/α_A . In particular, if

$$\frac{1}{c} \leq \frac{\alpha_Z}{\alpha_A} \leq \frac{1}{c} \left(\frac{3c^2 - \gamma + 2c\gamma - 2c}{c^2 + \gamma} \right),$$

then we have that the coefficient in front of $(\beta_*^\top \mathbf{u})^2$ is negative. Thus, when α_Z/α_A lies between these thresholds, the spike *helps*, but the spike is *harmful* outside this range. As $c \rightarrow \infty$, if $\gamma = o_c(c^2)$, the beneficial region shrinks and *alignment increasingly harms generalization*. On the other hand, if the spike is big enough ($\gamma = \omega_c(c^2)$), we have that the beneficial region limits to $0 \leq \frac{\alpha_Z}{\alpha_A} \leq 2$. Figures 3a and 3b plot the coefficient of $(\beta_*^\top \mathbf{u})^2$ for $c = 2$ and $c = 20$ for $\gamma = c$.

The upper bound on beneficial α_Z/α_A is surprising, as stronger target dependence on the spike might be expected to always favor alignment. Additionally, the dependence on the level of overparameterization c also offers new insights. Consider the example of $\gamma = c$, and $\alpha_Z/\alpha_A = 2$. Then when $c < 2$ or $c > (9 + \sqrt{57})/2$, we have that the ratio is outside the beneficial region. Figure 2a shows that in the beneficial region, the aligned risk is lower than the anti-aligned risk. However, outside the beneficial region, the aligned risk becomes strictly larger than the anti-aligned counterpart.

Next, in terms of benign vs. tempered vs. catastrophic overfitting, we have that

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = \begin{cases} \tau^2 [\gamma \alpha_Z^2 (\beta_*^\top \mathbf{u})^2 + \alpha_A^2 \|\beta_*\|^2] & \beta_* \not\perp \mathbf{u}, \gamma = \Theta_c(1) \\ \infty & \beta_* \not\perp \mathbf{u}, \omega_c(1) \leq \gamma \leq o_c(c^2) \\ \tau^2 \left[\alpha_A^2 \|\beta_*\|^2 + \left(\alpha_Z^2 \left(1 + \frac{1}{\phi}\right) - 2\alpha_Z \alpha_A \right) (\beta_*^\top \mathbf{u})^2 \right] & \beta_* \not\perp \mathbf{u}, \gamma = \phi c^2 \\ \tau^2 (\alpha_A^2 \|\beta_*\|^2 + (\alpha_Z^2 - 2\alpha_Z \alpha_A) (\beta_*^\top \mathbf{u})^2) & \beta_* \not\perp \mathbf{u}, \gamma = \omega_c(c^2) \\ \alpha_A^2 \tau^2 \|\beta_*\|^2 & \beta_* \perp \mathbf{u} \end{cases}.$$

For $\beta_* \not\perp \mathbf{u}$, if $\omega_c(1) \leq \gamma \leq o_c(c^2)$ we have *catastrophic overfitting*. If $\gamma = \Theta_c(c^2)$, overfitting is tempered, with benign overfitting precluded (Appendix Proposition 3). If $\gamma = \omega_c(c^2)$, overfitting is again tempered with benign requiring returning to the well-specified case ($\alpha_A = \alpha_Z$).

Ridge Extension. As done in Li and Sonthalia (2024), it is also possible to analyze the ridge regularized version by looking at the resolvent away from zero. Since the general proof structure remains the same, we leave this as a future direction and empirically test this case. In particular, we follow the setting of Figure 1a ($\theta^2 = c\tau^2$) and vary the regularization strengths ($\lambda \in \{0, 1, c, c^2, dc\}$). In all cases the catastrophic band in our phase diagram persists, indicating that generally regularization does not remove alignment-driven catastrophic overfitting. The plot can be seen in Appendix B.2.

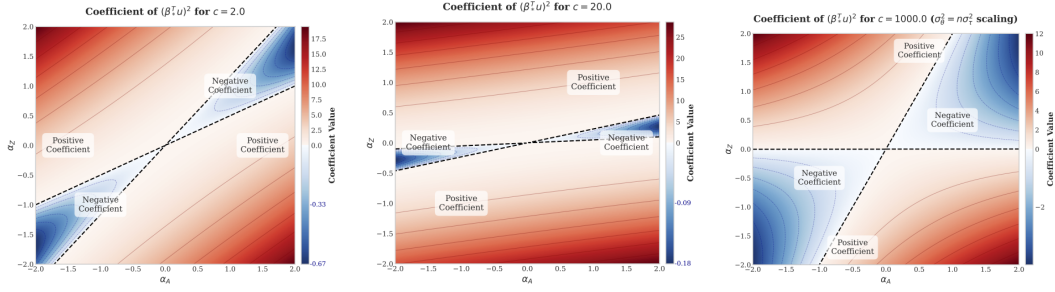
Equal Frobenius Norm Case. For $\theta^2 = d\tau^2$, the excess risk becomes:

$$\mathcal{R}_{c>1} = \alpha_A^2 \|\beta_*\|^2 \left(1 - \frac{1}{c}\right) + (\beta_*^\top \mathbf{u})^2 \left[\frac{c}{c-1} \left(\alpha_Z - \frac{\alpha_A}{c}\right)^2 - 2\alpha_A \left(\alpha_Z - \frac{\alpha_A}{c}\right) \right] + \frac{\tau_\varepsilon^2}{c-1}.$$

For $c > 1$, the beneficial region for the ratio α_Z/α_A is defined by: $\frac{1}{c} \leq \frac{\alpha_Z}{\alpha_A} \leq 2 - \frac{1}{c}$. The beneficial region expands with c , making alignment increasingly beneficial in extreme overparameterization (Figure 3c). Beneficial alignment can also be seen in Figure 2b. Here $\alpha_Z/\alpha_A = 1.1$, which is in the beneficial region for $c > 10/9$. Finally, the overfitting is tempered unless $\alpha_A = \alpha_Z$.

3.3 MISSPECIFIED TARGET AND COVARIATE SHIFT

Lastly, in addition to misspecification, we also have covariate shift between train and test. Specifically, $\alpha_Z \neq \tilde{\alpha}_Z$ or $\alpha_A \neq \tilde{\alpha}_A$, hence we have the spike/noise importance differ between train and test. For the **equal operator norm** case, we show the following.



(a) Operator norm scaling, $c = 2$. Large beneficial region. (b) Operator norm scaling, $c = 20$. Smaller beneficial region (c) Frobenius norm scaling, $c = 1000$. The beneficial region persists at extreme overparameterization.

Figure 3: Phase boundaries for spike alignment impact. Coefficient of $(\beta_*^\top \mathbf{u})^2$ as a function of α_Z/α_A , indicating whether alignment improves or harms generalization.

Theorem 3. Given data $\mathbf{Z}, \tilde{\mathbf{Z}}$ that satisfy Assumption 1, $\mathbf{A}, \tilde{\mathbf{A}}$ that satisfy Assumption 2 and $\mathbf{y}, \tilde{\mathbf{y}}$ according to Equation (2). If Assumption 3 holds, catastrophic overfitting occurs if $\tilde{\alpha}_Z = \alpha_Z$, $\beta_* \not\perp \mathbf{u}$, and $\omega_c(1) \leq \gamma \leq o_c(c^2)$. Additionally, if $\tilde{\alpha}_Z \neq \alpha_Z$ with $\gamma = \omega_c(1)$ and $\beta_* \not\perp \mathbf{u}$ we get catastrophic overfitting. Other scenarios yield tempered overfitting.

Different covariate shifts pose varying challenges. In particular, if $\alpha_Z \neq \tilde{\alpha}_Z$, (target’s spike dependence shifts), then catastrophic overfitting becomes unavoidable for sufficiently large spikes. This contradicts the earlier theoretical intuition, as increasing the spike size in this setting actually induces catastrophic overfitting instead of mitigating it.

Equal Frobenius Norm. In this case, we have the following theorem.

Theorem 4. Let $\mathbf{Z}, \tilde{\mathbf{Z}}$ satisfy Assumption 1, $\mathbf{A}, \tilde{\mathbf{A}}$ satisfy Assumption 2 and $\mathbf{y}, \tilde{\mathbf{y}}$ according to Equation (2). If Assumption 3 holds and $\alpha_Z \neq \tilde{\alpha}_Z$ then $\mathcal{R}_c = \infty$ for all $c \neq 1$. For $\alpha_Z = \tilde{\alpha}_Z$:

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = \tau^2 [(\beta_*^\top \mathbf{u})^2 (\alpha_Z^2 - 2\tilde{\alpha}_A \alpha_Z) + \|\beta_*\|^2 \tilde{\alpha}_A^2].$$

If $\alpha_Z \neq \tilde{\alpha}_Z$, catastrophic overfitting occurs. When β_* and \mathbf{u} are parallel, we have that $\tau^2 \|\beta_*\|^2 (\alpha_Z - \tilde{\alpha}_A)^2$. This is benign if and only if $\alpha_Z = \tilde{\alpha}_A$. Notably, if training data is misspecified ($\alpha_A \neq \alpha_Z$) but test data is well-specified and matches the training spike dependence ($\alpha_Z = \tilde{\alpha}_Z = \tilde{\alpha}_A$), benign overfitting becomes achievable.

3.4 GENERAL THEOREM

Prior results are special cases of our main theorem (Theorem 5). Its full form is complex (Appendix E). We present a high-level decomposition here.

Theorem 5 (Generalization Risk). Suppose Assumption 1, Assumption 2, and Assumption 3 hold.

$$\mathcal{R} = \mathbb{E} \left[\underbrace{\left\| \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} \right\|_F^2}_{\text{Bias}} + \underbrace{\tau^2 \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Variance}} + \underbrace{\tilde{\alpha}_A^2 \left\| \beta_*^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Data Noise}} + \underbrace{\left(-2\tilde{\alpha}_A \beta_*^\top \tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top \beta_{int} \right)}_{\text{Target Alignment}} \right].$$

- **Bias.** This is the squared error between the learned predictor β_{int} and the true parameter β_* projected onto the spike direction \mathbf{u} . In particular, the risk penalizes discrepancies only along the top eigen-direction of the population covariance Σ , reflecting the anisotropic influence of the spike.
- **Variance.** The variance is equivalent to $\tau^2 \|\beta_{int}\|_2$. This mirrors classical isotropic regression results (Hastie et al., 2022; Bartlett et al., 2020), but the norm $\|\beta_{int}\|^2$ itself is dependent upon the interaction between signal and noise, the alignment between β_* and \mathbf{u} , and the scaling parameters.
- **Data Noise.** The data noise term quantifies the contribution of the noise matrix \mathbf{A} to the target outputs y_i through α_A . Even in the absence of observation noise ($\tau_\varepsilon^2 = 0$), target corruption via data noise can create an irreducible error floor.

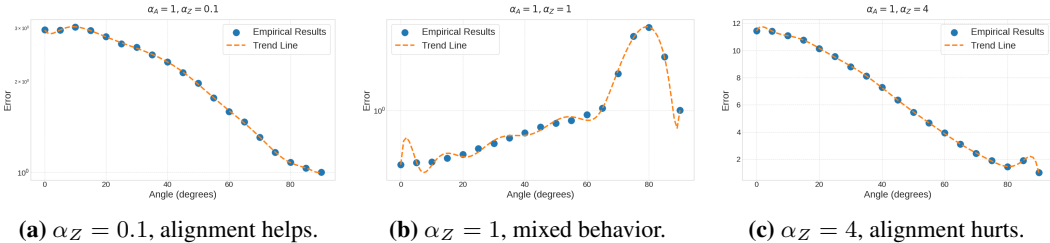


Figure 4: Alignment-phase transitions persist in deep networks. Generalization error vs. angle between spike direction \mathbf{u} and ground-truth parameter β_* when fitting data with a 3-layer ReLU networks. The effect of alignment switches as α_Z increases, consistent with the phase transitions predicted by our theory. Experimental details are in Appendix B.

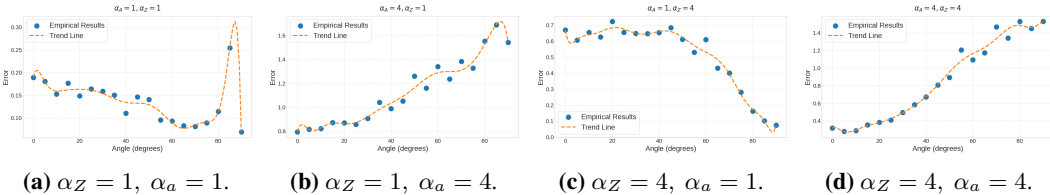


Figure 5: Generalization error vs. alignment for deep networks. Generalization error vs. the angle between spike direction \mathbf{u} and ground-truth parameter β_* when fitting MNIST-derived data with a ReLU network, for a sweep over $(\alpha_Z, \alpha_a) \in \{1, 4\}^2$.

- **Target Alignment.** The alignment term measures the inner product between β_{int} and β_* with respect to the sample noise covariance. This cross-term captures how mismatch between β_{int} and β_* , especially when mediated by \mathbf{A} , can amplify or dampen generalization error.

3.5 NONLINEAR MODELS

While our theoretical focus is on linear regression, key phenomena like α_Z dependent non-monotonic alignment effects appear in nonlinear models as well. We test this by training 3-layer ReLU networks to predict \mathbf{y} (Equation (2)) given \mathbf{X} , where we vary the alignment angle between spike \mathbf{u} and β_* and record the generalization error. Figure 4, shows our results for three α_Z values. For $\alpha_Z = 0.1$, increasing alignment with the spike is detrimental. For $\alpha_Z = 1$, alignment is beneficial, while for $\alpha_Z = 10$, alignment is detrimental again. This mirrors our theoretical findings that there is a region for beneficial alignment and a nuanced phase transition for different α_Z values.

To move beyond purely synthetic inputs, we perform an experiment on MNIST where we artificially inject a spiked direction into the input space. Please see Appendix B.1 for more details. We then train model as we vary $(\alpha_A, \alpha_Z) \in \{1, 4\}^2$ and the target–spike angle. We again observe a similar pattern: the phase transitions we recover potentially generalize to nonlinear deep networks trained on real data, broadening the theoretical robustness of our setting. Results can be seen in Figure 5.

4 CONCLUSION

This work provided a precise analytical characterization of the generalization error for minimum-norm interpolators in spiked covariance models. We decomposed the risk into interpretable components and comprehensively classified overfitting regimes based on spike strength, target alignment, and overparameterization. We reveal surprising phenomena, such as the potential for increasing spike strength to induce catastrophic overfitting before benign overfitting in well-specified aligned problems, and that strong target-spike alignment is not universally beneficial, especially under model misspecification. These alignment-dependent phase transitions, theoretically derived for linear models, were also empirically observed in nonlinear neural networks, suggesting broader relevance. Our results offer a more nuanced understanding of generalization in the presence of data anisotropy, challenging conventional intuitions and providing a detailed map of risk behaviors in overparameterized settings.

REFERENCES

- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-dimensional asymptotics of feature learning: How one gradient step improves the representation. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=akddwRG6EGi>.
- Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, and Denny Wu. Learning in the presence of low-dimensional structure: A spiked random matrix perspective. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HlIAoCHDWW>.
- Zhidong Bai and Wang Zhou. Large sample covariance matrices without independence structures in columns. 2008.
- Jinho Baik and Jack W Silverstein. Eigenvalues of large sample covariance matrices of spiked population models. *Journal of multivariate analysis*, 97(6):1382–1408, 2006.
- Jinho Baik, Gérard Ben Arous, and Sandrine Péché. Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. 2005.
- Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.
- Florent Benaych-Georges and Raj Rao Nadakuditi. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, 2011.
- Florent Benaych-Georges and Raj Rao Nadakuditi. The singular values and vectors of low rank perturbations of large rectangular random matrices. *Journal of Multivariate Analysis*, 111:120–135, 2012.
- Yuan Cao, Quanquan Gu, and Mikhail Belkin. Risk bounds for over-parameterized maximum margin classification on sub-gaussian mixtures. *Advances in Neural Information Processing Systems*, 34: 8407–8418, 2021.
- Niladri S. Chatterji and Philip M. Long. Finite-sample analysis of interpolating linear classifiers in the overparameterized regime. *Journal of Machine Learning Research*, 22(129):1–30, 2021. URL <http://jmlr.org/papers/v22/20-974.html>.
- Romain Couillet and Zhenyu Liao. *Random Matrix Methods for Machine Learning*. Cambridge University Press, 2022. doi: 10.1017/9781009128490. <https://zhenyu-liao.github.io/book/>.
- Alex Damian, Jason D. Lee, and Mahdi Soltanolkotabi. Neural networks can learn representations with gradient descent, 2022. URL <https://arxiv.org/abs/2206.15144>.
- Yatin Dandi, Florent Krzakala, Bruno Loureiro, Luca Pesce, and Ludovic Stephan. How two-layer neural networks learn, one (giant) step at a time. *Journal of Machine Learning Research*, 25(349): 1–65, 2024. URL <http://jmlr.org/papers/v25/23-1543.html>.
- Samet Demir and Zafer Dogan. Random features outperform linear models: Effect of strong input-label correlation in spiked covariance data. *arXiv preprint arXiv:2409.20250*, 2024.
- Samet Demir and Zafer Dogan. Asymptotic analysis of two-layer neural networks after one gradient step under gaussian mixtures data with structure. *arXiv preprint arXiv:2503.00856*, 2025.
- Edgar Dobriban and Stefan Wager. High-dimensional asymptotics of prediction: Ridge regression and classification. *The Annals of Statistics*, 46(1):247–279, 2018.

- Jake Freeman. Shrinkage to infinity: Reducing test error by inflating the minimum norm interpolator in linear models. *arXiv preprint arXiv:2510.19206*, 2025.
- Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-dimensional ridgeless least squares interpolation. *Annals of statistics*, 50(2):949, 2022.
- Kedar Karhadkar, Erin George, Michael Murray, Guido F Montufar, and Deanna Needell. Benign overfitting in leaky relu networks with moderate input dimension. *Advances in Neural Information Processing Systems*, 37:36634–36682, 2024.
- Chinmaya Kausik, Kashvi Srivastava, and Rishi Sonthalia. Double descent and overfitting under noisy inputs and distribution shift for linear denoisers. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=HxfqTdLIRF>.
- Dmitry Kobak, Jonathan Lomond, and Benoit Sanchez. The optimal ridge penalty for real-world high-dimensional data can be zero or negative due to the implicit ridge regularization. *Journal of Machine Learning Research*, 21(169):1–16, 2020. URL <http://jmlr.org/papers/v21/19-844.html>.
- Frederic Koehler, Lijia Zhou, Danica J Sutherland, and Nathan Srebro. Uniform convergence of interpolators: Gaussian width, norm bounds and benign overfitting. *Advances in Neural Information Processing Systems*, 34:20657–20668, 2021.
- Xinyue Li and Rishi Sonthalia. Least squares regression can exhibit under-parameterized double descent. *Advances in Neural Information Processing Systems*, 2024.
- Tengyuan Liang and Alexander Rakhlin. Just interpolate. *The Annals of Statistics*, 48(3):1329–1347, 2020.
- Neil Mallinar, James B Simon, Amirhesam Abedsoltan, Parthe Pandit, Mikhail Belkin, and Preetum Nakkiran. Benign, tempered, or catastrophic: A taxonomy of overfitting. *arXiv preprint arXiv:2207.06569*, 2022.
- V A Marchenko and Leonid A. Pastur. Distribution of eigenvalues for some sets of random matrices. *Mathematics of The Ussr-sbornik*, 1:457–483, 1967.
- Charles H. Martin and Michael W. Mahoney. Implicit self-regularization in deep neural networks: Evidence from random matrix theory and implications for learning. *Journal of Machine Learning Research*, 22(165):1–73, 2021. URL <http://jmlr.org/papers/v22/20-410.html>.
- Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Generalization error of random feature and kernel methods: hypercontractivity and kernel matrix concentration. *Applied and Computational Harmonic Analysis*, 59:3–84, 2022.
- Carl D. Meyer, Jr. Generalized Inversion of Modified Matrices. *SIAM Journal on Applied Mathematics*, 1973.
- Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks. *arXiv preprint arXiv:2310.07891*, 2023.
- Behrad Moniri, Donghwan Lee, Hamed Hassani, and Edgar Dobriban. A theory of non-linear feature learning with one gradient step in two-layer neural networks, 2024. URL <https://arxiv.org/abs/2310.07891>.
- Alireza Mousavi-Hosseini, Denny Wu, Taiji Suzuki, and Murat A Erdogdu. Gradient-based feature learning under structured data. *Advances in Neural Information Processing Systems*, 36:71449–71485, 2023.
- Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpolation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1): 67–83, 2020.

- Raj Rao Nadakuditi. Optshrink: An algorithm for improved low-rank signal matrix denoising by optimal, data-driven singular value shrinkage. *IEEE Transactions on Information Theory*, 60(5): 3002–3018, 2014.
- Ohad Shamir. The implicit bias of benign overfitting. In Po-Ling Loh and Maxim Raginsky, editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings of Machine Learning Research*, pages 448–478. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/shamir22a.html>.
- Rishi Sonthalia and Raj Rao Nadakuditi. Training data size induced double descent for denoising feed-forward neural networks and the role of training noise. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=FdMWtpVT1I>.
- Rishi Sonthalia, Michael Murray, and Guido’ Montúfar. Low rank gradients and where to find them. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=de3kwOXQ9e>.
- Alexander Tsigler and Peter L Bartlett. Benign overfitting in ridge regression. *Journal of Machine Learning Research*, 24(123):1–76, 2023.
- Yutong Wang, Rishi Sonthalia, and Wei Hu. Near-interpolators: Rapid norm growth and the trade-off between interpolation and generalization. In *International Conference on Artificial Intelligence and Statistics*, pages 4483–4491. PMLR, 2024a.
- Zhichao Wang, Denny Wu, and Zhou Fan. Nonlinear spiked covariance matrices and signal propagation in deep neural networks. *arXiv preprint arXiv:2402.10127*, 2024b.
- Yimin Wei. The weighted moore-penrose inverse of modified matrices. *Appl. Math. Comput.*, 122(1):1–13, June 2001. ISSN 0096-3003. doi: 10.1016/S0096-3003(00)00007-2. URL [https://doi.org/10.1016/S0096-3003\(00\)00007-2](https://doi.org/10.1016/S0096-3003(00)00007-2).
- Denny Wu and Ji Xu. On the optimal weighted ℓ_2 regularization in overparameterized linear regression. *Advances in Neural Information Processing Systems*, 33:10112–10123, 2020.

CONTENTS

1	Introduction	1
2	Problem Setting	3
3	Theoretical Results	5
3.1	Well Specified Problem	5
3.2	Misspecified Case and no Covariate Shift	7
3.3	Misspecified Target and Covariate Shift	8
3.4	General Theorem	9
3.5	Nonlinear Models	10
4	Conclusion	10
A	Notation	16
B	More Non-Linear Experiments and Settings	16
B.1	MNIST Data	16
B.2	Ridge Figure	16
C	Prior Works & Natural Extensions	17
C.1	New Techniques & Possible Theoretical Extensions	17
D	Spike Recovery Case	18
E	Proof of Theorem 5	19
E.1	Step 0: Rescaling	20
E.2	Step 1: Decompose Error	21
E.3	Step 2: Simplifying Terms	21
E.3.1	Bias	21
E.3.2	Variance	21
E.3.3	Data Noise	22
E.3.4	Target Alignment	22
E.3.5	Helper Lemmas	22
E.4	Step 3: Random Matrix Theory Estimates	27
E.4.1	Step 3(a): Showing that basic building blocks concentrate	28
E.4.2	Step 3(b): Bounding the Higher Moments	33
E.4.3	Step 3(c): Bounding γ_i moments.	36
E.5	Step 4: Bounding the Expectation of Products of Dependent Terms	40
E.5.1	Step 4: Bias	40
E.5.2	Step 4: Variance	41

E.5.3	Step 4: Data Noise	42
E.5.4	Step 4: Target Alignment	42
E.5.5	Bias: Helper Lemmas	42
E.5.6	Variance: Helper Lemmas	46
E.5.7	Target Alignment: Helper Lemmas	54
E.6	Step 5: Upscaling and Asymptotic Risk Formulas	58
F	Probability Lemmas	60
G	Proof of Specific Cases and Overfitting	70
G.1	Proof of Theorem 1.	70
G.2	Proof of Theorem 2.	70
G.3	Proof of Theorem 3.	70
G.4	Proof of Theorem 4.	71

A NOTATION

Symbol	Description / Role	Typical scaling / range	First used
d, n	Data dimension and sample size	$d, n \rightarrow \infty$ with $c = d/n$ fixed	Sec. 2
c	Aspect ratio d/n	$(0, \infty)$	Sec. 2
τ^2/d	Noise variance in ambient bulk A	$\tau^2 = \Theta(1)$	Sec. 2
θ^2	Spike (signal) variance	$\theta^2 = \gamma\tau^2$ (operator-norm) or $\theta^2 = d\tau^2$ (Frobenius)	Sec. 2
γ	Spike-to-noise ratio $\gamma = \theta^2/\tau^2$ (effective outlier eigenvalue)	$[0, \infty)$; critical line $\gamma = (1 + \sqrt{c})^2$	Sec. 2
α_Z, α_A	Coeffs. weighting spike vs. bulk in <i>targets</i> y	$\Theta(1)$	Eq. (2)
$\tilde{\alpha}_Z, \tilde{\alpha}_A$	Same coefficients for <i>test</i> data (covariate shift)	$\Theta(1)$	Sec. 3
β_*	True parameter vector	$\ \beta_*\ _2 = 1$	Sec. 2
\mathbf{u}	Spike direction in data covariance	$\ \mathbf{u}\ _2 = 1$	Sec. 2
\mathbf{A}, \mathbf{Z}	Bulk noise matrix, rank-one signal matrix	$A_{ij} \sim \mathcal{N}(0, \tau^2/d)$, $\mathbf{Z} = \theta \mathbf{u}\mathbf{v}^\top$	Sec. 2
$\varepsilon, \tau_\varepsilon^2$	Label noise and its variance	IID, $\mathcal{N}(0, \tau_\varepsilon^2)$	Sec. 2

Table 3: Glossary of recurrent parameters and symbols. All $\Theta(1)$ constants are independent of n, d .

Other Notations. We use lowercase a , lowercase bold \mathbf{a} , and uppercase bold \mathbf{A} letters to denote scalars, vectors, and matrices respectively. We use $\|\cdot\|_2$ to denote the Euclidean norm if the argument is a vector and the operator norm if the argument is a matrix. We use $\|\cdot\|_F$ to denote the Frobenius norm. When slicing one entry from a vector or matrix, we use both a_i, A_{ij} and $\mathbf{a}_i, \mathbf{A}_{ij}$, where the latter intends to emphasize the source of the scalar.

B MORE NON-LINEAR EXPERIMENTS AND SETTINGS

We used 500 data points in 750 dimensional space, with a hidden width of 1000. We used full batch gradient descent for 100 epochs with a learning rate of 1e-4. Each data point is averaged over 50 trials. Equal Frobenius norm scaling was used for the size of the spike.

B.1 MNIST DATA

For each trial, we construct a spiked design matrix by combining a rank-one signal with a ‘‘bulk’’ component drawn from MNIST. Fix sample size n and ambient dimension d (here $d = 784$). First, we sample a unit-norm ground-truth vector $\beta_* \in \mathbb{R}^d$ and, for each prescribed angle, construct a unit vector $\mathbf{u} \in \mathbb{R}^d$ at that angle to β_* . Independently sample unit vectors $\mathbf{v}, \mathbf{v}_t \in \mathbb{R}^n$ and set spike strengths $\theta = \theta_t = \sqrt{n} \tau$. The rank-one signal matrices are

$$\mathbf{Z} = \theta \mathbf{u}\mathbf{v}^\top \in \mathbb{R}^{d \times n}, \quad \mathbf{Z}_t = \theta_t \mathbf{u}\mathbf{v}_t^\top \in \mathbb{R}^{d \times n}.$$

The bulk matrices are generated by sampling n MNIST images (train split for training, test split for testing), vectorizing each image into \mathbb{R}^{784} , and stacking these as columns to obtain

$$\mathbf{A} \in \mathbb{R}^{d \times n} \quad \text{and} \quad \mathbf{A}_t \in \mathbb{R}^{d \times n},$$

with raw pixel intensities in $[0, 1]$. The training features are $\mathbf{X} = \mathbf{Z} + \mathbf{A}$, and responses are generated as a noisy linear measurement of a weighted combination of spike and bulk:

$$\mathbf{y} = \beta_*^\top (\alpha_Z \mathbf{Z} + \alpha_A \mathbf{A}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \tau_\varepsilon^2 \mathbf{I}_n).$$

Network architecture. We use the same network structure as previously, but trained for 1000 epochs.

B.2 RIDGE FIGURE

Figure 6 shows the risk that still manifests alignment-driven catastrophic overfitting with ridge regularization.

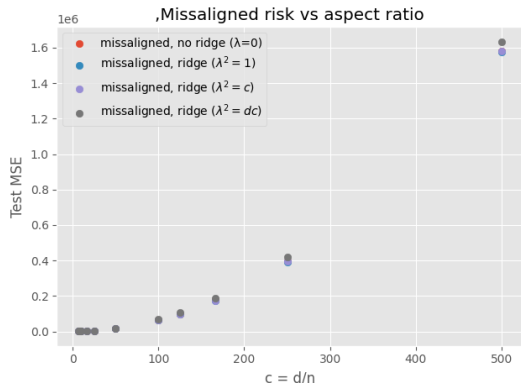


Figure 6: Catastrophic Overfitting even in the presence of Ridge Regularizer.

C PRIOR WORKS & NATURAL EXTENSIONS

Li and Sonthalia (2024). This paper mainly shows that least-squares regression can exhibit a *double-descent peak in the under-parameterized regime*, by carefully analyzing the ridge-regularized resolvent bounded away from zero. They provide two examples where these spectral/alignment properties move the peak to $c < 1$ and argue that existing double-descent explanations, which focus on $c \geq 1$, are incomplete. Our setting recovers theirs when the regularization $\mu = 0$, $\alpha_Z = \alpha_A = 1$ (equal dependence), and learning a denoiser (target is the noiseless data \mathbf{X}). Our model significantly generalizes beyond the data model with interpolating parameters α_A, α_Z , uncovering different insights about overfitting behaviors.

Demir and Dogan (2024) and Demir and Dogan (2025). Demir and Dogan (2025) analyzes the Gaussian-mixture/spiked data after one gradient step. In contrast, we provide closed-form generalization risk for the minimum-norm linear interpolator. In particular, this paper serves as the motivation for studying their setting, which is the analytically tractable “base case” that isolates spike–alignment mechanisms. On the other hand, Demir and Dogan (2024) considers whether substituting one activation function for another results in the same generalization error; it does not characterize the error itself, which is the goal of this paper.

In terms of proof techniques, our paper requires more sophisticated bounding approaches, as we argue in the subsequent sections.

C.1 NEW TECHNIQUES & POSSIBLE THEORETICAL EXTENSIONS

Our analysis combines ingredients that, to the best of our knowledge, have not been put together in prior work on benign overfitting or spiked regression.

In particular, we develop a *mixed spherical hypercontractivity* argument that controls products of random terms involving both bulk and spiked components. Instead of assuming strict Gaussianity, we only require zero mean, unit variance, and a mild rotational invariance; hypercontractivity up to order $2k$ then yields concentration for multilinear forms of the type that appear in our risk formula.

Two natural generalizations from our data model include (1) adding ridge regularization and (2) incorporating general low-rank data.

Ridge Regression. As discussed above for Li and Sonthalia (2024), how regularization affects the data resolvent can be tracked and a similar Stieltjes transform argument applies. In particular, check their Lemmas 13-15 for modified statements that can be generalized.

General Low-Rank Data. Our proof requires a Sherman-Morrison-style expansion for pseudo-inverse in Meyer (1973). This result from perturbation theory allows us to decouple spike and bulk and makes the subsequent analysis more tractable. A Woodbury-style formula for pseudo-inverse (extending Meyer (1973)) exists in Wei (2001), and the same proof structure follows. We can then bound contributions along each singular direction, which requires significantly more bookkeeping.

Due to the significant effort involved rigorously doing these calculations. We leave this as future work. However, we believe that an extension to ridge with multiple spikes is feasible.

D SPIKE RECOVERY CASE

We consider the special case where the goal is to recover the spike direction \mathbf{u} . In this setting, the target \mathbf{y} depends only on the spike component \mathbf{Z} , with no contribution from the noise \mathbf{A} :

$$\alpha_A = \tilde{\alpha}_A = 0, \quad \alpha_Z = \tilde{\alpha}_Z = \alpha > 0.$$

Thus, the target \mathbf{y} is proportional to the signal \mathbf{Z} plus possible observation noise ε .

Equal Operator Norm In this regime, we have that the risk is

$$\mathcal{R}_{c<1} = \frac{\gamma\alpha_Z^2\tau^2}{(1-c)(\gamma+1)}(\boldsymbol{\beta}^\top \mathbf{u})^2 + \frac{c}{1-c}\tau_\varepsilon^2, \quad \mathcal{R}_{c>1} = \frac{\gamma c(c^2 + \gamma)\alpha_Z^2\tau^2}{(c-1)(\gamma+c)^2}(\boldsymbol{\beta}^\top \mathbf{u})^2 + \frac{1}{c-1}\tau_\varepsilon^2$$

Here again, we see that when $\gamma = \Theta_c(1)$, we have tempered overfitting and $\omega_c(1) \leq \gamma \leq o_c(c^2)$, we have catastrophic overfitting and for $\gamma = \Omega_c(c^2)$ we get tempered overfitting again.

Equal Frobenius Norm . In this regime, we have that

$$R_{c<1} = \frac{\alpha_Z^2\tau^2}{1-c}(\boldsymbol{\beta}^\top \mathbf{u})^2 + \frac{c}{1-c}\tau_\varepsilon^2 \quad R_{c>1} = \frac{c\alpha_Z^2\tau^2}{c-1}(\boldsymbol{\beta}^\top \mathbf{u})^2 + \frac{1}{c-1}\tau_\varepsilon^2.$$

This generalizes the spike recovery setting studied in Sonthalia and Nadakuditi (2023), which assumed noiseless targets ($\tau_\varepsilon = 0$) and the equal Frobenius norm scaling. Our formula allows for observation noise and thus captures the more realistic case where the target \mathbf{y} itself contains randomness not aligned with the spike. Here we see that we have tempered overfitting unless $\tau^2 = o(1)$, which is the case considered in Sonthalia and Nadakuditi (2023).

E PROOF OF THEOREM 5

Theorem 5 (Generalization Risk). *Suppose Assumption 1, Assumption 2, and Assumption 3 hold.*

$$\mathcal{R} = \mathbb{E} \left[\underbrace{\left\| \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} \right\|_F^2}_{\text{Bias}} + \underbrace{\tau^2 \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Variance}} + \underbrace{\tilde{\alpha}_A^2 \left\| \beta_*^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Data Noise}} + \underbrace{\left(-2\tilde{\alpha}_A \beta_*^\top \tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top \beta_{int} \right)}_{\text{Target Alignment}} \right].$$

In particular, as $n, d \rightarrow \infty$ with $d/n \rightarrow c \in (0, \infty)$, we have the following expressions for each term.

Bias: For $c < 1$, we have that the bias term is

$$\tilde{\theta}^2 \left[(\beta_*^\top \mathbf{u})^2 \left(\tilde{\alpha}_Z - \alpha_Z + (\alpha_Z - \alpha_A) + \frac{\tau^2}{\theta^2 + \tau^2} \right)^2 + \tau_\varepsilon^2 \frac{c}{1-c} \frac{1}{d(\theta^2 + \tau^2)} \right].$$

If $c > 1$, we that the bias term is

$$\tilde{\theta}^2 (\beta_*^\top \mathbf{u})^2 \left(\tilde{\alpha}_Z - \alpha_Z + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\tau^2 c}{\theta^2 + \tau^2 c} \right)^2 + \tilde{\theta}^2 \left[\alpha_A^2 \frac{\|\beta_*\|^2}{d} \frac{c-1}{c} \frac{\theta^2 \tau^2 c}{(\theta^2 + \tau^2 c)^2} + \tau_\varepsilon^2 \frac{c}{c-1} \frac{\theta^2 + \tau^2}{n(\theta^2 + \tau^2 c)^2} \right].$$

Variance: For $c < 1$, we have that the variance term is

$$\alpha_A^2 \tilde{\tau}^2 \|\beta_*\|^2 + \tilde{\tau}^2 (\beta_*^\top \mathbf{u})^2 \left[\frac{1}{1-c} \frac{\theta^4 + \theta^2 \tau^2 c}{(\theta^2 + \tau^2)^2} (\alpha_Z - \alpha_A)^2 + 2\alpha_A (\alpha_Z - \alpha_A) \frac{\theta^2}{\theta^2 + \tau^2} \right] + \tau_\varepsilon^2 \frac{\tilde{\tau}^2}{\tau^2} \left[\frac{c}{1-c} - \frac{\theta^2}{d(\theta^2 + \tau^2)} \frac{c}{1-c} \right].$$

For $c > 1$, we have that the variance term is

$$\tilde{\tau}^2 \|\beta_*\|^2 \left(\frac{\alpha_A^2}{c} - \frac{\alpha_A^2}{d} \frac{\theta^2}{\theta^2 + \tau^2 c} \right) + \tilde{\tau}^2 (\beta_*^\top \mathbf{u})^2 \frac{c}{(c-1)} \frac{\theta^2}{\theta^2 + \tau^2 c} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + \tau_\varepsilon^2 \frac{\tilde{\tau}^2}{\tau^2} \left(\frac{1}{c-1} - \frac{\theta^2}{d(\theta^2 + \tau^2 c)} \frac{c}{c-1} \right).$$

Data Noise: For all c , we have that

$$\tilde{\alpha}_A^2 \tilde{\tau}^2 \|\beta_*\|^2.$$

Target Alignment: For $c < 1$, we have that the alignment term is

$$-2\tilde{\alpha}_A \tilde{\tau}^2 \left((\alpha_Z - \alpha_A) \frac{\theta^2}{\theta^2 + \tau^2} (\beta_*^\top \mathbf{u})^2 + \alpha_A \|\beta_*\|^2 \right).$$

For $c > 1$, we have that the alignment term is

$$-2\tilde{\alpha}_A \tilde{\tau}^2 \left(\left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\theta^2}{\theta^2 + \tau^2 c} (\beta_*^\top \mathbf{u})^2 + \alpha_A \|\beta_*\|^2 \left(\frac{1}{c} - \frac{1}{d} \frac{\theta^2}{\theta^2 + \tau^2 c} \right) \right).$$

Error terms: The largest error terms for all c are:

$$o(1) + O\left(\frac{1}{n}\right) = o(1).$$

Remark: We note that the above theorem is very general and captures all of the theorems in the main text as special cases. It is worth noting that the theorem also incorporates different signal and bulk strengths for test data, namely for $\tilde{\theta}$ and $\tilde{\tau}$.

The proof will be broken up into roughly 6 steps

0. **Rescale the problem** To apply standard results we rescale the problem. Section E.1

1. **Decompose the error into four terms.** We shall refer to these terms as the 1) bias, 2) variance, 3) data noise, and 4) target alignment. Section E.2
2. **Simplify the expressions.** We shall then use the result from Meyer (1973) to simplify the expression for each of the four terms. In particular, we shall express each term as the product of dependent functions of the eigenvalues of \mathbf{X} . Section E.3
3. **Random matrix theory estimate.** We then use standard results from random matrix theory such as Marchenko and Pastur (1967); Bai and Zhou (2008); Baik and Silverstein (2006) to obtain a closed-form formula of the building blocks for the risk. Section E.4
4. **Bound Products.** We then show that products of our building blocks concentrate. Step 4 (Section E.5) then collects the final terms.
5. **Undo Scaling** Step 5 (Section E.6) gives us back the correct scaling.

Section F has some generic probability lemmas that we need.

E.1 STEP 0: RESCALING

In order to better align with existing results and use them accordingly, we change our scalings for now and switch back after our derivation. That is, we divide everything by \sqrt{d} . Hence, we shall use

$$\frac{\theta}{\sqrt{d}} \mathbf{u} \mathbf{w}^\top = \theta \frac{\|\mathbf{w}\|}{\sqrt{d}} \mathbf{u} \frac{\mathbf{w}^\top}{\|\mathbf{w}\|}$$

as the spike. We shall let

$$\eta^2 := \theta^2 \frac{\|\mathbf{w}\|^2}{d} \quad \text{and} \quad \mathbf{v} := \frac{\mathbf{w}^\top}{\|\mathbf{w}\|}$$

Here, we treat \mathbf{v} as fixed unit norm vector and our spike is

$$\mathbf{Z}_r := \eta \mathbf{u} \mathbf{v}^\top$$

The \mathbf{A} noise after dividing by \sqrt{d} is

$$\mathbf{A}_r := \frac{\tau}{\sqrt{d}} N$$

where N are mean zero variance 1 entries. Here the appendix, we shall use the letter ρ for τ . Finally let

$$\mathbf{X}_r = \mathbf{Z}_r + \mathbf{A}_r$$

We can note that β_{int} , is still the solution to

$$\left\| \frac{\mathbf{y}}{\sqrt{d}} - \beta^\top \mathbf{X}_r \right\|^2, \quad \text{where} \quad \frac{\mathbf{y}}{\sqrt{d}} = \beta_*^\top (\mathbf{Z}_r + \mathbf{A}_r) + \frac{\boldsymbol{\varepsilon}}{\sqrt{d}}.$$

We define

$$\frac{\boldsymbol{\varepsilon}}{\sqrt{d}} =: \boldsymbol{\varepsilon}_r \sim \mathcal{N}\left(0, \frac{\tau_\varepsilon^2}{d}\right), \quad \tau_{\varepsilon,r}^2 := \frac{\tau_\varepsilon^2}{d}.$$

Then when we want to test, we shall look at the rescaled error

$$\frac{1}{\tilde{n}} \left\| \beta_*^\top (\tilde{\alpha}_Z \tilde{\mathbf{Z}}_r + \tilde{\alpha}_A \tilde{\mathbf{A}}_r) - \beta_{int}^\top (\tilde{\mathbf{Z}}_r + \tilde{\mathbf{A}}_r) \right\|_F^2$$

Through Steps 1 - 4, we shall drop the subscript r .

E.2 STEP 1: DECOMPOSE ERROR

Using the fact that $\tilde{\mathbf{A}}$ has been zero entries and is independent of $\tilde{\mathbf{Z}}$, we see that we can decompose the error as follows. Again here we consider \tilde{n} samples of test data and take the average (in expectation, this is the same as one test point).

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \beta_*^\top (\tilde{\alpha}_z \tilde{\mathbf{Z}} + \tilde{\alpha}_A \tilde{\mathbf{A}}) - \beta_{int}^\top (\tilde{\mathbf{Z}} + \tilde{\mathbf{A}}) \right\|_F^2 \right] \\ &= \mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} \right\|_F^2 \right] + \mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \tilde{\alpha}_A \beta_*^\top \tilde{\mathbf{A}} - \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2 \right] \\ &= \mathbb{E} \left[\underbrace{\frac{1}{\tilde{n}} \left\| \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} \right\|_F^2}_{\text{Bias}} + \underbrace{\frac{1}{\tilde{n}} \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Variance}} + \underbrace{\frac{1}{\tilde{n}} \tilde{\alpha}_A^2 \left\| \beta_*^\top \tilde{\mathbf{A}} \right\|_F^2}_{\text{Data Noise}} + \underbrace{\left(-\frac{2}{\tilde{n}} \tilde{\alpha}_A \beta_*^\top \tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top \beta_{int} \right)}_{\text{Target Alignment}} \right]. \end{aligned}$$

We compute these four terms one by one in the following sections.

E.3 STEP 2: SIMPLIFYING TERMS

This section simplifies the four terms. We begin by recalling results from prior work. We state them here for completeness.

Theorem 6 (Theorems 3, 5 of Meyer (1973)). *Define the following helper functions $\mathbf{h} = \mathbf{v}^\top \mathbf{A}^\dagger$, $\mathbf{k} = \mathbf{A}^\dagger \mathbf{u}$, $\mathbf{t} = \mathbf{v}^\top (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})$, $\xi = 1 + \eta \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u}$, $\mathbf{s} = (\mathbf{I} - \mathbf{A} \mathbf{A}^\dagger) \mathbf{u}$, $\gamma_1 = \eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2 + \xi^2$, $\gamma_2 = \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2$ and*

$$\begin{aligned} \mathbf{p}_1 &= -\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top - \eta \mathbf{k}, & \mathbf{q}_1^\top &= -\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger - \mathbf{h}, \\ \mathbf{p}_2 &= -\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top - \eta \mathbf{k}, & \mathbf{q}_2^\top &= -\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s}^\top - \mathbf{h}, \end{aligned}$$

Then we have that

$$(\mathbf{Z} + \mathbf{A})^\dagger = \begin{cases} \mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger - \frac{\xi}{\gamma_1} \mathbf{p}_1 \mathbf{q}_1^\top, & c < 1 \\ \mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top - \frac{\xi}{\gamma_2} \mathbf{p}_2 \mathbf{q}_2^\top, & c > 1 \end{cases}.$$

The following subsections - Bias E.3.1, Variance E.3.2, Data Noise E.3.3, and Target Alignment E.3.4 - present the linear algebraic simplifications of the results. To derive this results. We shall need some helper results that are presented in Section E.3.5.

E.3.1 BIAS

Using Lemma 5, we have that if $c < 1$

$$\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} = \left[\tilde{\alpha}_z - \alpha_z + \frac{\xi}{\gamma_1} (\alpha_z - \alpha_A) \right] \beta_*^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_1} \boldsymbol{\varepsilon}^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top,$$

and if $c > 1$

$$\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} = \beta_*^\top \left[(\tilde{\alpha}_z - \alpha_z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} - \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top.$$

The bias equals the expected squared norm of this term (divided by \tilde{n}).

E.3.2 VARIANCE

Lemma 8 gives us that

$$\begin{aligned} \mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2 \right] &= \mathbb{E} \left[\frac{\tilde{\tau}^2 \alpha_z^2}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{Z} \beta_* + \frac{\tilde{\tau}^2 \alpha_A^2}{d} \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_* \right. \\ &\quad \left. + \frac{2 \tilde{\tau}^2 \alpha_A \alpha_z}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_* + \frac{\tilde{\tau}^2}{d} \boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \boldsymbol{\varepsilon} \right]. \end{aligned}$$

E.3.3 DATA NOISE

The data noise term is the simplest to understand. Preliminary calculation gives us:

$$\frac{1}{\tilde{n}} \tilde{\alpha}_A^2 \mathbb{E}_{\tilde{A}} \left[\left\| \beta_*^\top \tilde{A} \right\|_F^2 \right] = \frac{\tilde{\alpha}_A^2 \tilde{\rho}^2 \tilde{n}}{\tilde{n} d} \|\beta_*\|^2 = \frac{\tilde{\alpha}_A^2 \tilde{\rho}^2}{d} \|\beta_*\|^2.$$

E.3.4 TARGET ALIGNMENT

To understand this term, we first note that \tilde{A} is independent of everything else. Hence we replace $\tilde{A}\tilde{A}^\top$ with its expectation $\frac{\tilde{\rho}^2 \tilde{n}}{d} \mathbf{I}$.

$$\mathbb{E}_{\tilde{A}} \left[-\frac{2}{\tilde{n}} \tilde{\alpha}_A \beta_*^\top \tilde{A} \tilde{A}^\top \beta_{int} \right] = -\frac{2}{\tilde{n}} \frac{\tilde{\rho}^2 \tilde{n}}{d} \tilde{\alpha}_A \beta_*^\top \beta_{int} = -\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \beta_*^\top \beta_{int}.$$

Since ε has mean-zero entries that are independent of everything else. We see that

$$\mathbb{E}_\varepsilon [\beta_*^\top \beta_{int}] = \mathbb{E}_\varepsilon \left[\beta_*^\top \left((\alpha_z \beta_*^\top \mathbf{Z} + \varepsilon^\top) (\mathbf{Z} + \mathbf{A})^\dagger + \alpha_A \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger \right)^\top \right] \quad (5)$$

$$= \beta_*^\top \left(\alpha_z \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger - \alpha_A \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger \right)^\top \quad (6)$$

$$= \alpha_z \beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{Z}^\top \beta_* + \alpha_A \beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_*. \quad (7)$$

E.3.5 HELPER LEMMAS

Proposition 1 (Proposition 2 from Sonthalia and Nadakuditi (2023)). *In the setting from Section 2*

$$\mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger = \begin{cases} \frac{\eta \xi}{\gamma_1} \mathbf{u} \mathbf{h} + \frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger, & c < 1 \\ \frac{\eta \xi}{\gamma_2} \mathbf{u} \mathbf{h} + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top, & c > 1 \end{cases}.$$

Lemma 1. *If $\xi \neq 0$ and \mathbf{A} has full rank, we have:*

$$\varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} = \begin{cases} -\frac{\tilde{\eta} \xi}{\eta \gamma_1} \varepsilon^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top & c < 1 \\ -\frac{\tilde{\eta} \xi}{\eta \gamma_2} \varepsilon^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top & c > 1 \end{cases}.$$

Proof. After substitutions, Proposition 1 implies that for $c < 1$, $\varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}}$ becomes:

$$\begin{aligned} & \varepsilon^\top \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger - \frac{\xi}{\gamma_1} \mathbf{p}_1 \left(-\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger - \mathbf{h} \right) \right) \tilde{\mathbf{Z}} \\ &= \tilde{\eta} \varepsilon^\top \left(\mathbf{A}^\dagger \mathbf{u} \tilde{\mathbf{v}}^\top + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{u} \tilde{\mathbf{v}}^\top - \frac{\xi}{\gamma_1} \mathbf{p}_1 \left(-\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{u} - \mathbf{h} \mathbf{u} \right) \tilde{\mathbf{v}}^\top \right) \quad \text{by } \tilde{\mathbf{Z}} = \tilde{\eta} \mathbf{u} \tilde{\mathbf{v}}^\top. \end{aligned}$$

Since $\mathbf{k} = \mathbf{A}^\dagger \mathbf{u}$ and $\mathbf{h} \mathbf{u} = \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u} = \frac{\xi - 1}{\eta}$, we then have that

$$\begin{aligned} & \tilde{\eta} \varepsilon^\top \left(\mathbf{A}^\dagger \mathbf{u} \tilde{\mathbf{v}}^\top + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{u} \tilde{\mathbf{v}}^\top - \frac{\xi}{\gamma_1} \mathbf{p}_1 \left(-\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{u} - \mathbf{h} \mathbf{u} \right) \tilde{\mathbf{v}}^\top \right) \\ &= \tilde{\eta} \varepsilon^\top \left(\mathbf{k} \tilde{\mathbf{v}}^\top + \frac{\eta \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top \tilde{\mathbf{v}}^\top + \frac{\xi}{\gamma_1} \mathbf{p}_1 \left(\frac{\eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2 + \xi^2 - \xi}{\xi \eta} \right) \tilde{\mathbf{v}}^\top \right) \\ &= \tilde{\eta} \varepsilon^\top \left(\mathbf{k} \tilde{\mathbf{v}}^\top + \frac{\eta \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top \tilde{\mathbf{v}}^\top + \frac{1}{\gamma_1} \mathbf{p}_1 \left(\frac{\gamma_1 - \xi}{\eta} \right) \tilde{\mathbf{v}}^\top \right) \\ &= \tilde{\eta} \varepsilon^\top \left(\frac{1}{\eta} \left(\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top + \eta \mathbf{k} \right) \tilde{\mathbf{v}}^\top + \frac{1}{\eta} \mathbf{p}_1 \tilde{\mathbf{v}}^\top - \frac{\xi}{\eta \gamma_1} \mathbf{p}_1 \tilde{\mathbf{v}}^\top \right) \\ &= \varepsilon^\top \left(-\frac{\tilde{\eta}}{\eta} \mathbf{p}_1 \tilde{\mathbf{v}}^\top + \frac{\tilde{\eta}}{\eta} \mathbf{p}_1 \tilde{\mathbf{v}}^\top - \frac{\tilde{\eta} \xi}{\eta \gamma_1} \mathbf{p}_1 \tilde{\mathbf{v}}^\top \right) \\ &= -\frac{\tilde{\eta} \xi}{\eta \gamma_1} \varepsilon^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top. \end{aligned}$$

For $c > 1$, we note that the calculation is exactly the same. An example of such a calculation can be seen in the proof of Lemma 4. \square

Lemma 2. *In the setting of Section 2, we have:*

$$\mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger = \begin{cases} \mathbf{I} - \frac{\eta\xi}{\gamma_1}\mathbf{u}\mathbf{h} + \frac{\eta^2\|\mathbf{t}\|^2}{\gamma_1}\mathbf{u}\mathbf{k}^\top\mathbf{A}^\dagger, & c < 1 \\ \mathbf{A}\mathbf{A}^\dagger + \frac{\eta\xi}{\gamma_2}\mathbf{h}^\top\mathbf{s}^\top - \frac{\eta^2\|\mathbf{s}\|^2}{\gamma_2}\mathbf{h}^\top\mathbf{h} - \frac{\eta^2\|\mathbf{h}\|^2}{\gamma_2}\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{s}^\top - \frac{\eta\xi}{\gamma_2}\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{h}, & c > 1 \end{cases}.$$

Proof. For $c < 1$, \mathbf{Z}, \mathbf{A} are $d \times n$ with $d < n$. Since \mathbf{A} is assumed to have full rank, $\mathbf{Z} + \mathbf{A}$ has full rank with probability 1, and hence

$$(\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger = \mathbf{I}.$$

Thus, from Proposition 1,

$$\mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger = (\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger - \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger = \mathbf{I} - \frac{\eta\xi}{\gamma_1}\mathbf{u}\mathbf{h} - \frac{\eta^2\|\mathbf{t}\|^2}{\gamma_1}\mathbf{u}\mathbf{k}^\top\mathbf{A}^\dagger.$$

For $c > 1$, since $(\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger$ is no longer the identity matrix, we directly expand using Theorem 6:

$$\begin{aligned} \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger &= \mathbf{A} \left(\mathbf{A}^\dagger + \frac{\eta}{\xi}\mathbf{A}^\dagger\mathbf{h}^\top\mathbf{s}^\top - \frac{\xi}{\gamma_2} \left(\frac{\eta^2\|\mathbf{s}\|^2}{\xi}\mathbf{A}^\dagger\mathbf{h}^\top + \eta\mathbf{k} \right) \left(\frac{\eta\|\mathbf{h}\|^2}{\xi}\mathbf{s}^\top + \mathbf{h} \right) \right) \\ &= \mathbf{A}\mathbf{A}^\dagger + \frac{\eta}{\xi}\mathbf{A}\mathbf{A}^\dagger\mathbf{h}^\top\mathbf{s}^\top - \frac{\xi}{\gamma_2} \left(\frac{\eta^2\|\mathbf{s}\|^2}{\xi}\mathbf{A}\mathbf{A}^\dagger\mathbf{h}^\top + \eta\mathbf{A}\mathbf{A}^\dagger\mathbf{u} \right) \left(\frac{\eta\|\mathbf{h}\|^2}{\xi}\mathbf{s}^\top + \mathbf{h} \right). \end{aligned}$$

Noting that $\mathbf{A}\mathbf{A}^\dagger\mathbf{h}^\top = \mathbf{A}\mathbf{A}^\dagger\mathbf{A}^{\dagger\top}\mathbf{v} = \mathbf{A}^{\dagger\top}\mathbf{v} = \mathbf{h}^\top$, we have

$$\begin{aligned} \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger &= \mathbf{A}\mathbf{A}^\dagger + \frac{\eta}{\xi}\mathbf{h}^\top\mathbf{s}^\top - \frac{\xi}{\gamma_2} \left(\frac{\eta^2\|\mathbf{s}\|^2}{\xi}\mathbf{h}^\top + \eta\mathbf{A}\mathbf{A}^\dagger\mathbf{u} \right) \left(\frac{\eta\|\mathbf{h}\|^2}{\xi}\mathbf{s}^\top + \mathbf{h} \right) \\ &= \mathbf{A}\mathbf{A}^\dagger + \frac{\eta}{\xi}\mathbf{h}^\top\mathbf{s}^\top - \frac{\eta^3\|\mathbf{s}\|^2\|\mathbf{h}\|^2}{\xi\gamma_2}\mathbf{h}^\top\mathbf{s}^\top - \frac{\eta^2\|\mathbf{s}\|^2}{\gamma_2}\mathbf{h}^\top\mathbf{h} - \frac{\eta^2\|\mathbf{h}\|^2}{\gamma_2}\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{s}^\top - \frac{\eta\xi}{\gamma_2}\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{h}. \end{aligned}$$

We can combine the coefficients in front of $\mathbf{h}^\top\mathbf{s}^\top$ to get

$$\frac{\eta}{\xi} - \frac{\eta^3\|\mathbf{s}\|^2\|\mathbf{h}\|^2}{\xi\gamma_2} = \frac{\eta(\eta^2\|\mathbf{s}\|^2\|\mathbf{h}\|^2 + \xi^2) - \eta^3\|\mathbf{s}\|^2\|\mathbf{h}\|^2}{\xi\gamma_2} = \frac{\eta\xi}{\gamma_2}.$$

The statement follows from here. \square

Lemma 3. *If $\xi \neq 0$ and \mathbf{A} has full rank, we have:*

$$\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} = \begin{cases} \left(1 - \frac{\xi}{\gamma_1}\right) \beta_*^\top \tilde{\mathbf{Z}} & c < 1 \\ \left(1 - \frac{\xi}{\gamma_2}\right) \beta_*^\top \tilde{\mathbf{Z}} & c > 1 \end{cases}.$$

Proof. Using Proposition 1 for $c < 1$ and $\tilde{\mathbf{Z}} = \tilde{\eta}\mathbf{u}\tilde{\mathbf{v}}^\top$, we have that

$$\begin{aligned} \beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} &= \beta_*^\top \left(\frac{\eta\xi}{\gamma_1}\mathbf{u}\mathbf{h} + \frac{\eta^2\|\mathbf{t}\|^2}{\gamma_1}\mathbf{u}\mathbf{k}^\top\mathbf{A}^\dagger \right) \tilde{\mathbf{Z}} \\ &= \tilde{\eta}\beta_*^\top \left(\frac{\eta\xi}{\gamma_1}\mathbf{u}\mathbf{h}\mathbf{u}\tilde{\mathbf{v}}^\top + \frac{\eta^2\|\mathbf{t}\|^2}{\gamma_1}\mathbf{u}\mathbf{k}^\top\mathbf{A}^\dagger\mathbf{u}\tilde{\mathbf{v}}^\top \right) \\ &= \tilde{\eta}\beta_*^\top \left(\frac{\eta\xi}{\gamma_1}\mathbf{u}\mathbf{v}^\top\mathbf{A}^\dagger\mathbf{u}\tilde{\mathbf{v}}^\top + \frac{\eta^2\|\mathbf{t}\|^2}{\gamma_1}\mathbf{u}\mathbf{k}^\top\mathbf{A}^\dagger\mathbf{u}\tilde{\mathbf{v}}^\top \right). \end{aligned}$$

Note $\xi - 1 = \eta\mathbf{v}^\top\mathbf{A}^\dagger\mathbf{u}$, $\mathbf{k}\mathbf{A}^\dagger\mathbf{u} = \mathbf{k}^\top\mathbf{k} = \|\mathbf{k}\|^2$. The above equation becomes

$$\tilde{\eta}\beta_*^\top \left(\frac{\xi(\xi - 1)}{\gamma_1} + \frac{\eta^2\|\mathbf{t}\|^2\|\mathbf{k}\|^2}{\gamma_1} \right) \mathbf{u}\tilde{\mathbf{v}}^\top = \beta_*^\top \left(\frac{\xi(\xi - 1)}{\gamma_1} + \frac{\eta^2\|\mathbf{t}\|^2\|\mathbf{k}\|^2}{\gamma_1} \right) \tilde{\mathbf{Z}}^\top.$$

Using $\gamma_1 = \eta^2\|\mathbf{t}\|^2\|\mathbf{k}\|^2 + \xi^2$ to combine the coefficients, we have that

$$\frac{\xi(\xi - 1)}{\gamma_1} + \frac{\eta^2\|\mathbf{t}\|^2\|\mathbf{k}\|^2}{\gamma_1} = \frac{-\xi + \xi^2 + \eta^2\|\mathbf{t}\|^2\|\mathbf{k}\|^2}{\gamma_1} = \frac{-\xi + \gamma_1}{\gamma_1} = 1 - \frac{\xi}{\gamma_1}.$$

This completes the proof for $c < 1$. Similarly, for $c > 1$, we obtain

$$\begin{aligned}\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} &= \beta_*^\top \left(\frac{\eta^\xi \mathbf{u} \mathbf{h}}{\gamma_2} + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top \right) \tilde{\mathbf{Z}} \\ &= \tilde{\eta} \beta_*^\top \left(\frac{\eta^\xi \mathbf{u} \mathbf{h} \mathbf{u} \tilde{\mathbf{v}}^\top}{\gamma_2} + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top \mathbf{u} \tilde{\mathbf{v}}^\top \right) \\ &= \tilde{\eta} \beta_*^\top \left(\frac{\eta^\xi \mathbf{u} \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u} \tilde{\mathbf{v}}^\top}{\gamma_2} + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top \mathbf{u} \tilde{\mathbf{v}}^\top \right).\end{aligned}$$

Note $\xi - 1 = \eta \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u}$, $\mathbf{s}^\top \mathbf{u} = \|\mathbf{s}\|^2$. The above equation becomes

$$\tilde{\eta} \beta_*^\top \left(\frac{\xi(\xi - 1)}{\gamma_2} + \frac{\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2} \right) \mathbf{u} \tilde{\mathbf{v}}^\top = \beta_*^\top \left(\frac{\xi(\xi - 1)}{\gamma_2} + \frac{\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2} \right) \tilde{\mathbf{Z}}^\top.$$

Using $\gamma_2 = \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2$ to combine the coefficients, we have that

$$\frac{\xi(\xi - 1)}{\gamma_2} + \frac{\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2} = \frac{-\xi + \xi^2 + \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2} = \frac{-\xi + \gamma_2}{\gamma_2} = 1 - \frac{\xi}{\gamma_2}.$$

The target expression follows. □

Lemma 4. *If $\xi \neq 0$ and \mathbf{A} has full rank, we have:*

$$\beta_*^\top \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} = \begin{cases} \frac{\xi}{\gamma_1} \beta_*^\top \tilde{\mathbf{Z}} & c < 1 \\ \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\xi}{\gamma_2} \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \tilde{\mathbf{Z}} & c > 1 \end{cases}.$$

Proof. We begin with $c < 1$. Since \mathbf{A} is assumed to have full rank, $\mathbf{Z} + \mathbf{A}$ has full column rank with probability 1, and hence

$$(\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger = \mathbf{I}.$$

It follows from Lemma 3 that

$$\begin{aligned}\beta_*^\top \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} &= \beta_*^\top (\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} - \beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} \\ &= \beta_*^\top \tilde{\mathbf{Z}} - \left(1 - \frac{\xi}{\gamma_1}\right) \beta_*^\top \tilde{\mathbf{Z}} = \frac{\xi}{\gamma_1} \beta_*^\top \tilde{\mathbf{Z}}.\end{aligned}$$

For $c > 1$, $\mathbf{Z} + \mathbf{A}$ now has full row rank instead of full column rank. Hence, we do not have $(\mathbf{Z} + \mathbf{A})(\mathbf{Z} + \mathbf{A})^\dagger = \mathbf{I}$ and need to directly expand it using Theorem 6 and its helper variables:

$$\begin{aligned}
\beta_*^\top \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} &= \beta_*^\top \mathbf{A} \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top - \frac{\xi}{\gamma_2} \mathbf{p}_2 \mathbf{q}_2^\top \right) \tilde{\mathbf{Z}} \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(\mathbf{k} \tilde{\mathbf{v}}^\top + \frac{\eta \|\mathbf{s}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top \tilde{\mathbf{v}}^\top - \frac{\xi}{\gamma_2} \mathbf{p}_2 \mathbf{q}_2^\top \mathbf{u} \tilde{\mathbf{v}}^\top \right) \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(-\frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top - \frac{\xi}{\gamma_2} \mathbf{p}_2 \left(-\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s}^\top - \mathbf{h} \right) \mathbf{u} \tilde{\mathbf{v}}^\top \right) \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(-\frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top + \frac{\xi}{\gamma_2} \mathbf{p}_2 \left(\frac{\eta \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\xi} + \frac{\xi - 1}{\eta} \right) \tilde{\mathbf{v}}^\top \right) \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(-\frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top + \frac{\xi}{\gamma_2} \mathbf{p}_2 \left(\frac{\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2 - \xi}{\xi \eta} \right) \tilde{\mathbf{v}}^\top \right) \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(-\frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top + \frac{\xi}{\gamma_2} \mathbf{p}_2 \left(\frac{\gamma_2 - \xi}{\xi \eta} \right) \tilde{\mathbf{v}}^\top \right) \\
&= \tilde{\eta} \beta_*^\top \mathbf{A} \left(-\frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top + \frac{1}{\eta} \mathbf{p}_2 \tilde{\mathbf{v}}^\top - \frac{\xi}{\eta \gamma_2} \mathbf{p}_2 \tilde{\mathbf{v}}^\top \right) \\
&= -\frac{\tilde{\eta} \xi}{\eta \gamma_2} \beta_*^\top \mathbf{A} \mathbf{p}_2 \tilde{\mathbf{v}}^\top \\
&= \frac{\tilde{\eta} \xi}{\eta \gamma_2} \beta_*^\top \left(\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h}^\top + \eta \mathbf{A} \mathbf{k} \right) \tilde{\mathbf{v}}^\top \quad \text{by plugging in the expression of } \mathbf{p}_2 \\
&= \frac{\tilde{\eta} \eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \tilde{\mathbf{v}}^\top + \frac{\xi}{\gamma_2} \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \tilde{\mathbf{Z}} \quad \text{by } \tilde{\eta} \mathbf{k} \tilde{\mathbf{v}}^\top = \mathbf{A}^\dagger \tilde{\eta} \mathbf{u} \tilde{\mathbf{v}}^\top = \mathbf{A}^\dagger \tilde{\mathbf{Z}}.
\end{aligned}$$

Noting that $\beta_*^\top \mathbf{h}^\top$ is a scalar, we then introduce $\mathbf{1} = \mathbf{u}^\top \mathbf{u}$ and get that

$$\frac{\tilde{\eta} \eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \mathbf{u} \tilde{\mathbf{v}}^\top = \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} \quad \text{since } \tilde{\eta} \mathbf{u} \tilde{\mathbf{v}}^\top = \tilde{\mathbf{Z}}.$$

Thus, the final expression is

$$\frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\xi}{\gamma_2} \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \tilde{\mathbf{Z}}.$$

□

Lemma 5 (Bias Term). *In the setting of Section 2, we have that if $c < 1$,*

$$\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} = \left[\tilde{\alpha}_z - \alpha_z + \frac{\xi}{\gamma_1} (\alpha_z - \alpha_A) \right] \beta_*^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_1} \boldsymbol{\varepsilon}^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top,$$

and if $c > 1$,

$$\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} = \beta_*^\top \left[(\tilde{\alpha}_z - \alpha_z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} - \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top.$$

Proof. To simplify the bias term, we first need the following expansion:

$$\begin{aligned}
\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \beta_{int}^\top \tilde{\mathbf{Z}} &= \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - (\beta_*^\top (\alpha_z \mathbf{Z} + \alpha_A \mathbf{A}) + \boldsymbol{\varepsilon}^\top) (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} \\
&= \tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \alpha_z \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger - \alpha_A \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}} - \boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}}.
\end{aligned}$$

From Lemmas 1, 3, 4, we get simplified expressions for $\boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}}$, $\beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger \tilde{\mathbf{Z}}$, $\beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger$ and plug them in. For $c < 1$, we get

$$\begin{aligned}
&\tilde{\alpha}_z \beta_*^\top \tilde{\mathbf{Z}} - \alpha_z \left(1 - \frac{\xi}{\gamma_1} \right) \beta_*^\top \tilde{\mathbf{Z}} - \alpha_A \frac{\xi}{\gamma_1} \beta_*^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_1} \boldsymbol{\varepsilon}^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top \\
&= \left[\tilde{\alpha}_z - \alpha_z + \frac{\xi}{\gamma_1} (\alpha_z - \alpha_A) \right] \beta_*^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_1} \boldsymbol{\varepsilon}^\top \mathbf{p}_1 \tilde{\mathbf{v}}^\top.
\end{aligned}$$

On the other hand, for $c > 1$, we have

$$\begin{aligned} & \tilde{\alpha}_Z \beta_*^\top \tilde{\mathbf{Z}} - \alpha_Z \left(1 - \frac{\xi}{\gamma_2}\right) \beta_*^\top \tilde{\mathbf{Z}} - \alpha_A \left[\frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\xi}{\gamma_2} \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \tilde{\mathbf{Z}} \right] + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top \\ &= \beta_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} - \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top. \end{aligned}$$

□

Lemma 6 (Squared Norms of \mathbf{p}_1 and \mathbf{p}_2). Recall $\mathbf{p}_1 = -\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top - \eta \mathbf{k}$ and $\mathbf{p}_2 = -\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{h} - \eta \mathbf{k}$.

1. $\|\mathbf{p}_1\|^2 = \frac{\eta^2 \|\mathbf{k}\|^2}{\xi^2} \gamma_1$.
2. $\|\mathbf{p}_2\|^2 = \frac{\eta^4 \|\mathbf{s}\|^4}{\xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top + \frac{2\eta^3 \|\mathbf{s}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \eta^2 \|\mathbf{k}\|^2$.

Proof. For \mathbf{p}_1 , we have

$$\|\mathbf{p}_1\|^2 = \left(-\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t} - \eta \mathbf{k} \right) \left(-\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top - \eta \mathbf{k}^\top \right) = \left(\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \right)^2 \|\mathbf{t}\|^2 + 2 \frac{\eta^3 \|\mathbf{k}\|^2}{\xi} \mathbf{t} \mathbf{k} + \eta^2 \|\mathbf{k}\|^2.$$

Using $\mathbf{t} \mathbf{k} = \mathbf{0}$ yields the first result, which we can further simplify as

$$\frac{\eta^2 \|\mathbf{k}\|^2}{\xi^2} (\eta^2 \|\mathbf{k}\|^2 \|\mathbf{t}\|^2 + \xi^2) = \frac{\eta^2 \|\mathbf{k}\|^2}{\xi^2} \gamma_1.$$

For \mathbf{p}_2 , similarly, we have

$$\begin{aligned} \|\mathbf{p}_2\|^2 &= \left(-\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h} \mathbf{A}^\dagger - \eta \mathbf{k} \right) \left(-\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top - \eta \mathbf{k}^\top \right) \\ &= \frac{\eta^4 \|\mathbf{s}\|^4}{\xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top + \frac{2\eta^3 \|\mathbf{s}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \eta^2 \|\mathbf{k}\|^2. \end{aligned}$$

□

Lemma 7 (Squared Norms of \mathbf{q}_1 and \mathbf{q}_2). Let $\mathbf{q}_1^\top = -\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger - \mathbf{h}$ and $\mathbf{q}_2^\top = -\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s}^\top - \mathbf{h}$.

1. $\|\mathbf{q}_1\|^2 = \frac{\eta^2 \|\mathbf{t}\|^4}{\xi^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} + \frac{2\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \|\mathbf{h}\|^2$.
2. $\|\mathbf{q}_2\|^2 = \frac{\|\mathbf{h}\|^2}{\xi^2} \gamma_2$.

Proof. Similar to Lemma 6, we directly expand the two terms:

$$\|\mathbf{q}_1\|^2 = \left(-\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger - \mathbf{h} \right) \left(-\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{k} - \mathbf{h}^\top \right) = \frac{\eta^2 \|\mathbf{t}\|^4}{\xi^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} + \frac{2\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \|\mathbf{h}\|^2.$$

$$\begin{aligned} \|\mathbf{q}_2\|^2 &= \left(-\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s}^\top - \mathbf{h} \right) \left(-\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s} - \mathbf{h}^\top \right) = \frac{\eta^2 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\xi^2} + \|\mathbf{h}\|^2 \quad \text{since } \mathbf{h} \mathbf{s} = \mathbf{0} \\ &= \frac{\|\mathbf{h}\|^2 (\eta^2 \|\mathbf{h}\|^2 \|\mathbf{s}\|^2 + \xi^2)}{\xi^2} \\ &= \frac{\|\mathbf{h}\|^2}{\xi^2} \gamma_2. \end{aligned}$$

□

Lemma 8 (Preliminary Expansion of Variance). *In the setting of Section 2, we have*

$$\mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2 \right] = \mathbb{E} \left[\frac{\tilde{\tau}^2 \alpha_z^2}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{Z} \beta_* + \frac{\tilde{\tau}^2 \alpha_A^2}{d} \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}^\top \beta_* \right. \\ \left. + \frac{2\tilde{\tau}^2 \alpha_A \alpha_z}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}^\top \beta_* + \frac{\tilde{\tau}^2}{d} \varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \varepsilon \right].$$

Proof. Since $\tilde{\mathbf{A}}$ is independent of the other terms, we replace $\tilde{\mathbf{A}}\tilde{\mathbf{A}}^\top$ with its expectation $\frac{\tilde{\tau}^2 \tilde{n}}{d} \mathbf{I}$.

$$\mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \beta_{int}^\top \tilde{\mathbf{A}} \right\|_F^2 \right] = \mathbb{E} \left[\frac{1}{\tilde{n}} \beta_{int}^\top \tilde{\mathbf{A}} \tilde{\mathbf{A}}^\top \beta_{int} \right] = \frac{1}{\tilde{n}} \frac{\tilde{\tau}^2 \tilde{n}}{d} \mathbb{E} [\beta_{int}^\top \beta_{int}] = \frac{\tilde{\tau}^2}{d} \mathbb{E} [\|\beta_{int}\|^2].$$

We now plug in the expression for β_{int} . Since ε is a zero-mean vector and independent from other random variables, terms with only one ε have zero expectation. A straightforward expansion gives:

$$\frac{\tilde{\tau}^2}{d} \|\beta_{int}\|_F^2 = \frac{\tilde{\tau}^2}{d} (\beta_*^\top (\alpha_z \mathbf{Z} + \alpha_A \mathbf{A}) + \varepsilon^\top) (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} (\beta_*^\top (\alpha_z \mathbf{Z} + \alpha_A \mathbf{A}) + \varepsilon^\top)^\top.$$

After eliminating zero expectations as above, the expectation becomes:

$$\mathbb{E} \left[\frac{\tilde{\tau}^2}{d} \|\beta_{int}\|_F^2 \right] = \mathbb{E} \left[\frac{\tilde{\tau}^2 \alpha_z^2}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{Z} \beta_* + \frac{\tilde{\tau}^2 \alpha_A^2}{d} \beta_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}^\top \beta_* \right. \\ \left. + \frac{2\tilde{\tau}^2 \alpha_A \alpha_z}{d} \beta_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}^\top \beta_* + \frac{\tilde{\tau}^2}{d} \varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \varepsilon \right].$$

□

E.4 STEP 3: RANDOM MATRIX THEORY ESTIMATES

To do the estimates we recall the set up. In particular, we have that

$$\mathbf{Z} = \eta \mathbf{u} \mathbf{v}^\top, \quad \text{where } \theta = \frac{\eta}{\sqrt{n}} \text{ and } \|\mathbf{v}\| = 1,$$

and the entries of

$$A_{ij} = \mathcal{N} \left(0, \frac{\rho^2}{d} \right)$$

Recall the following definition $\mathbf{h} = \mathbf{v}^\top \mathbf{A}^\dagger$, $\mathbf{k} = \mathbf{A}^\dagger \mathbf{u}$, $\mathbf{t} = \mathbf{v}^\top (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})$, $\xi = 1 + \eta \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u}$, $\mathbf{s} = (\mathbf{I} - \mathbf{A} \mathbf{A}^\dagger) \mathbf{u}$, $\gamma_1 = \eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2 + \xi^2$, $\gamma_2 = \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2$ and

$$\mathbf{p}_1 = -\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t}^\top - \eta \mathbf{k}, \quad \mathbf{q}_1^\top = -\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger - \mathbf{h}, \\ \mathbf{p}_2 = -\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top - \eta \mathbf{k}, \quad \mathbf{q}_2^\top = -\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s}^\top - \mathbf{h},$$

To show that each of the four terms, bias, variance, data noise, and target alignment concentrate in the limit, we do this in two steps.

- First, we compute the mean and variance for basic building blocks such as $\|\mathbf{h}\|^2$ and other variables. Section E.4.1.
- Second, we provide bounds on the higher moments. Section E.4.2.
- Next, we prove bounds on the moments of γ_i . Section E.4.3.

E.4.1 STEP 3(A): SHOWING THAT BASIC BUILDING BLOCKS CONCENTRATE

We begin by bounding the mean and variance.

Lemma 9 (Generalized version of Lemma 7 from Sonthalia and Nadakuditi (2023)). *Suppose A_{ij} have mean 0 and variance ρ^2/d , the entries are uncorrelated, have finite fourth moment, the distribution is invariant under left and right orthogonal transformation and the empirical spectral distribution of $\frac{1}{\rho^2} \mathbf{A} \mathbf{A}^\top$ converges to the Marchenko-Pastur law. Additionally, if \mathbf{u} and \mathbf{v} are fixed unit norm vectors. Then we have that*

1. $\mathbb{E}[\|\mathbf{h}\|^2] = \begin{cases} \frac{1}{\rho^2} \frac{c^2}{1-c} & c < 1 \\ \frac{1}{\rho^2} \frac{c}{c-1} & c > 1 \end{cases} + o\left(\frac{1}{\rho^2}\right)$ and $\text{Var}(\|\mathbf{h}\|^2) = O\left(\frac{1}{\rho^4 n}\right)$.
2. $\mathbb{E}[\|\mathbf{k}\|^2] = \frac{1}{\rho^2} \frac{c}{1-c} + o\left(\frac{1}{\rho^2}\right)$ and $\text{Var}(\|\mathbf{k}\|^2) = O\left(\frac{1}{\rho^4 n}\right)$.
3. $\mathbb{E}[\|\mathbf{s}\|^2] = 1 - \frac{1}{c}$ and $\text{Var}(\|\mathbf{s}\|^2) = O\left(\frac{1}{d}\right)$.
4. $\mathbb{E}[\|\mathbf{t}\|^2] = 1 - c$ and $\text{Var}(\|\mathbf{t}\|^2) = O\left(\frac{1}{n}\right)$.
5. $\mathbb{E}\left[\frac{\xi}{\eta}\right] = \frac{1}{\eta}$ and $\text{Var}\left(\frac{\xi}{\eta}\right) = O\left(\frac{1}{\max(n, d) \rho^2}\right)$.
6. $\mathbb{E}\left[\frac{\xi^2}{\eta^2}\right] = \frac{1}{\eta^2} + \frac{1}{\max(n, d)} \frac{c}{\rho^2 |1-c|} + o\left(\frac{1}{\max(n, d) \rho^2}\right) = \frac{1}{\eta^2} + O\left(\frac{1}{\max(n, d) \rho^2}\right)$
and $\text{Var}\left(\frac{\xi^2}{\eta^2}\right) = O\left(\frac{1}{\max(d, n) \rho^4}\right)$.

Note that here $\max(d, n)$, d , n are interchangeable in the variance big-Oh terms since they only differ by an absolute constant c . We include the details for completion.

Proof. Items 1 – 5 come from the original statement, which assumes unit variance. Here our variance parameter ρ simply induces a multiplicative change. We now focus on item 6.

Let $\zeta = \xi/\eta = 1/\eta + \mathbf{v}^\top \mathbf{A}^\dagger \mathbf{u}$. With $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ (SVD), $\mathbf{A} \in \mathbb{R}^{d \times n}$ having i.i.d. $\mathcal{N}(0, \rho^2/d)$ entries, and \mathbf{u}, \mathbf{v} fixed unit vectors, we have $\zeta = \frac{1}{\eta} + \sum_{i=1}^r \frac{1}{\sigma_i} b_i a_i$, where $r = \min(d, n)$, $\mathbf{a} = \mathbf{V}^\top \mathbf{v}$, $\mathbf{b} = \mathbf{U}^\top \mathbf{u}$ are uniformly random on S^{n-1} and S^{d-1} respectively since \mathbf{U}, \mathbf{V} are random rotations.

Since \mathbf{A} has zero-mean entries, only the non-cross terms remain in the expectation, and the fourth moment is

$$\mathbb{E}[\zeta^4] = \frac{1}{\eta^4} + \frac{6}{\eta^2} \sum_{i,j} \mathbb{E}\left[\frac{1}{\sigma_i \sigma_j}\right] \mathbb{E}[b_i b_j] \mathbb{E}[a_i a_j] + \sum_{i,j,k,l} \mathbb{E}\left[\frac{1}{\sigma_i \sigma_j \sigma_k \sigma_l}\right] \mathbb{E}[b_i b_j b_k b_l] \mathbb{E}[a_i a_j a_k a_l].$$

Furthermore, non-zero expectation terms require paired indices (since odd moments of the uniformly random vector on the sphere equals 0). In particular, using exact spherical moments, we have $\mathbb{E}[a_i^4] = \frac{3}{n(n+2)}$, $\mathbb{E}[a_i^2] = \frac{1}{n}$, $\mathbb{E}[a_i^2 a_j^2] = \frac{1}{n(n+2)}$ ($i \neq j$), $\mathbb{E}[b_i^4] = \frac{3}{d(d+2)}$, $\mathbb{E}[b_i^2] = \frac{1}{d}$, $\mathbb{E}[b_i^2 b_j^2] = \frac{1}{d(d+2)}$ ($i \neq j$):

$$\begin{aligned} \mathbb{E}[\zeta^4] &= \frac{1}{\eta^4} + \frac{6}{\eta^2} \sum_{i=1}^r \mathbb{E}\left[\frac{1}{\sigma_i^2}\right] \frac{1}{dn} + \sum_{i=1}^r \mathbb{E}\left[\frac{1}{\sigma_i^4}\right] \frac{9}{d(d+2)n(n+2)} + 3 \sum_{i \neq k} \mathbb{E}\left[\frac{1}{\sigma_i^2 \sigma_k^2}\right] \frac{1}{d(d+2)n(n+2)} \\ &= \frac{1}{\eta^4} + \underbrace{\frac{9 \sum_{i=1}^r \mathbb{E}[1/\sigma_i^4]}{d(d+2)n(n+2)}}_{I_1} + \underbrace{\frac{3 \sum_{i \neq k} \mathbb{E}[1/(\sigma_i^2 \sigma_k^2)]}{d(d+2)n(n+2)}}_{I_2} + \underbrace{\frac{6 \sum_{i=1}^r \mathbb{E}[1/\sigma_i^2]}{\eta^2 dn}}_{I_3}. \end{aligned}$$

Leading Order Scaling and Mean. Let $N = \max(d, n)$, assume $n, d \rightarrow \infty$ with $d/n \rightarrow c \neq 1$. Lemma 5 from Sonthalia and Nadakuditi (2023) implies that if \mathbf{A} has unit variance entries, the moments of its inverse eigenvalue are expressions of c and are hence $O(1)$. In our case, it will just scale with ρ instead:

$$\mathbb{E}[1/\sigma_i^4] = O(1/\rho^4), \quad \mathbb{E}[1/(\sigma_i^2 \sigma_k^2)] = O(1/\rho^4), \quad \text{and} \quad \mathbb{E}[1/\sigma_i^8] = O(1/\rho^8) \quad \text{etc.}$$

In particular, we also need the following exact expectation from the same lemma:

$$\mathbb{E}\left[\frac{1}{\sigma_i^2}\right] = \frac{c}{\rho^2|1-c|} + o\left(\frac{1}{\rho^2}\right) = O\left(\frac{1}{\rho^2}\right). \quad (8)$$

Since the above I_1, I_3 have $r = \min(d, n)$ summands, this implies

$$I_1 = O\left(\frac{r}{N^4 \rho^4}\right) = O\left(\frac{1}{N^3 \rho^4}\right), \quad I_3 = O\left(\frac{r}{\eta^2 N^2 \rho^2}\right) = O\left(\frac{1}{N \rho^4}\right).$$

Similarly, I_2 has $r(r-1) \approx r^2$ summands, and

$$I_2 = O\left(\frac{r^2}{N^4 \rho^4}\right) = O\left(\frac{1}{N^2 \rho^4}\right)$$

$$\implies \mathbb{E}[\zeta^4] = \frac{1}{\eta^4} + I_1 + I_2 + I_3 = \frac{1}{\eta^4} + O\left(\frac{1}{\max(d, n) \rho^4}\right) \quad \text{since } I_3 \text{ dominates.} \quad (9)$$

With a similar expansion for the second moment and taking spherical moments, we get that

$$\begin{aligned} \mathbb{E}[\zeta^2] &= \frac{1}{\eta^2} + \sum_{i,j} \mathbb{E}\left[\frac{1}{\sigma_i \sigma_j}\right] \mathbb{E}[b_i b_j] \mathbb{E}[a_i a_j] = \frac{1}{\eta^2} + \frac{\sum_{i=1}^r \mathbb{E}[1/\sigma_i^2]}{dn} \\ &= \frac{1}{\eta^2} + \frac{\min(d, n)}{dn} \left(\frac{c}{\rho^2|1-c|} + o\left(\frac{1}{\rho^2}\right) \right) \quad \text{by Equation 8} \\ &= \frac{1}{\eta^2} + \frac{1}{\max(d, n)} \frac{c}{\rho^2|1-c|} + o\left(\frac{1}{\max(d, n) \rho^2}\right). \end{aligned}$$

This gives us the mean. Furthermore,

$$(\mathbb{E}[\zeta^2])^2 = \frac{1}{\eta^4} + \frac{2}{\eta^2} \frac{\sum_{i=1}^r \mathbb{E}[1/\sigma_i^2]}{dn} + \frac{(\sum_{i=1}^r \mathbb{E}[1/\sigma_i^2])^2}{d^2 n^2} = \frac{1}{\eta^4} + O\left(\frac{1}{\max(d, n) \rho^4}\right). \quad (10)$$

Variance. $\text{Var}(\zeta^2) = \mathbb{E}[\zeta^4] - (\mathbb{E}[\zeta^2])^2$. From Equations 9, 10, the overall scaling is determined by the dominant term:

$$\text{Var}\left(\left(\frac{\zeta}{\eta}\right)^2\right) = O\left(\frac{1}{\max(d, n) \rho^4}\right).$$

□

Lemma 10 (General Terms). *In the setting of Section 2 we have the following expectations:*

1. For $c < 1$, $\mathbb{E}[\beta_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \beta_*] = \frac{c}{\rho^2(1-c)} (\beta_*^\top \mathbf{u})^2 + o\left(\frac{1}{\rho^2}\right)$ and the variance is $O(1/(\rho^4 d))$.
2. For $c < 1$, $\mathbb{E}[\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k}] = \frac{c^2}{\rho^4(1-c)^3} + o\left(\frac{1}{\rho^4}\right)$ and the variance is $O(1/(\rho^8 d))$.
3. For $c > 1$, $\mathbb{E}[\beta_*^\top \mathbf{s} \mathbf{u}^\top \beta_*] = \frac{c-1}{c} (\beta_*^\top \mathbf{u})^2$ and the variance is $O(1/d)$.
4. For $c > 1$, $\mathbb{E}[\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top \beta_*] = \frac{c-1}{c^2} (\beta_*^\top \mathbf{u})^2 + o(1)$ and the variance is $O(1/d)$.
5. For $c > 1$, $\mathbb{E}[\beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_*] = \frac{\|\beta_*\|^2}{d} \frac{c}{\rho^2(c-1)} + o\left(\frac{1}{\rho^2 d}\right)$ and the variance is $O(1/(\rho^4 d^2))$.
6. For $c > 1$, $\mathbb{E}[\mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top] = \frac{1}{\rho^4} \frac{c^3}{(c-1)^3} + o\left(\frac{1}{\rho^4}\right)$ and the variance is $O(1/(\rho^8 d))$.

7. For $c > 1$, $\mathbb{E}[\|\mathbf{k}\|^2] = \frac{1}{\rho^2} \frac{1}{c-1} + o\left(\frac{1}{\rho^2}\right)$ and the variance is $O(1/(\rho^4 n))$

Proof. For all these terms, we evaluate the expectation using the SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, with $\mathbf{A}^\dagger = \mathbf{V}\mathbf{\Sigma}^\dagger\mathbf{U}^\top$, and important expectations from Lemma 5 of Sonthalia and Nadakuditi (2023) regarding the spectrum of \mathbf{A} : suppose $\tilde{\mathbf{A}}$ has unit variance (general ρ^2 is a multiplicative change), and let $\sigma_i(\tilde{\mathbf{A}})$ denote the i -th singular value. We have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{\sigma_i^2(\tilde{\mathbf{A}})}\right] &= \begin{cases} \frac{c}{1-c} + o(1) & c < 1 \\ \frac{c}{c-1} + o(1) & c > 1 \end{cases}, & \mathbb{E}\left[\frac{1}{\sigma_i^4(\tilde{\mathbf{A}})}\right] &= \begin{cases} \frac{c^2}{(1-c)^3} + o(1) & c < 1 \\ \frac{c^3}{(c-1)^3} + o(1) & c > 1 \end{cases}. \\ \mathbb{E}\left[\frac{1}{\sigma_i^2(\mathbf{A})}\right] &= \begin{cases} \frac{1}{\rho^2} \frac{c}{1-c} + o\left(\frac{1}{\rho^2}\right) & c < 1 \\ \frac{1}{\rho^2} \frac{c}{c-1} + o\left(\frac{1}{\rho^2}\right) & c > 1 \end{cases}, & \mathbb{E}\left[\frac{1}{\sigma_i^4(\mathbf{A})}\right] &= \begin{cases} \frac{1}{\rho^4} \frac{c^2}{(1-c)^3} + o\left(\frac{1}{\rho^4}\right) & c < 1 \\ \frac{1}{\rho^4} \frac{c^3}{(c-1)^3} + o\left(\frac{1}{\rho^4}\right) & c > 1 \end{cases}. \end{aligned} \quad (11)$$

For the first term, we note that

$$\begin{aligned} \beta_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \beta_* &= (\beta_*^\top \mathbf{u}) \mathbf{u}^\top \mathbf{A}^{\dagger\top} \mathbf{A}^\dagger \beta_* \\ &= (\beta_*^\top \mathbf{u}) \mathbf{u}^\top \mathbf{U} \mathbf{\Sigma}^\dagger \mathbf{\Sigma}^\dagger \mathbf{U}^\top \beta_* \\ &= (\beta_*^\top \mathbf{u}) \sum_{i=1}^d (\mathbf{u}^\top \mathbf{U})_i (\mathbf{U}^\top \beta_*)_i \frac{1}{\sigma_i^2(\mathbf{A})} \\ &= (\beta_*^\top \mathbf{u}) \sum_{i=1}^d (\mathbf{u}^\top \mathbf{u}_i) (\beta_*^\top \mathbf{u}_i) \frac{1}{\sigma_i^2(\mathbf{A})}, \end{aligned}$$

where \mathbf{u}_i denotes the i -th column of \mathbf{U} . We further note that $\mathbf{u}^\top \beta_* = \mathbf{u}^\top \mathbf{U} \mathbf{U}^\top \beta_*$. Since permuting columns of an orthogonal matrix does not break orthogonality and \mathbf{U} is uniformly random, we have that the marginals \mathbf{u}_i are identical. Thus, we have that

$$\mathbb{E}[\mathbf{u}^\top \mathbf{u}_1 \beta_*^\top \mathbf{u}_1] = \dots = \mathbb{E}[\mathbf{u}^\top \mathbf{u}_d \beta_*^\top \mathbf{u}_d] = \frac{1}{d} (\mathbf{u}^\top \beta_*) \quad \text{since } \mathbb{E}[\mathbf{u}_i \mathbf{u}_i^\top] = \frac{1}{d} \mathbf{I}.$$

It follows from here that

$$\begin{aligned} \mathbb{E}[\beta_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \beta_*] &= (\beta_*^\top \mathbf{u}) \sum_{i=1}^d \mathbb{E}[\mathbf{u}^\top \mathbf{u}_i \beta_*^\top \mathbf{u}_i] \mathbb{E}\left[\frac{1}{\sigma_i^2(\mathbf{A})}\right] \\ &= \frac{1}{\rho^2} (\beta_*^\top \mathbf{u})^2 \sum_{i=1}^d \frac{1}{d} \left(\frac{c}{1-c} + o(1)\right) \quad \text{by Equation 11} \\ &= \frac{1}{\rho^2} \frac{c}{1-c} (\beta_*^\top \mathbf{u})^2 + o\left(\frac{1}{\rho^2}\right). \end{aligned}$$

Since \mathbf{A} is isotropic Gaussian, we have that \mathbf{U}, \mathbf{V} are uniformly random orthogonal matrices. Thus, $\mathbf{u}^\top \mathbf{U}$ and $\mathbf{U}^\top \beta_*$ are uniformly random vectors on the spheres of radius $\|\mathbf{u}\|$ and $\|\beta_*\|$ respectively.

Hence, when we consider the squared terms to compute the variance, the term from the two uniform vectors will contribute $O(1/d^2)$. Together with the singular value term (now squared to have $O(1/\rho^4)$) and the summation, the variance is of order $O(1/(\rho^4 d))$.

For the second term, we have that by Equation 11,

$$\begin{aligned} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} &= \mathbf{u}^\top ((\mathbf{A}\mathbf{A}^\top)^\dagger)^2 \mathbf{u} = \mathbf{u}^\top \mathbf{U} ((\mathbf{\Sigma}\mathbf{\Sigma}^\top)^\dagger)^2 \mathbf{U}^\top \mathbf{u} = \sum_{i=1}^d (\mathbf{u}^\top \mathbf{u}_i)^2 \frac{1}{\sigma_i^4(\mathbf{A})}, \\ \mathbb{E}[\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}] &= \sum_{i=1}^d \mathbb{E}[(\mathbf{u}^\top \mathbf{u}_i)^2] \mathbb{E}\left[\frac{1}{\sigma_i^4(\mathbf{A})}\right] = \sum_{i=1}^d \frac{1}{\rho^4} \frac{1}{d} \left(\frac{c^2}{(1-c)^3} + o(1)\right) = \frac{1}{\rho^4} \frac{c^2}{(1-c)^3} + o\left(\frac{1}{\rho^4}\right), \end{aligned}$$

where we again use $\mathbb{E}[(\mathbf{u}^\top \mathbf{u}_i)^2] = 1/d$ since it is the entry of a uniformly random vector of length $\|\mathbf{u}\| = 1$.

Similarly, the variance is $O(1/(\rho^8 d))$ from the summation of d independent variances each of $O(1/(\rho^8 d^2))$.

For the third term, we have that

$$\beta_*^\top \mathbf{s} \mathbf{u}^\top \beta_* = \beta_*^\top (\mathbf{I} - \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} (\mathbf{u}^\top \beta_*) = (\beta_*^\top \mathbf{u})^2 - (\beta_*^\top \mathbf{u}) \sum_{i=1}^n (\beta_*^\top \mathbf{u}_i) (\mathbf{u}^\top \mathbf{u}_i).$$

Similarly, we take the expectation (in particular, $\mathbb{E}[(\beta_*^\top \mathbf{u}_i) (\mathbf{u}^\top \mathbf{u}_i)] = 1/d (\beta_*^\top \mathbf{u})$) and have

$$(\beta_*^\top \mathbf{u})^2 \left[1 - \sum_{i=1}^n \frac{1}{d} \right] = \left(1 - \frac{1}{c} \right) (\beta_*^\top \mathbf{u})^2.$$

The variance for this term is $O(1/d)$ from summation of $n = d/c$ terms of $O(1/d^2)$.

For the fourth term, we plug in $\mathbf{s} = (\mathbf{I} - \mathbf{A} \mathbf{A}^\dagger) \mathbf{u}$ and have

$$\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top \beta_* = (\beta_*^\top \mathbf{u}) \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} - (\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u})^2.$$

From previous calculations, we have that

$$\mathbb{E}[\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u}] = \mathbb{E} \left[\sum_{i=1}^n (\beta_*^\top \mathbf{u}_i) (\mathbf{u}^\top \mathbf{u}_i) \right] = \frac{1}{c} (\beta_*^\top \mathbf{u}).$$

Using Proposition 2 and this result, we can then show

$$\mathbb{E}[(\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u})^2] = \frac{1}{c^2} (\beta_*^\top \mathbf{u})^2 + o(1).$$

It follows that

$$\mathbb{E}[\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top \beta_*] = \frac{c-1}{c^2} (\beta_*^\top \mathbf{u})^2 + o(1).$$

The variance for this term is $O(1/d)$, where the dominant term is a summation of $n = d/c$ terms of $O(1/d^2)$.

For the fifth term, we have

$$\beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* = (\beta_*^\top \mathbf{A}^\dagger \mathbf{v})^2 = \sum_{i,j} (\beta_*^\top \mathbf{U})_i (\beta_*^\top \mathbf{U})_j \frac{1}{\sigma_i(\mathbf{A}) \sigma_j(\mathbf{A})} (\mathbf{V}^\top \mathbf{v})_i (\mathbf{V}^\top \mathbf{v})_j.$$

Since $\beta_*^\top \mathbf{U}$ (and $\mathbf{V}^\top \mathbf{v}$) are uniformly random and independent of everything else, we only have the diagonal terms when we take the expectation. By Equation 11,

$$\mathbb{E}[\beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_*] = \sum_{i=1}^n \frac{\|\beta_*\|^2}{d} \frac{1}{n} \frac{1}{\rho^2} \left(\frac{c}{c-1} + o(1) \right) = \frac{\|\beta_*\|^2}{d} \frac{1}{\rho^2} \frac{c}{c-1} + o\left(\frac{1}{\rho^2 d}\right)$$

The variance for this term is $O(1/(\rho^4 d^2))$ from $O(d^2)$ terms of individual variances of $O(1/(\rho^4 d^4))$.

For the sixth term, by expansion and Equation 11, similar to above,

$$\mathbb{E}[\mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top] = \sum_{i=1}^n \mathbb{E}[(\mathbf{V}^\top \mathbf{v})_i^2] \mathbb{E} \left[\frac{1}{\sigma_i^4(\mathbf{A})} \right] = \sum_{i=1}^n \frac{1}{n} \mathbb{E} \left[\frac{1}{\sigma_i^4(\mathbf{A})} \right] = \frac{1}{\rho^4} \frac{c^3}{(c-1)^3} + o\left(\frac{1}{\rho^4}\right).$$

The variance is $O(1/(\rho^8 d))$.

For the final term, by expansion and Equation 11,

$$\mathbb{E}[\|\mathbf{k}\|^2] = \sum_{i=1}^n \mathbb{E}[(\mathbf{u}^\top \mathbf{U})_i^2] \mathbb{E} \left[\frac{1}{\sigma_i^2(\mathbf{A})} \right] = \frac{1}{\rho^2} \frac{n}{d} \frac{c}{c-1} + o\left(\frac{1}{\rho^2}\right) = \frac{1}{\rho^2} \frac{1}{c-1} + o\left(\frac{1}{\rho^2}\right)$$

The variance is $O(1/(\rho^4 n))$. \square

Lemma 11 (Zero Expectation). *In the setting of Section 2, we have the following expectations for*

1. $\forall c, \mathbb{E}[\beta_*^\top \mathbf{u} \mathbf{h} \beta_*] = 0$ and $\text{Var}(\beta_*^\top \mathbf{u} \mathbf{h} \beta_*) = O(1/(\rho^2 d))$
2. If $c > 1$, $\mathbb{E}[\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \beta_*] = 0$ and $\text{Var}(\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \beta_*) = O(1/(\rho^2 d^2))$
3. If $c > 1$, $\mathbb{E}[\beta_*^\top \mathbf{s} \mathbf{h} \beta_*] = 0$ and $\text{Var}(\beta_*^\top \mathbf{s} \mathbf{h} \beta_*) = O(1/(\rho^2 d))$
4. $\forall c, \mathbb{E}[\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top] = 0$ and $\text{Var}(\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top) = O(1/(\rho^6 d))$
5. If $c > 1$, $\mathbb{E}[\mathbf{h} \mathbf{A} \mathbf{A}^\dagger \beta_*] = 0$ and $\text{Var}(\mathbf{h} \mathbf{A} \mathbf{A}^\dagger \beta_*) = O(1/(\rho^2 d))$

Proof. Similar to Lemma 10, for all these terms, we evaluate the expectation using the SVD $\mathbf{A} = \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^\top$, with $\mathbf{A}^\dagger = \mathbf{V} \boldsymbol{\Sigma}^\dagger \mathbf{U}^\top$.

For the first term, we note that

$$\beta_*^\top \mathbf{u} \mathbf{h} \beta_*^\top = (\beta_*^\top \mathbf{u}) \mathbf{v}^\top \mathbf{A}^\dagger \beta_* = (\beta_*^\top \mathbf{u}) \mathbf{v}^\top \mathbf{V} \boldsymbol{\Sigma}^\dagger \mathbf{U}^\top \beta_* = (\beta_*^\top \mathbf{u}) \sum_{i=1}^{\min(n,d)} (\mathbf{v}^\top \mathbf{V})_i (\mathbf{U}^\top \beta_*)_i \frac{1}{\sigma_i(\mathbf{A})}.$$

Since \mathbf{A} is isotropic Gaussian, again we have that \mathbf{U}, \mathbf{V} are uniformly random orthogonal matrices. Thus, $\mathbf{v}^\top \mathbf{V}$ and $\mathbf{U}^\top \beta_*$ are uniformly random vectors on a spheres of radius $\|\mathbf{v}\|$ and $\|\beta_*\|$ respectively. In particular, they are independent and have mean zero, which implies

$$\mathbb{E}[\beta_*^\top \mathbf{u} \mathbf{h} \beta_*^\top] = 0.$$

The variance will be $O(1/(\rho^2 d))$ as a summation of $O(d)$ terms of $O(1/(\rho^2 d^2))$.

For the second term, we note that

$$\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} = \sum_{i=1}^{\min(n,d)} (\beta_*^\top \mathbf{U})_i (\mathbf{U}^\top \mathbf{u})_i \quad \text{and} \quad \mathbf{h} \beta_* = \sum_{i=1}^{\min(n,d)} (\mathbf{v}^\top \mathbf{V})_i (\mathbf{U}^\top \beta_*)_i \frac{1}{\sigma_i(\mathbf{A})}$$

Multiplying the two together yields

$$\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \beta_* = \sum_{i,j}^{\min(n,d)} (\beta_*^\top \mathbf{U})_i (\mathbf{U}^\top \mathbf{u})_i (\mathbf{v}^\top \mathbf{V})_j (\mathbf{U}^\top \beta_*)_j \frac{1}{\sigma_i(\mathbf{A})}.$$

We note that $\mathbf{v}^\top \mathbf{V}$ is a uniformly random mean zero vector independent of everything else in the summation. Hence, the expectation is equal to zero, and similar to Lemma ??, the variance of this term is $O(1/(\rho^2 d^2))$ (a summation of $O(d^2)$ terms of $O(1/(\rho^2 d^4))$).

For the third term, we have that

$$\beta_*^\top \mathbf{s} \mathbf{h} \beta_* = \beta_*^\top (\mathbf{I} - \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} \mathbf{h} \beta_* = \beta_*^\top \mathbf{u} \mathbf{h} \beta_* - \beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \beta_*.$$

Then using the previous two parts, we get that each term has mean zero. Thus, we get the needed result. Using Lemma 34 and the first two terms, the variance of this term is $O(1/(\rho^2 d))$.

For the fourth term, we have that:

$$\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top = \mathbf{u} \mathbf{U} \boldsymbol{\Sigma}^\dagger \mathbf{V}^\top \boldsymbol{\Sigma}^\dagger \boldsymbol{\Sigma}^\dagger \mathbf{V}^\top \mathbf{v} = \sum_{i=1}^{\min(n,d)} (\mathbf{u}^\top \mathbf{U})_i (\mathbf{V}^\top \mathbf{v})_i \frac{1}{\sigma_i(\mathbf{A})^3}.$$

Similarly, using the independence of $\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{V}$ and uniformly random entries, we get mean zero and variance $O(1/(\rho^6 d))$.

For the last term, we have that:

$$\mathbf{h} \mathbf{A} \mathbf{A}^\dagger \beta_* = \sum_{i=\min(n,d)}^r (\mathbf{V}^\top \mathbf{v})_i (\mathbf{U}^\top \beta_*)_i \frac{1}{\sigma_i(\mathbf{A})}.$$

Using the independence of $\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{V}$ and uniformly random entries, we get mean zero and variance $O(1/(\rho^2 d))$. \square

E.4.2 STEP 3(B): BOUNDING THE HIGHER MOMENTS

To bound the higher moments, we will use the following Gaussian hypercontractivity lemma.

Lemma 12 (Gaussian Hypercontractivity Inequality). *Let $G \sim \mathcal{N}(0, 1)$ be a standard Gaussian random variable. Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be a degree k polynomial. Then, for any $q \geq 2$, the L_q norm of $f(G)$ is bounded by its L_2 norm as follows:*

$$\|f(G)\|_{L_q} \leq (q-1)^{k/2} \|f(G)\|_{L_2},$$

where the L_p norm of a random variable X is defined as $\|X\|_{L_p} = (\mathbb{E}[|X|^p])^{1/p}$.

Proof. Follows directly from (Mei et al., 2022, Lemma 20). \square

Lemma 13 (Multivariate Gaussian Hypercontractivity). *Let $G = (G_1, \dots, G_M) \sim \mathcal{N}(0, I_M)$ and let $P : \mathbb{R}^M \rightarrow \mathbb{R}$ be a polynomial of total degree r . Consider the Hermite expansion of P*

$$P(x) = \sum_{\alpha \in \mathbb{N}^m, |\alpha| \leq r} c_\alpha \mathbf{H}_\alpha(x).$$

with coefficient random and independent of G . Then there exists a constant C that is only dependent on M, r such that for any $q \geq 2$,

$$\|P(G)\|_{L_q} \leq C(q-1)^{r/2} \left(\sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q}^2 \alpha! \right)^{1/2}$$

Further, if for all $|\alpha| \leq r$, we have that $\|c_\alpha\|_{L_q}^2 \leq C_q^2 \|c_\alpha\|_{L_2}^2$, then

$$\|P(G)\|_{L_q} \leq C(q-1)^{r/2} \|P(G)\|_{L_2}$$

Where the L_p norm is over all of the randomness. Furthermore,

Proof. Let $H_k : \mathbb{R} \rightarrow \mathbb{R}$ be the probabilist Hermite polynomial. Given $\alpha \in \mathbb{N}^M$, define

$$\mathbf{H}_\alpha(x) := \prod_{j=1}^M H_{\alpha_j}(x_j)$$

Then since P is degree r , then we can decompose

$$P(x) = \sum_{\alpha \in \mathbb{N}^m, |\alpha| \leq r} c_\alpha \mathbf{H}_\alpha(x).$$

Here $|\alpha| = \sum_j \alpha_j$. Since the Hermite polynomials are orthogonal, we can see that

$$\int_{\mathbb{R}^M} \mathbf{H}_\alpha(x) \mathbf{H}_{\bar{\alpha}}(x) \gamma_M(x) dx = \delta_{\alpha\bar{\alpha}} \prod_{j=1}^M \alpha_j!,$$

where γ_M is the density for an M -dimensional standard normal distribution.

$$\begin{aligned} \|P(x)\|_{L_2}^2 &= \mathbb{E}_\Sigma \left[\int_{\mathbb{R}^M} |P(x)|^2 \gamma_M(x) dx \right] \\ &= \sum_{|\alpha| \leq r} \sum_{|\bar{\alpha}| \leq r} \mathbb{E}_\Sigma [c_\alpha c_{\bar{\alpha}}] \int \mathbf{H}_\alpha(x) \mathbf{H}_{\bar{\alpha}}(x) \gamma_M(x) dx \\ &= \sum_{|\alpha| \leq r} \|c_\alpha\|_{L_2}^2 \alpha! \end{aligned}$$

where $\alpha! := \prod_{j=1}^M \alpha_j!$.

Then using the 1D Gaussian Hypercontractivity (Lemma 12, we see that

$$\begin{aligned}\|\mathbf{H}_\alpha(x)\|_{L_q} &= \prod_{j=1}^M \|H_{\alpha_j}(x_j)\|_{L_q} \\ &\leq \prod_{j=1}^M (q-1)^{\alpha_j/2} \|H_{\alpha_j}(x_j)\|_{L_2} \\ &= (q-1)^{|\alpha|/2} \prod_{j=1}^M \sqrt{\alpha_j!} \\ &= (q-1)^{|\alpha|/2} \sqrt{\alpha!}\end{aligned}$$

Thus, using the triangle inequality we get that

$$\|P(x)\|_{L_q} \leq \sum_{|\alpha| \leq r} \|c_\alpha \mathbf{H}_\alpha(x)\|_{L_q} = \sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q} \|\mathbf{H}_\alpha(x)\|_{L_q}$$

Thus

$$\|P(x)\|_{L_q} \leq \sum_{|\alpha| \leq r} \|c_\alpha \mathbf{H}_\alpha(x)\|_{L_q} \leq \sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q} (q-1)^{|\alpha|/2} \sqrt{\alpha!} \leq (q-1)^{r/2} \sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q} \sqrt{\alpha!}$$

Then using Cauchy-Schwartz, we get that

$$\sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q} \sqrt{\alpha!} \leq \left(\sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q}^2 \alpha! \right)^{1/2} \left(\sum_{|\alpha| \leq r} 1 \right)^{1/2}.$$

Finally, we note that

$$C_{M,r} := \left(\sum_{|\alpha| \leq r} 1 \right)^{1/2}$$

is some universal constant that only depends on M, r . Thus, we get that

$$\|P(x)\|_{L_q} \leq C_{M,r} (q-1)^{r/2} \left(\sum_{|\alpha| \leq r} \|c_\alpha\|_{L_q}^2 \alpha! \right)^{1/2}$$

Using the assumption

$$\|c_\alpha\|_{L_q}^2 \leq C_q^2 \|c_\alpha\|_{L_2}^2$$

Then we get

$$\|P(x)\|_{L_q} \leq C_{M,r} C_q (q-1)^{r/2} \|P(x)\|_{L_2}$$

□

Lemma 14 (Product Spherical Hypercontractivity). *Let $l_1, l_2, l_3 \geq 0$, let $\Theta_1 \sim \text{Unif}(S^{l_1})$, $\Theta_2 \sim \text{Unif}(S^{l_2})$, $\Theta_3 \sim \text{Unif}(S^{l_3})$ be independent, and let $H : \mathbb{R}^{l_1+1} \times \mathbb{R}^{l_2+1} \times \mathbb{R}^{l_3+1} \rightarrow \mathbb{R}$ be a multi-homogeneous polynomial of total degree r . Then for every $q \geq 2$,*

$$\|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_q} \leq C_{r,q} (q-1)^{r/2} \|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_2},$$

where the norms are with respect to the product measure. For homogeneous polynomials, the constant is independent of the dimension.

Proof. H is multi-homogeneous of degrees r_1, r_2, r_3 with $r_1 + r_2 + r_3 = r$. Let $G_1 \sim \mathcal{N}(0, I_{l_1+1})$, $G_2 \sim \mathcal{N}(0, I_{l_2+1})$, $G_3 \sim \mathcal{N}(0, I_{l_3+1})$ be independent with polar decompositions $G_i = R_i \Theta_i$, where the R_i 's are independent of each other and of the Θ_i 's. Then

$$H(G_1, G_2, G_3) = R_1^{r_1} R_2^{r_2} R_3^{r_3} H(\Theta_1, \Theta_2, \Theta_3),$$

so for any $p > 0$,

$$\mathbb{E} [|H(G_1, G_2, G_3)|^p] = \left(\prod_{i=1}^3 \mathbb{E} [R_i^{pr_i}] \right) \mathbb{E} [|H(\Theta_1, \Theta_2, \Theta_3)|^p]$$

Then we have that

$$\|H(G_1, G_2, G_3)\|_{L_p} = \left(\prod_i (\mathbb{E} [R_i^{pr_i}])^{1/p} \right) \|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_p}. \quad (12)$$

Apply Gaussian hypercontractivity (Lemma 12) to $H(G_1, G_2, G_3)$ (total degree r):

$$\|H(G_1, G_2, G_3)\|_{L_q} \leq C(q-1)^{r/2} \|H(G_1, G_2, G_3)\|_{L_2}, \quad q \geq 2.$$

Using Equation 12 with $p = q$ and $p = 2$ yields

$$\|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_q} \leq C(q-1)^{r/2} \left(\prod_i \frac{(\mathbb{E} [R_i^{2r_i}])^{1/2}}{(\mathbb{E} [R_i^{qr_i}])^{1/q}} \right) \|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_2}.$$

For each i , since $q \geq 2$ and $R_i \geq 0$, monotonicity of L_p norms implies $(\mathbb{E} [R_i^{qr_i}])^{1/(qr_i)} \geq (\mathbb{E} [R_i^{2r_i}])^{1/(2r_i)}$, hence

$$\frac{(\mathbb{E} [R_i^{2r_i}])^{1/2}}{(\mathbb{E} [R_i^{qr_i}])^{1/q}} \leq 1.$$

Thus the product is less than 1, so

$$\|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_q} \leq C(q-1)^{r/2} \|H(\Theta_1, \Theta_2, \Theta_3)\|_{L_2}. \quad \square$$

Lemma 15 (Product spherical hypercontractivity with random coefficients). *Let $l_1, l_2, l_3 \geq 0$ and let $\Theta_i \sim \text{Unif}(S^{l_i})$ be independent. Let $r \in \mathbb{N}$ and let $H : \mathbb{R}^{l_1+1} \times \mathbb{R}^{l_2+1} \times \mathbb{R}^{l_3+1} \rightarrow \mathbb{R}$ be a multi-homogeneous polynomial of total degree at most r . Suppose the coefficients of P are random on an auxiliary probability space and are independent of $(\Theta_1, \Theta_2, \Theta_3)$. If the random coefficients satisfy $\|c_\alpha\|_{L_q} \leq K_q \|c_\alpha\|_{L_2}$ in the Hermite basis expansion, then for all $q \geq 2$:*

$$\|H\|_{L_q} \leq C_{r,q} (q-1)^{r/2} \|H\|_{L_2}.$$

Proof. The proof is identical to that of Lemma 14, except we begin with the version of Gaussian hypercontractivity that handles random coefficients satisfying the stated assumption. \square

Recall

$$\mathbf{a} := \mathbf{V}^\top \mathbf{v} \in \mathbb{R}^n \quad \mathbf{b} := \mathbf{U}^\top \mathbf{u} \in \mathbb{R}^d, \quad \text{and} \quad \mathbf{u}_\beta = \mathbf{U}^\top \boldsymbol{\beta}_*$$

Then, since \mathbf{u}, \mathbf{u} are fixed, and \mathbf{U}, \mathbf{V} are independent Haar orthogonal matrices, we have that \mathbf{a}, \mathbf{b} are all uniformly random vectors on their respective spheres. Additionally, using the assumption that $\boldsymbol{\beta}_*$ is uniformly random such that $\boldsymbol{\beta}_*^\top \mathbf{u}$ is constant. \mathbf{u}_β is uniformly random on a sphere \mathbb{S}^{d-2} .

Consider the following centered versions and polynomial representations.

1. $Y_h := \|\mathbf{h}\|^2 - \mathbb{E} [\|\mathbf{h}\|^2] = \mathbf{a}^\top (\boldsymbol{\Sigma}^\dagger \boldsymbol{\Sigma}^{\dagger\top} - \mu_h) \mathbf{a}$
2. $Y_k := \|\mathbf{k}\|^2 - \mathbb{E} [\|\mathbf{k}\|^2] = \mathbf{b}^\top (\boldsymbol{\Sigma}^{\dagger\top} \boldsymbol{\Sigma}^\dagger - \mu_k) \mathbf{b}$
3. $Y_t := \|\mathbf{t}\|^2 - \mathbb{E} [\|\mathbf{t}\|^2] = \mathbf{a}^\top ((I - \boldsymbol{\Sigma}^\dagger \boldsymbol{\Sigma}) - \mu_t)$
4. $Y_s := \|\mathbf{s}\|^2 - \mathbb{E} [\|\mathbf{s}\|^2] = \mathbf{b}^\top ((I - \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger) - \mu_t) \mathbf{b}$
5. $Y_\xi := \frac{\xi}{\eta} - \mathbb{E} \left[\frac{\xi}{\eta} \right] = \mathbf{a}^\top \boldsymbol{\Sigma} \mathbf{b} = \mathbf{a}^\top \boldsymbol{\Sigma}^\dagger \mathbf{b}$
6. $\tilde{T}_1 := \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \boldsymbol{\beta}_*] = (\boldsymbol{\beta}_*^\top \mathbf{u}) \mathbf{b}^\top (\boldsymbol{\Sigma}^{\dagger\top} \boldsymbol{\Sigma}^\dagger) \mathbf{u}_\beta - \mu_{\tilde{T}_1} (\mathbf{b}^\top \mathbf{b})$

7. $\tilde{T}_2 := \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} - \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}] = \mathbf{b}^\top \left((\boldsymbol{\Sigma}^{\dagger\top} \boldsymbol{\Sigma}^\dagger)^2 - \mu_{\tilde{T}_2} \right) \mathbf{b}$
8. $\tilde{T}_3 := \boldsymbol{\beta}_*^\top \mathbf{s} \mathbf{u}^\top \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{s} \mathbf{u}^\top \boldsymbol{\beta}_*] = (\boldsymbol{\beta}_*^\top \mathbf{u}) \mathbf{u}_\beta^\top (I - \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger) \mathbf{b} - \mu_{\tilde{T}_3} (\mathbf{u}_\beta^\top \mathbf{u}_\beta)$
9. $\tilde{T}_4 := \boldsymbol{\beta}_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top \boldsymbol{\beta}_*] = \mathbf{u}_\beta^\top \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger \mathbf{b} \mathbf{b}^\top (I - \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger) \mathbf{u}_\beta - \mu_{\tilde{T}_4} (\mathbf{b}^\top \mathbf{b}) (\mathbf{u}_\beta^\top \mathbf{u}_\beta)$
10. $\tilde{T}_5 := \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_*] = (\mathbf{u}_\beta \boldsymbol{\Sigma}^{\dagger\top} \mathbf{a})^2 - \mu_{\tilde{T}_5} (\mathbf{a}^\top \mathbf{a}) (\mathbf{u}_\beta^\top \mathbf{u}_\beta)$
11. $\tilde{T}_6 := \mathbf{h} (\mathbf{A}^\dagger)^\top \mathbf{A}^\dagger \mathbf{h}^\top - \mathbb{E} [\mathbf{h} (\mathbf{A}^\dagger)^\top \mathbf{A}^\dagger \mathbf{h}^\top] = \mathbf{a}^\top \left((\boldsymbol{\Sigma}^\dagger \boldsymbol{\Sigma}^{\dagger\top})^2 - \mu_{\tilde{T}_6} \right) \mathbf{a}$
12. $\tilde{S}_1 := \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{h} \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{h} \boldsymbol{\beta}_*] = (\boldsymbol{\beta}_*^\top \mathbf{u}) \mathbf{a}^\top \boldsymbol{\Sigma}^\dagger \mathbf{u}_\beta$
13. $\tilde{S}_2 := \boldsymbol{\beta}_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \boldsymbol{\beta}_*] = \mathbf{u}_\beta^\top \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger \mathbf{b} \mathbf{a}^\top \boldsymbol{\Sigma}^\dagger \mathbf{u}_\beta$
14. $\tilde{S}_3 := \boldsymbol{\beta}_*^\top \mathbf{s} \mathbf{h} \boldsymbol{\beta}_* - \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{s} \mathbf{h} \boldsymbol{\beta}_*] = \mathbf{u}_\beta (I - \boldsymbol{\Sigma} \boldsymbol{\Sigma}^\dagger) \mathbf{b} \mathbf{a}^\top \boldsymbol{\Sigma}^\dagger \mathbf{u}_\beta$
15. $\tilde{S}_4 := \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top - \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top] = \mathbf{b}^\top \boldsymbol{\Sigma}^{\dagger\top} \boldsymbol{\Sigma}^\dagger \boldsymbol{\Sigma}^{\dagger\top} \mathbf{a}$

Hence we see that these are all homogeneous polynomials in uniformly random spherical variables. Thus, we can use Lemma 14, we get bounds on the higher moments. In particular, since the coefficients are only dependent on constants and $\boldsymbol{\Sigma}$, we see that the coefficients are independent of \mathbf{a} , \mathbf{b} , \mathbf{u}_β . Then using a change of basis we see that that coefficients of the decomposition are also random and independent of the input variables. Finally, since the spectrum converges to the Marchenko-Pastur, we have that the coefficients have bounded moments. Hence the second assumption is satisfied.

E.4.3 STEP 3(C): BOUNDING γ_i MOMENTS.

Lemma 16 (Moments of γ_i/η^2). *We have:*

(i) For γ_1/η^2 ,

$$\mathbb{E} \left[\frac{\gamma_1}{\eta^2} \right] = \frac{c}{\rho^2} + \frac{1}{\eta^2} + o\left(\frac{1}{\rho^2}\right), \quad \text{Var} \left(\frac{\gamma_1}{\eta^2} \right) = O\left(\frac{1}{\rho^4 n}\right).$$

(ii) For γ_2/η^2 ,

$$\mathbb{E} \left[\frac{\gamma_2}{\eta^2} \right] = \frac{1}{\rho^2} + \frac{1}{\eta^2} + o\left(\frac{1}{\rho^2}\right), \quad \text{Var} \left(\frac{\gamma_2}{\eta^2} \right) = O\left(\frac{1}{\rho^4 n}\right).$$

Proof. We decompose

$$\frac{\gamma_i}{\eta^2} = \zeta_i + \frac{\xi^2}{\eta^2}, \quad i = 1, 2, \quad \text{where } \zeta_1 = \|\mathbf{t}\|^2 \|\mathbf{k}\|^2, \quad \zeta_2 = \|\mathbf{s}\|^2 \|\mathbf{h}\|^2.$$

Expectation Estimates: We begin by noting that $\|\mathbf{t}\|^2$ depends only on \mathbf{V} and is independent of \mathbf{U} , $\boldsymbol{\Sigma}$. $\|\mathbf{s}\|^2$ depends only on \mathbf{U} and is independent of \mathbf{V} , $\boldsymbol{\Sigma}$. Additionally, $\|\mathbf{k}\|^2$ depends on \mathbf{U} and $\boldsymbol{\Sigma}$, hence is independent of \mathbf{V} . Also $\|\mathbf{h}\|^2$ depends on \mathbf{V} and $\boldsymbol{\Sigma}$ and is independent of \mathbf{U} , hence is independent of \mathbf{U} .

Thus, we have have that $\|\mathbf{t}\|^2$ and $\|\mathbf{k}\|^2$ are independent and $\|\mathbf{s}\|^2$ and $\|\mathbf{h}\|^2$ are independent. Thus, we see that

$$\mathbb{E}[\zeta_1] = \mathbb{E}[\|\mathbf{t}\|^2 \|\mathbf{k}\|^2] = \mathbb{E}[\|\mathbf{t}\|^2] \mathbb{E}[\|\mathbf{k}\|^2].$$

Using Lemma 9 again,

$$\mathbb{E}[\|\mathbf{t}\|^2] = 1 - c, \quad \mathbb{E}[\|\mathbf{k}\|^2] = \frac{1}{\rho^2} \frac{c}{1 - c} + o\left(\frac{1}{\rho^2}\right).$$

We plug them into the expectation and get:

$$\mathbb{E}[\zeta_1] = (1 - c) \left[\left(\frac{1}{\rho^2} \frac{c}{1 - c} \right) + o\left(\frac{1}{\rho^2}\right) \right] = \frac{c}{\rho^2} + o\left(\frac{1}{\rho^2}\right).$$

Finally, we also have that from Lemma 9,

$$\mathbb{E} \left[\frac{\xi^2}{\eta^2} \right] = \frac{1}{\eta^2} + O \left(\frac{1}{\rho^2 n} \right), \quad \text{Var} \left(\frac{\xi^2}{\eta^2} \right) = O \left(\frac{1}{\rho^4 n} \right),$$

Hence,

$$\mathbb{E} \left[\frac{\gamma_1}{\eta^2} \right] = \mathbb{E}[\zeta_1] + \mathbb{E} \left[\frac{\xi^2}{\eta^2} \right] = \frac{c}{\rho^2} + \frac{1}{\eta^2} + o \left(\frac{1}{\rho^2} \right).$$

A similar argument applies for γ_2/η^2 , using the corresponding results for $\|\mathbf{s}\|^2$, $\|\mathbf{h}\|^2$.

Variance Estimates:

Again using independence, we have that

$$\begin{aligned} \text{Var}(\|\mathbf{t}\|^2 \|\mathbf{k}\|^2) &= \text{Var}(\|\mathbf{t}\|^2) \text{Var}(\|\mathbf{k}\|^2) + \mathbb{E}[\|\mathbf{t}\|^2]^2 \text{Var}(\|\mathbf{k}\|^2) + \mathbb{E}[\|\mathbf{k}\|^2]^2 \text{Var}(\|\mathbf{t}\|^2) \\ &= O \left(\frac{1}{n} \right) O \left(\frac{1}{\rho^4 n} \right) + (1-c)^2 O \left(\frac{1}{\rho^4 n} \right) + \frac{1}{\rho^4} \frac{c^2}{(1-c)^2} O \left(\frac{1}{n} \right) \\ &= O \left(\frac{1}{\rho^4 n} \right). \end{aligned}$$

We then use Lemma 34 to compute the variance of the sum:

$$\begin{aligned} \text{Var} \left(\zeta_1 + \frac{\xi^2}{\eta^2} \right) &\leq \left(\sqrt{\text{Var}(\zeta_1)} + \sqrt{\text{Var} \left(\frac{\xi^2}{\eta^2} \right)} \right)^2 \\ &= \left(\sqrt{O \left(\frac{1}{\rho^4 n} \right)} + \sqrt{O \left(\frac{1}{\rho^4 n} \right)} \right)^2 \\ &= O \left(\frac{1}{\rho^4 n} \right). \end{aligned}$$

This proof is similar to the other case. □

Lemma 17 (Moments of $(\gamma_i/\eta^2)^2$). *We have, as $n, d \rightarrow \infty$ with $d/n \rightarrow c \neq 1$,*

(i) For γ_1/η^2 ,

$$\mathbb{E} \left[\left(\frac{\gamma_1}{\eta^2} \right)^2 \right] = \left(\frac{c}{\rho^2} + \frac{1}{\eta^2} \right)^2 + O \left(\frac{1}{\rho^4} \right), \quad \text{Var} \left(\left(\frac{\gamma_1}{\eta^2} \right)^2 \right) = O \left(\frac{1}{\rho^4 n} \right).$$

(ii) For γ_2/η^2 ,

$$\mathbb{E} \left[\left(\frac{\gamma_2}{\eta^2} \right)^2 \right] = \left(\frac{1}{\rho^2} + \frac{1}{\eta^2} \right)^2 + O \left(\frac{1}{\rho^4} \right), \quad \text{Var} \left(\left(\frac{\gamma_2}{\eta^2} \right)^2 \right) = O \left(\frac{1}{\rho^4 n} \right).$$

Proof. Write, for $i \in \{1, 2\}$,

$$\frac{\gamma_i}{\eta^2} = \zeta_i + \frac{\xi^2}{\eta^2}, \quad \zeta_1 := \|\mathbf{t}\|^2 \|\mathbf{k}\|^2, \quad \zeta_2 := \|\mathbf{s}\|^2 \|\mathbf{h}\|^2.$$

Means. Using Lemma 16 and the fact that for any random variable

$$\mathbb{E}[Y^2] = \mathbb{E}[Y]^2 + \text{Var}(Y)$$

we get the means.

Variances. Using

$$Y^2 = \mathbb{E}[Y]^2 + 2(\mathbb{E}[Y])(Y - \mathbb{E}[Y]) + (Y - \mathbb{E}[Y])^2,$$

Thus, using Lemma 34 we have that

$$\text{Var}(Y^2) \leq \left(\sqrt{4(\mathbb{E}[X])^2 \text{Var}(X_i)} + \sqrt{\text{Var}((Y - \mathbb{E}[Y])^2)} \right)^2.$$

By spherical hypercontractivity for degree-4 polynomials,

$$\mathbb{E} \left[\left(\frac{\gamma_i^2}{\eta^4} - \mathbb{E} \left[\frac{\gamma_i^2}{\eta^4} \right] \right)^4 \right] \lesssim \text{Var} \left(\frac{\gamma_i^2}{\eta^4} \right)^2,$$

hence

$$\text{Var} \left(\left(\frac{\gamma_i^2}{\eta^4} - \mathbb{E} \left[\frac{\gamma_i^2}{\eta^4} \right] \right)^2 \right) \mathbb{E} \left[\left(\frac{\gamma_i^2}{\eta^4} - \mathbb{E} \left[\frac{\gamma_i^2}{\eta^4} \right] \right)^4 \right] \lesssim \text{Var} \left(\frac{\gamma_i^2}{\eta^4} \right)^2.$$

Using $\mathbb{E} \left[\frac{\gamma_i^2}{\eta^2} \right]^2 = O(1)$ and $\text{Var} \left(\frac{\gamma_i^2}{\eta^2} \right) = O(\rho^{-4} n^{-1})$ gives

$$\text{Var} \left(\frac{\gamma_i^2}{\eta^4} \right) = O \left(\frac{1}{\rho^4 n} \right),$$

as claimed. \square

Lemma 18 (Finite Negative Moments of γ_i). *Fix $p > 0$. There exists an $N(p)$ such that for all $n, d \geq N(p)$, we have that for $c < 1$*

$$\mathbb{E} [\gamma_1^{-p}] \leq \eta^{-2p} \mathbb{E} [\sigma_1^{2p}] \mathbb{E} [T^{-p}] \leq \frac{\rho^{2p}}{\eta^{2p}} M^p$$

and for $c > 1$, we have that

$$\mathbb{E} [\gamma_2^{-p}] \leq \eta^{-2p} \mathbb{E} [\sigma_1^{2p}] \mathbb{E} [S^{-p}] \leq \frac{\rho^{2p}}{\eta^{2p}}, M^p$$

where σ_1 is the largest singular value of A , $T := \|\mathbf{t}\|^2 \sim \text{Beta}(\frac{n-d}{2}, \frac{d}{2})$, and $S := \|\mathbf{s}\|^2 \sim \text{Beta}(\frac{d-n}{2}, \frac{n}{2})$.

Proof. Recall our SVD $A = U\Sigma V^\top$ and that

$$\gamma_1 = \eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2 + \xi^2 \quad \text{and} \quad \gamma_2 = \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2.$$

Then we have that

$$\|\mathbf{k}\|^2 = \sum_{i=1}^d \frac{\mathbf{b}_i^2}{\sigma_i^2} \geq \frac{1}{\sigma_1^2} \|\mathbf{b}\|^2 = \frac{1}{\sigma_1^2}$$

Similarly,

$$\|\mathbf{h}\|^2 = \sum_{i=1}^n \frac{\mathbf{a}_i^2}{\sigma_i^2} \geq \frac{1}{\sigma_1^2} \|\mathbf{a}\|^2 = \frac{1}{\sigma_1^2}$$

Thus, we see that

$$\gamma_1 \geq \eta^2 \|\mathbf{t}\|^2 \frac{1}{\sigma_1^2} \quad \text{and} \quad \gamma_2 \geq \eta^2 \|\mathbf{s}\|^2 \frac{1}{\sigma_1^2}.$$

$\|\mathbf{t}\|^2$ depends only on \mathbf{V} and is independent of \mathbf{U}, Σ . $\|\mathbf{s}\|^2$ depends only on \mathbf{U} and is independent of \mathbf{V}, Σ . σ_1 depends only on Σ and is independent of \mathbf{U}, \mathbf{V} . Therefore, σ_1 is independent of $T := \|\mathbf{t}\|^2$ and of $S := \|\mathbf{s}\|^2$.

Thus, we get that

$$\frac{1}{\gamma_1^p} \leq \frac{1}{\eta^{2p}} \frac{\sigma_1^{2p}}{\|\mathbf{t}\|^{2p}} \quad \text{and} \quad \frac{1}{\gamma_2^p} \leq \frac{1}{\eta^{2p}} \frac{\sigma_1^{2p}}{\|\mathbf{s}\|^{2p}}$$

Then taking the expectation and using the independence, we get that

$$\mathbb{E} \left[\frac{1}{\gamma_1^p} \right] \leq \frac{1}{\eta^{2p}} \mathbb{E} \left[\frac{1}{\|\mathbf{t}\|^{2p}} \right] \mathbb{E} [\sigma_1^{2p}] \quad \text{and} \quad \mathbb{E} \left[\frac{1}{\gamma_2^p} \right] \leq \frac{1}{\eta^{2p}} \mathbb{E} \left[\frac{1}{\|\mathbf{s}\|^{2p}} \right] \mathbb{E} [\sigma_1^{2p}]$$

For $c < 1$ (where $d < n$), the right null space of \mathbf{A} (dimension $n - d$) is a uniformly random $(n - d)$ -dimensional subspace of \mathbb{R}^n . The squared norm $\|\mathbf{t}\|^2$ represents the squared length of the projection of the fixed unit vector $\mathbf{v} \in \mathbb{R}^n$ onto this random subspace. The distribution of such a squared projection norm is Beta $\left(\frac{n-d}{2}, \frac{d}{2}\right)$, as it can be represented as the ratio of two independent chi-squared random variables: $\sum_{i=1}^{n-d} G_i^2 / \sum_{i=1}^n G_i^2$, where $G_i \sim \mathcal{N}(0, 1)$ IID, which follows the desired Beta distribution. Similarly for $c > 1$.

Since the eigenvalue distribution converges to the compactly supported distribution. We can see that for sufficiently large n, d , we have that there exists an $M \geq 1$ such that $\sigma_1 \leq \rho M$ almost surely.

For $Y \sim \text{Beta}(\alpha, \beta)$ and $p < \alpha$,

$$\mathbb{E}[Y^{-p}] = \frac{\Gamma(\alpha - p) \Gamma(\alpha + \beta)}{\Gamma(\alpha) \Gamma(\alpha + \beta - p)}.$$

Moreover, using Stirling on the Γ ratio,

$$\mathbb{E}[T^{-p}] \rightarrow_{n, d \rightarrow \infty} \left(\frac{\alpha_1 + \beta_1}{\alpha_1}\right)^p = \left(\frac{1}{1 - c}\right)^p \quad (c < 1),$$

and

$$\mathbb{E}[S^{-p}] \rightarrow_{n, d \rightarrow \infty} \left(\frac{\alpha_2 + \beta_2}{\alpha_2}\right)^p = \left(\frac{c}{c - 1}\right)^p \quad (c > 1).$$

Thus, there is an M such that

$$\mathbb{E}\left[\frac{1}{\gamma_1^p}\right] \leq \left(\frac{\rho}{\eta}\right)^{2p} M^p \quad \text{and} \quad \mathbb{E}\left[\frac{1}{\gamma_2^p}\right] \leq \left(\frac{\rho}{\eta}\right)^{2p} M^p$$

□

Lemma 19 (Moments of η^2/γ_i). *We have:*

(i) For η^2/γ_1 ,

$$\mathbb{E}\left[\frac{\eta^2}{\gamma_1}\right] = \frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right), \quad \text{Var}\left(\frac{\eta^2}{\gamma_1}\right) = O\left(\frac{1}{n}\right).$$

(ii) For η^2/γ_2 ,

$$\mathbb{E}\left[\frac{\eta^2}{\gamma_2}\right] = \frac{\rho^2 \eta^2}{\eta^2 + \rho^2} + o\left(\frac{1}{\rho^2}\right), \quad \text{Var}\left(\frac{\eta^2}{\gamma_2}\right) = O\left(\frac{1}{n}\right).$$

Proof. By Lemmas 32 and 16, the expectation of η^2/γ_1 can be computed by:

$$\mathbb{E}\left[\frac{\eta^2}{\gamma_1}\right] = \frac{1}{\mathbb{E}[\gamma_1/\eta^2]} 1 + o\left(\frac{1}{\rho^2 d}\right) = \frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right).$$

By Lemmas 33 and 16, the variance of η^2/γ_1 can be computed by:

$$\begin{aligned} \text{Var}\left(\frac{\eta^2}{\gamma_1}\right) &= \frac{1}{\mathbb{E}[\gamma_1/\eta^2]^4} O\left(\text{Var}\left(\frac{\gamma_1}{\eta^2}\right)\right) + o\left(\text{Var}\left(\frac{\gamma_1}{\eta^2}\right)\right) \\ &= \frac{\rho^8 \eta^8}{(\eta^2 c + \rho^2)^4} O\left(\frac{1}{n}\right) + o\left(\frac{1}{n}\right) \\ &= O\left(\frac{1}{n}\right) \quad \text{by the scalings of } \eta \text{ and } \rho. \end{aligned}$$

The proof is similar for the other term. □

Lemma 20 (Moments of η^4/γ_i^2). *We have:*

(i) For η^4/γ_1^2 ,

$$\mathbb{E}\left[\frac{\eta^4}{\gamma_1^2}\right] = \frac{\rho^4 \eta^4}{(\eta^2 c + \rho^2)^2} + o(1), \quad \text{Var}\left(\frac{\eta^4}{\gamma_1^2}\right) = O\left(\frac{1}{n}\right).$$

(ii) For η^4/γ_2^2 ,

$$\mathbb{E} \left[\frac{\eta^4}{\gamma_2^2} \right] = \frac{\rho^4 \eta^4}{(\eta^2 + \rho^2)^2} + o(1), \quad \text{Var} \left(\frac{\eta^4}{\gamma_2^2} \right) = O \left(\frac{1}{n} \right).$$

Proof. The expectation of η^4/γ_1^2 can be computed by Lemma 19. By definition we have that

$$\mathbb{E} \left[\frac{\eta^4}{\gamma_1^2} \right] = \left(\mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \right)^2 + \text{Var} \left(\frac{\eta^2}{\gamma_1} \right) = \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o \left(\frac{1}{\rho^2} \right) \right)^2 + O \left(\frac{1}{n} \right).$$

The variance follows Lemma 33 and Lemma 17:

$$\text{Var} \left(\frac{\eta^4}{\gamma_1^2} \right) = O \left(\frac{1}{n} \right),$$

since the mean is $O(1)$.

The proof is similar for the other term. \square

Lemma 21. Suppose $\varepsilon \in \mathbb{R}^n$ whose entries have mean 0, variance τ_ε , and follow our noise assumptions. Then for any independent random matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$, we have

$$\mathbb{E}_{\varepsilon, \mathbf{Q}} [\varepsilon^\top \mathbf{Q} \varepsilon] = \tau_\varepsilon^2 \mathbb{E} [\text{Tr}(\mathbf{Q})].$$

Proof. We have that

$$\varepsilon^\top \mathbf{Q} \varepsilon = \sum_{i=1}^n \sum_{j=1}^n \varepsilon_i \varepsilon_j Q_{ij}.$$

We take the expectation of this sum. By the independence assumption and assumption $\mathbb{E}[\varepsilon_i \varepsilon_j] = 0$ when $i \neq j$, we then have

$$\mathbb{E}_{\varepsilon, \mathbf{Q}} [\varepsilon^\top \mathbf{Q} \varepsilon] = \sum_{i=1}^n \mathbb{E} [\varepsilon_i^2] \mathbb{E} [Q_{ii}] = \tau_\varepsilon^2 \mathbb{E} \left[\sum_{i=1}^n Q_{ii} \right] = \tau_\varepsilon^2 \mathbb{E} [\text{Tr}(\mathbf{Q})].$$

\square

E.5 STEP 4: BOUNDING THE EXPECTATION OF PRODUCTS OF DEPENDENT TERMS

In Section E.2 we decomposed the error into four terms – Bias, Variance, Data Noise and Target alignment. In Section E.3, we wrote each of these terms as the sum and product of various “elementary building blocks”. In Section E.4, we showed that these elementary building blocks concentrate. In this section, since we have tight concentration (i.e., the higher moment bounds). We can use Lemma 36 and Lemma 37, which shows that the expectation of the product can be approximated by the product of the expectations. In this section, we do that calculation for our different terms.

E.5.1 STEP 4: BIAS

We begin with the bias term. Recall that for $c < 1$, the expected bias by Lemma 5 is equal to

$$\mathbb{E}[\text{Bias}] = \mathbb{E} \left[\left[\tilde{\alpha}_Z - \alpha_Z + \frac{\xi}{\gamma_1} (\alpha_Z - \alpha_A) \right]^2 \tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 + \frac{\tilde{\eta}^2 \xi^2}{\eta^2 \gamma_1^2} \tau_\varepsilon^2 \|\mathbf{p}_1\|^2 \right],$$

where the cross term equals 0 due to ε having mean zero entries. These two remaining expectations are given by Lemmas 22, 23, informally via:

$$\text{Lemma 22} + \tau_\varepsilon^2 \frac{\tilde{\eta}^2}{\eta^2} \times \text{Lemma 23}.$$

For $c < 1$, we can plug in the value to get that the expected first term is given by

$$\tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 c + \rho^2} (\alpha_Z - \alpha_A) \right]^2 + o(1) + O \left(\frac{\eta}{n} \right)$$

and the second is given by

$$\tau_\varepsilon^2 \frac{\tilde{\eta}^2}{\eta^2} \left(\frac{c}{c-1} \frac{\eta^2}{\eta^2 c + \rho^2} + o(1) + O\left(\frac{1}{\rho^2 n}\right) \right).$$

Adding them, we then have the desired result:

$$\frac{\tilde{\eta}^2}{\tilde{n}} \left(\left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 c + \rho^2} (\alpha_Z - \alpha_A) \right]^2 (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + \tau_\varepsilon^2 \frac{c}{1-c} \frac{1}{\eta^2 c + \rho^2} \right) + o\left(\frac{1}{\tilde{n}}\right) + O\left(\frac{\eta}{n^2}\right).$$

For $c > 1$, we instead have the following expansion:

$$\underbrace{\boldsymbol{\beta}_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right]}_{\mathbf{t}_1} \tilde{\mathbf{Z}} - \underbrace{\alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top}_{\mathbf{t}_2} \tilde{\mathbf{Z}} + \underbrace{\frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top}_{\mathbf{t}_3}$$

The bias equals the expectation of the norm of this vector. Taking the Frobenius norm, we have the six terms. Among the cross-terms, $\langle \mathbf{t}_1, \mathbf{t}_3 \rangle$ and $\langle \mathbf{t}_2, \mathbf{t}_3 \rangle$ have zero mean since \mathbf{t}_3 contains $\boldsymbol{\varepsilon}$ whose entries have mean 0. We now look at the other terms

$$\mathbb{E}[\|\mathbf{t}_3\|^2] = \mathbb{E} \left[\left\| \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top \right\|^2 \right] = \tau_\varepsilon^2 \frac{\tilde{\eta}^2}{\eta^2} \mathbb{E} \left[\frac{\xi^2}{\gamma_2^2} \|\mathbf{p}_2\|^2 \right] \quad \text{by Lemma 21}$$

The expectation is given by Lemma 23. Subsequently, Lemmas 22, 24, 25 give $\mathbb{E}[\|\mathbf{t}_1\|^2]$, $\mathbb{E}[\|\mathbf{t}_2\|^2]$, $\mathbb{E}[\langle \mathbf{t}_1, \mathbf{t}_3 \rangle]$ respectively. Informally, we can compute the bias via:

$$\begin{aligned} \mathbb{E}[\mathbf{Bias}] &= \mathbb{E} \left[\left\| \boldsymbol{\beta}_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} - \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} + \frac{\tilde{\eta}}{\eta} \frac{\xi}{\gamma_2} \boldsymbol{\varepsilon}^\top \mathbf{p}_2 \tilde{\mathbf{v}}^\top \right\|^2 \right] \\ &= \mathbb{E}[\|\mathbf{t}_1\|^2] + \mathbb{E}[\|\mathbf{t}_2\|^2] + \mathbb{E}[\|\mathbf{t}_3\|^2] - 2\mathbb{E}[\langle \mathbf{t}_1, \mathbf{t}_3 \rangle] \\ &= \text{Lemma 22} + \tau_\varepsilon^2 \frac{\tilde{\eta}^2}{\eta^2} \text{Lemma 23} + \text{Lemma 24} - 2 \times \text{Lemma 25}. \end{aligned}$$

Similar to $c < 1$, adding them together and dividing by \tilde{n} , we get

$$\begin{aligned} \frac{\tilde{\eta}^2}{\tilde{n}} \left[(\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left((\tilde{\alpha}_Z - \alpha_Z)^2 + \frac{\rho^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 \right) + \alpha_A^2 \frac{\|\boldsymbol{\beta}_*\|^2}{d} \left(\frac{c-1}{c} \right) \frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} + \frac{\tau_\varepsilon^2}{c-1} \frac{\eta^2 c + \rho^2}{(\eta^2 + \rho^2)^2} \right] \\ + o\left(\frac{1}{\tilde{n}}\right) + O\left(\frac{\eta}{n^2}\right). \end{aligned}$$

E.5.2 STEP 4: VARIANCE

Recall that for the variance, we have the following expression (Section E.3.2).

$$\begin{aligned} \mathbb{E} \left[\frac{1}{\tilde{n}} \left\| \boldsymbol{\beta}_{int}^\top \tilde{\mathbf{A}} \right\|_F^2 \right] &= \mathbb{E} \left[\frac{\tilde{\tau}^2 \alpha_z^2}{d} \boldsymbol{\beta}_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{Z} \boldsymbol{\beta}_* + \frac{\tilde{\tau}^2 \alpha_A^2}{d} \boldsymbol{\beta}_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \boldsymbol{\beta}_* \right. \\ &\quad \left. + \frac{2\tilde{\tau}^2 \alpha_A \alpha_z}{d} \boldsymbol{\beta}_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \boldsymbol{\beta}_* + \frac{\tilde{\tau}^2}{d} \boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \boldsymbol{\varepsilon} \right]. \end{aligned}$$

In particular that the expectation will be the weighted sum of the expressions from Lemmas 26, 27, 28, 29. Informally,

$$\frac{\tilde{\rho}^2}{d} (\alpha_Z^2 \times \text{Lemma 26} + 2\alpha_Z \alpha_A \times \text{Lemma 28} + \alpha_A^2 \times \text{Lemma 27} + \text{Lemma 29}).$$

This yields that for $c < 1$, after simplification, the variance is

$$\begin{aligned} \frac{\tilde{\rho}^2}{d} \left[\alpha_A^2 \|\boldsymbol{\beta}_*\|^2 + (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left[(\alpha_Z - \alpha_A)^2 \frac{\eta^2 (\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} + 2\alpha_A (\alpha_Z - \alpha_A) \frac{\eta^2 c}{\eta^2 c + \rho^2} \right] \right. \\ \left. + \tau_\varepsilon^2 \left(\frac{c}{1-c} \frac{d}{\rho^2} - \frac{\eta^2}{\rho^2 (\eta^2 c + \rho^2)} \frac{c^2}{1-c} \right) \right] + o(1) + O\left(\frac{1}{n}\right). \end{aligned}$$

For $c > 1$, we similarly simplify it to:

$$\begin{aligned} & \frac{\tilde{\rho}^2}{d} \left[\|\boldsymbol{\beta}_*\|^2 \left(\frac{\alpha_A^2}{c} - \frac{\alpha_A^2}{d} \frac{\eta^2}{\eta^2 + \rho^2} \right) + (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{c}{c-1} \frac{\eta^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 \right. \\ & \left. + \tau_\varepsilon^2 \left(\frac{d}{\rho^2} \frac{1}{c-1} - \frac{\eta^2}{\rho^2(\eta^2 + \rho^2)} \frac{c}{c-1} \right) \right] + o(1) + O\left(\frac{1}{n}\right). \end{aligned}$$

E.5.3 STEP 4: DATA NOISE

Recall that for the data noise, we have the following expression

$$\frac{\tilde{\alpha}_A^2 \tilde{\rho}^2}{d} \|\boldsymbol{\beta}_*\|^2$$

Noting that $\|\boldsymbol{\beta}_*\|^2 = \Theta(1)$, we see that this term has no more randomness and we do not need to estimate anything.

E.5.4 STEP 4: TARGET ALIGNMENT

Recall from Section E.3.4 that the alignment is given by

$$-\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \mathbb{E} [\alpha_Z \boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{Z}^\top \boldsymbol{\beta}_* + \alpha_A \boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \boldsymbol{\beta}_*]$$

From Lemma 30, we have that

$$\mathbb{E} [\boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{Z}^\top \boldsymbol{\beta}_*] = \begin{cases} \frac{\eta^2 c}{\rho^2 + \eta^2 c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right) & c < 1 \\ \frac{\eta^2}{\eta^2 + \rho^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right) & c > 1 \end{cases}.$$

and from Lemma 31, we have that

$$\mathbb{E} [\boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \boldsymbol{\beta}_*] = \begin{cases} \|\boldsymbol{\beta}_*\|^2 - \frac{\eta^2 c}{\rho^2 + \eta^2 c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o\left(\frac{1}{\rho^2}\right) + O\left(\frac{1}{n}\right), & c < 1 \\ \frac{1}{c} \|\boldsymbol{\beta}_*\|^2 - \frac{\eta^2}{\eta^2 + \rho^2} \left(\frac{\|\boldsymbol{\beta}_*\|^2}{d} + \frac{1}{c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \right) + o(1) + O\left(\frac{1}{n}\right), & c > 1 \end{cases}.$$

Thus for $c < 1$, the entire interaction term now becomes

$$-\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \left(\alpha_A \|\boldsymbol{\beta}_*\|^2 + (\alpha_Z - \alpha_A) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{\eta^2 c}{\rho^2 + \eta^2 c} + o(1) \right).$$

For $c > 1$, instead we have

$$-\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \left(\frac{\alpha_A}{c} \|\boldsymbol{\beta}_*\|^2 - \frac{\alpha_A}{d} \frac{\eta^2}{\eta^2 + \rho^2} \|\boldsymbol{\beta}_*\|^2 + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\eta^2}{\eta^2 + \rho^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) \right).$$

E.5.5 BIAS: HELPER LEMMAS

Lemma 22. *In the same setting as Section 2, we have that for $c < 1$,*

$$\begin{aligned} & \mathbb{E} \left[\left(\tilde{\alpha}_Z - \alpha_Z + \frac{\xi}{\gamma_1} (\alpha_Z - \alpha_A) \right)^2 \tilde{\eta}^2 (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \right] \\ & = \tilde{\eta}^2 (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 c + \rho^2} (\alpha_Z - \alpha_A) \right]^2 + o(1) + O\left(\frac{\eta}{n}\right). \end{aligned}$$

For $c > 1$,

$$\begin{aligned} & \mathbb{E} \left[\left\| \boldsymbol{\beta}_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} \right\|^2 \right] \\ & = \tilde{\eta}^2 (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right) \right]^2 + o(1) + O\left(\frac{\eta}{n}\right). \end{aligned}$$

Proof. For $c < 1$, we first expand the square and get:

$$\left(\tilde{\alpha}_Z - \alpha_Z + \frac{\xi}{\gamma_1}(\alpha_Z - \alpha_A)\right)^2 = (\tilde{\alpha}_Z - \alpha_Z)^2 + \frac{1}{\eta^2} \frac{\eta^2 \xi^2}{\gamma_1^2} (\alpha_Z - \alpha_A)^2 + \frac{2}{\eta} \frac{\eta \xi}{\gamma_1} (\alpha_Z - \alpha_A)(\tilde{\alpha}_Z - \alpha_Z).$$

By Lemmas 9 and 20, then we see that, using the square root of the covariance to bound the difference between the expectation of the product and the product of the expectation.

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \xi^2}{\gamma_1^2} \right] &= \mathbb{E} \left[\frac{\eta^4}{\gamma_1^2} \right] \mathbb{E} \left[\frac{\xi^2}{\eta^2} \right] + \sqrt{\text{Var} \left(\frac{\eta^4}{\gamma_1^2} \right) \text{Var} \left(\frac{\xi^2}{\eta^2} \right)} \\ &= \left(\frac{\rho^4 \eta^4}{(\eta^2 c + \rho^2)^2} + o(1) \right) \left(\frac{1}{\eta^2} + O \left(\frac{1}{\rho^2 n} \right) \right) + O \left(\frac{1}{n} \right) \\ &= \frac{\rho^4 \eta^2}{(\eta^2 c + \rho^2)^2} + o \left(\frac{1}{\eta^2} \right) + O \left(\frac{1}{n} \right). \\ \mathbb{E} \left[\frac{\eta \xi}{\gamma_1} \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} \left[\frac{\xi}{\eta} \right] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_1} \right) \text{Var} \left(\frac{\xi}{\eta} \right)} \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o(1) \right) \left(\frac{1}{\eta} \right) + O \left(\frac{1}{n} \right) \\ &= \frac{\rho^2 \eta}{\eta^2 c + \rho^2} + o \left(\frac{1}{\eta} \right) + O \left(\frac{1}{n} \right). \end{aligned}$$

Combining these terms together, we have that

$$\begin{aligned} &\mathbb{E} \left[\left(\tilde{\alpha}_Z - \alpha_Z + \frac{\xi}{\gamma_1}(\alpha_Z - \alpha_A) \right)^2 \tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \right] \\ &= (\beta_*^\top \mathbf{u})^2 \left[\tilde{\eta}^2 (\tilde{\alpha}_Z - \alpha_Z)^2 + \frac{\tilde{\eta}^2}{\eta^2} \left(\frac{\rho^4 \eta^2}{(\eta^2 c + \rho^2)^2} + o \left(\frac{1}{\eta^2} \right) + O \left(\frac{1}{n} \right) \right) (\alpha_Z - \alpha_A)^2 \right. \\ &\quad \left. + \frac{2\tilde{\eta}^2}{\eta} \left(\frac{\rho^2 \eta}{\eta^2 c + \rho^2} + o \left(\frac{1}{\eta} \right) + O \left(\frac{1}{n} \right) \right) (\alpha_Z - \alpha_A)(\tilde{\alpha}_Z - \alpha_Z) \right] \\ &= \tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \left(\left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 c + \rho^2} (\alpha_Z - \alpha_A) \right]^2 + o \left(\frac{1}{\eta^2} \right) + O \left(\frac{1}{\eta n} \right) \right). \end{aligned}$$

We now consider $c > 1$. Recalling that $\tilde{\mathbf{Z}} = \tilde{\eta} \mathbf{u} \tilde{\mathbf{v}}^\top$, we let $c_1 = \tilde{\alpha}_Z - \alpha_Z$ and expand:

$$\begin{aligned} &\left\| \beta_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} \right\|^2 \\ &= \beta_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right]^\top \beta_* \\ &= \tilde{\eta}^2 \beta_*^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right] \mathbf{u} \mathbf{u}^\top \left[(\tilde{\alpha}_Z - \alpha_Z) \mathbf{I} + \frac{\xi}{\gamma_2} (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \right]^\top \beta_* \\ &= c_1^2 \tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 + \tilde{\eta}^2 \frac{\xi^2}{\gamma_2^2} \beta_*^\top (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} \mathbf{u}^\top ((\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger)^\top \beta_* + 2c_1 \tilde{\eta}^2 \frac{\xi}{\gamma_2} \beta_*^\top (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} \mathbf{u}^\top \beta_*. \end{aligned}$$

Not that for the second and third terms, we have that ξ, γ_2 only depend on the singular values of \mathbf{A} and the rest only depend on the singular vectors. Hence, these terms are independent.

First note that when $d > n$, the number of singular values equals n , which is less than the dimension d . As a result,

$$\mathbf{A} \mathbf{A}^\dagger = \mathbf{U} \Sigma \mathbf{V}^\top \mathbf{V} \Sigma^\dagger \mathbf{U}^\top = \mathbf{U} \begin{bmatrix} \mathbf{I}_{n \times n} & \mathbf{0}_{n \times (d-n)} \\ \mathbf{0}_{(d-n) \times n} & \mathbf{0}_{(d-n) \times (d-n)} \end{bmatrix} \mathbf{U}^\top.$$

Then we have that

$$\mathbb{E} [\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \beta_*^\top] = \sum_{i=1}^n \mathbb{E} [(\beta_*^\top \mathbf{U})_i^2] = \frac{n}{d} \|\beta_*\|^2 = \frac{1}{c} \|\beta_*\|^2, \quad (13)$$

since $\beta_*^\top U$ is a uniformly random vector of length $\|\beta_*\|$ in \mathbb{R}^d after the rotation U .

For the middle term, by Proposition 2 and the above Equation 13, we have

$$\begin{aligned} & \mathbb{E} [\beta_*^\top (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} \mathbf{u}^\top ((\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \beta_*)] \\ &= \alpha_Z^2 (\beta_*^\top \mathbf{u})^2 - 2\alpha_A \alpha_Z \mathbb{E} [\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{u}^\top \beta_*] + \alpha_A^2 \mathbb{E} [(\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u})^2] \\ &= \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 (\beta_*^\top \mathbf{u})^2 + o(1). \end{aligned}$$

Similarly, for the last term, we have

$$\mathbb{E} [\beta_*^\top (\alpha_Z \mathbf{I} - \alpha_A \mathbf{A} \mathbf{A}^\dagger) \mathbf{u} \mathbf{u}^\top \beta_*] = \left(\alpha_Z - \frac{\alpha_A}{c} \right) (\beta_*^\top \mathbf{u})^2 + o(1).$$

Thus putting these expectations together, we get

$$\mathbb{E} \left[\tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \left[c_1^2 + \frac{\xi^2}{\gamma_2^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + 2c_1 \frac{\xi}{\gamma_2} \left(\alpha_Z - \frac{\alpha_A}{c} \right) \right] \right] = \mathbb{E} \left[\tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \left[c_1 + \frac{\xi}{\gamma_2} \left(\alpha_Z - \frac{\alpha_A}{c} \right) \right]^2 \right].$$

Similar to the $c < 1$ case, we take the expectation for terms involving $\frac{\xi}{\gamma_2}$ and get:

$$\tilde{\eta}^2 (\beta_*^\top \mathbf{u})^2 \left[\left((\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right) \right)^2 + o\left(\frac{1}{\eta^2}\right) + O\left(\frac{1}{\eta n}\right) \right].$$

□

Lemma 23 (Expectations involving p_1 and p_2). *In the setting of Section 2, we have that*

1. For $c = d/n < 1$:

$$\mathbb{E} \left[\frac{\xi^2}{\gamma_1^2} \|\mathbf{p}_1\|^2 \right] = \frac{c}{1-c} \frac{\eta^2}{\eta^2 c + \rho^2} + o(1) + O\left(\frac{1}{\rho^2 n}\right).$$

2. For $c = d/n > 1$:

$$\mathbb{E} \left[\frac{\xi^2}{\gamma_2^2} \|\mathbf{p}_2\|^2 \right] = \frac{\eta^2}{c-1} \frac{\eta^2 c + \rho^2}{(\eta^2 + \rho^2)^2} + o(1) + O\left(\frac{1}{\rho^2 n}\right).$$

Proof. First, Lemma 6 tells us that

$$\frac{\xi^2}{\gamma_1^2} \|\mathbf{p}_1\|^2 = \frac{\eta^2 \|\mathbf{k}\|^2}{\gamma_1}.$$

Then recall from Lemma 9 that

$$\mathbb{E}[\|\mathbf{k}\|^2] = \frac{1}{\rho^2} \frac{c}{1-c} + o\left(\frac{1}{\rho^2}\right) \quad \text{and} \quad \text{Var}(\|\mathbf{k}\|^2) = O\left(\frac{1}{\rho^4 n}\right)$$

and Lemma 19 tells us

$$\mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] = \frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right) \quad \text{and} \quad \text{Var} \left(\frac{\eta^2}{\gamma_i} \right) = O\left(\frac{1}{n}\right)$$

Again Section E.4.2 tells us that the assumption of Lemma 37 is satisfied and that

$$\begin{aligned} \mathbb{E} \left[\frac{\xi^2}{\gamma_1^2} \|\mathbf{p}_1\|^2 \right] &= \mathbb{E} \left[\frac{\eta^2 \|\mathbf{k}\|^2}{\gamma_1} \right] = \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} [\|\mathbf{k}\|^2] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_1} \right) \text{Var} (\|\mathbf{k}\|^2)} \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right) \right) \left(\frac{1}{\rho^2} \frac{c}{1-c} + o\left(\frac{1}{\rho^2}\right) \right) + O\left(\frac{1}{\rho^2 n}\right) \\ &= \frac{c}{1-c} \frac{\eta^2}{\eta^2 c + \rho^2} + o(1) + O\left(\frac{1}{\rho^2 n}\right). \end{aligned}$$

Using Lemma 6 for \mathbf{p}_2 ,

$$\frac{\xi^2}{\gamma_2^2} \|\mathbf{p}_2\|^2 = \frac{1}{\gamma_2^2} (\eta^4 \|\mathbf{s}\|^4 \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top} + 2\eta^3 \xi \|\mathbf{s}\|^2 \mathbf{k}^{\top} \mathbf{A}^{\dagger} \mathbf{h}^{\top} + \eta^2 \xi^2 \|\mathbf{k}\|^2).$$

To begin, we start estimating

$$\mathbb{E} \left[\frac{\eta^4 \|\mathbf{s}\|^4}{\gamma_2^2} \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top} \right].$$

Using our Spherical Hypercontractivity, we have that $\|\mathbf{s}\|^2$ and $\mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}$ satisfy the assumptions for Lemma 36. Then using Lemmas 9 and 10 we first have that

$$\begin{aligned} \mathbb{E} [\|\mathbf{s}\|^2] &= 1 - \frac{1}{c} \quad \text{and} \quad \text{Var} (\|\mathbf{s}\|^2) = O\left(\frac{1}{d}\right) \\ \mathbb{E} [\mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}] &= \frac{1}{\rho^4} \frac{c^3}{(c-1)^3} + o\left(\frac{1}{\rho^4}\right) \quad \text{and} \quad \text{Var} (\beta_*^{\top} \mathbf{h}^{\top} \mathbf{u}^{\top} \beta_*) = O\left(\frac{1}{\rho^8 d}\right). \end{aligned}$$

Thus, using Lemma 37, we have that

$$\begin{aligned} \mathbb{E} [\|\mathbf{s}\|^4 \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}] &= (\mathbb{E} [\|\mathbf{s}\|^2])^2 \mathbb{E} [\mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}] + O\left(\max\left(\frac{1}{d}, \frac{1}{\rho^8 d}\right)\right) \\ &= \left(1 - \frac{1}{c}\right)^2 \left(\frac{1}{\rho^4} \frac{c^3}{(c-1)^3} + o\left(\frac{1}{\rho^4}\right)\right) + O\left(\frac{1}{n}\right) \\ &= \frac{1}{\rho^4} \frac{c}{c-1} + o\left(\frac{1}{\rho^4}\right) + O\left(\frac{1}{n}\right). \end{aligned}$$

and using Lemma 36, since all the means are $O(1)$, we have that

$$\text{Var} (\|\mathbf{s}\|^4 \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}) = O(\max(\text{Var} (\|\mathbf{s}\|^2), \text{Var} (\mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}))) = O\left(\frac{1}{n}\right).$$

Then Lemma 20 gives mean and variance of $\frac{\eta^4}{\gamma_2^2}$. Since $\frac{\eta^4}{\gamma_2^2}$ does not satisfy the higher moment bound, and cannot be directly included in the product, we can include it via the classical bound:

$$\mathbb{E} \left[\frac{\eta^4 \|\mathbf{s}\|^4}{\gamma_2^2} \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top} \right] = \mathbb{E} \left[\frac{\eta^4}{\gamma_2^2} \right] \mathbb{E} [\|\mathbf{s}\|^4 \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}] + \sqrt{\text{Var} (\|\mathbf{s}\|^4 \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}) \text{Var} \left(\frac{\eta^4}{\gamma_2^2}\right)} \quad (14)$$

$$= \left(\frac{\rho^4 \eta^4}{(\eta^2 + \rho^2)^2} + o(1)\right) \left(\frac{1}{\rho^4} \frac{c}{c-1} + o\left(\frac{1}{\rho^4}\right)\right) + O\left(\frac{1}{n}\right) \quad (15)$$

$$= \frac{c}{c-1} \frac{\eta^4}{(\eta^2 + \rho^2)^2} + o(1) + O\left(\frac{1}{n}\right). \quad (16)$$

Similarly, we can do the same thing for the other term. For the middle term we note that from Lemma 11

$$\mathbb{E} [\mathbf{k}^{\top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}] = 0 \quad \text{and} \quad \text{Var} (\mathbf{k}^{\top} \mathbf{A}^{\dagger} \mathbf{h}^{\top}) = O\left(\frac{1}{\rho^6 d}\right)$$

and Lemma 9 tells us

$$\mathbb{E} [\|\mathbf{s}\|^2] = 1 - \frac{1}{c} \quad \text{and} \quad \text{Var} (\|\mathbf{s}\|^2) = O\left(\frac{1}{d}\right)$$

and

$$\mathbb{E} \left[\frac{\xi}{\eta} \right] = \frac{1}{\eta} \quad \text{and} \quad \text{Var} \left(\frac{\xi}{\eta} \right) = O\left(\frac{1}{\rho^2 n}\right)$$

Thus using Lemma 37, we have that

$$\mathbb{E} \left[\frac{\xi}{\eta} \|\mathbf{s}\|^2 \mathbf{k}^{\top} \mathbf{A}^{\dagger} \mathbf{h}^{\top} \right] = 0 + O\left(\frac{1}{d}\right)$$

Thus using the standard covariance bound for the expectation of product versus product of expectation, we have that

$$\mathbb{E} \left[\frac{\eta^3 \xi \|\mathbf{s}\|^2}{\gamma_2^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right] = 0 + \sqrt{\text{Var} \left(\frac{\eta^4}{\gamma_2^2} \right) O \left(\frac{1}{n} \right)} = O \left(\frac{1}{n} \right).$$

For the last term, we have that, using Lemma 37

$$\begin{aligned} \mathbb{E} \left[\frac{\xi^2}{\eta^2} \|\mathbf{k}\|^2 \right] &= \frac{1}{\eta^2} \cdot \left(\frac{1}{\rho^2} \frac{1}{c-1} + o \left(\frac{1}{\rho^2} \right) \right) + O \left(\frac{1}{\rho^4 n} \right) \\ &= \frac{1}{\eta^2 \rho^2} \frac{1}{c-1} + o \left(\frac{1}{\eta^2 \rho^2} \right) + O \left(\frac{1}{\rho^4 n} \right) \end{aligned}$$

and from Lemma 36

$$\text{Var} \left(\frac{\xi^2}{\eta^2} \|\mathbf{k}\|^2 \right) = O \left(\frac{1}{\rho^4 n} \right)$$

Then using the standard bound, we have that

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \xi^2 \|\mathbf{k}\|^2}{\gamma_2^2} \right] &= \mathbb{E} \left[\frac{\eta^4}{\gamma_2^2} \right] \mathbb{E} \left[\frac{\xi^2}{\eta^2} \|\mathbf{k}\|^2 \right] + \sqrt{\text{Var} \left(\frac{\eta^4}{\gamma_2^2} \right) O \left(\frac{1}{\rho^4 n} \right)} \\ &= \left(\frac{\rho^4 \eta^4}{(\eta^2 + \rho^2)^2} + o(1) \right) \left(\frac{1}{\eta^2 \rho^2} \frac{1}{c-1} + o \left(\frac{1}{\eta^2 \rho^2} \right) + O \left(\frac{1}{\rho^4 n} \right) \right) + O \left(\frac{1}{\rho^2 n} \right) \\ &= \frac{1}{c-1} \frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} + o \left(\frac{1}{\eta^2 \rho^2} \right) + O \left(\frac{1}{\rho^2 n} \right). \end{aligned}$$

Finally, putting all three terms together we get

$$\begin{aligned} \mathbb{E} \left[\frac{\xi^2}{\gamma_2^2} \|\mathbf{p}_2\|^2 \right] &= \frac{c}{c-1} \frac{\eta^4}{(\eta^2 + \rho^2)^2} + o(1) + \frac{1}{c-1} \frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} + o \left(\frac{1}{\rho^2 \eta^2} \right) + O \left(\frac{1}{\rho^2 n} \right) \\ &= \frac{\eta^2}{c-1} \frac{\eta^2 c + \rho^2}{(\eta^2 + \rho^2)^2} + o(1) + O \left(\frac{1}{\rho^2 n} \right). \end{aligned}$$

□

From the above proofs, we make an important observation that the individual terms from Lemmas 9, 10, 11, 16 all have means $O(1)$ and variances $O(1/n)$. Hence, by Lemma 36, we can bound the variance of a product of terms by $O(1/n)$, given that these terms satisfy the lemma assumptions. Essentially, only η^2/γ_i and η^4/γ_i^2 fail the assumption on higher moment bound, so we deal with them via the classical bound after carrying out the product. This simplification ensures proper concentration and will be used at times in the following proofs without reference.

E.5.6 VARIANCE: HELPER LEMMAS

Lemma 24. *In the setting of Section 2, we have that for $c > 1$:*

$$\mathbb{E} \left[\left\| \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} \right\|^2 \right] = \tilde{\eta}^2 \alpha_A^2 \frac{\|\boldsymbol{\beta}_*\|^2}{d} \left(\frac{c-1}{c} \right) \frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} + O \left(\frac{1}{n} \right).$$

Proof. Since $\tilde{\mathbf{Z}} = \tilde{\eta} \mathbf{u} \tilde{\mathbf{v}}^\top$, we have that

$$\left\| \alpha_A \frac{\eta \|\mathbf{s}\|^2}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \tilde{\mathbf{Z}} \right\|^2 = \tilde{\eta}^2 \alpha_A^2 \frac{\eta^2 \|\mathbf{s}\|^4}{\gamma_2^2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_* = \alpha_A^2 \frac{\tilde{\eta}^2}{\eta^2} \frac{\eta^4 \|\mathbf{s}\|^4}{\gamma_2^2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_*.$$

Similar to last lemma, using Lemmas 37, 9, 10, 20, we get

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^4 \|\mathbf{s}\|^4}{\gamma_2^2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_* \right] &= \mathbb{E} \left[\frac{\eta^4}{\gamma_2^2} \right] (\mathbb{E} [\|\mathbf{s}\|^2])^2 \mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{h} \boldsymbol{\beta}_*] + \sqrt{\text{Var} \left(\frac{\eta^4}{\gamma_2^2} \right) O \left(\frac{1}{n} \right)} \\ &= \left(\frac{\rho^4 \eta^4}{(\rho^2 + \eta^2)^2} + o(1) \right) \left(1 - \frac{1}{c} \right)^2 \left(\frac{\|\boldsymbol{\beta}_*\|^2}{d} \frac{c}{\rho^2(c-1)} + o \left(\frac{1}{\rho^2 d} \right) \right) + O \left(\frac{1}{n} \right) \\ &= \frac{\|\boldsymbol{\beta}_*\|^2}{d} \left(\frac{c-1}{c} \right) \frac{\eta^4 \rho^2}{(\eta^2 + \rho^2)^2} + O \left(\frac{1}{n} \right). \end{aligned}$$

Hence, it directly follows from here that

$$\begin{aligned}\mathbb{E}\left[\left\|\alpha_A\frac{\eta\|\mathbf{s}\|^2}{\gamma_2}\boldsymbol{\beta}_*^\top\mathbf{h}^\top\mathbf{u}^\top\tilde{\mathbf{Z}}\right\|^2\right] &= \alpha_A^2\frac{\tilde{\eta}^2}{\eta^2}\mathbb{E}\left[\frac{\eta^4\|\mathbf{s}\|^4}{\gamma_2^2}\boldsymbol{\beta}_*^\top\mathbf{h}^\top\mathbf{h}\boldsymbol{\beta}_*\right] \\ &= \tilde{\eta}^2\alpha_A^2\frac{\|\boldsymbol{\beta}_*\|^2}{d}\left(\frac{c-1}{c}\right)\frac{\eta^2\rho^2}{(\eta^2+\rho^2)^2}+O\left(\frac{1}{n}\right).\end{aligned}$$

□

Lemma 25. *In the setting of Section 2, we have that for $c > 1$:*

$$\mathbb{E}\left[\frac{\eta\|\mathbf{s}\|^2}{\gamma_2}\boldsymbol{\beta}_*^\top\left[(\tilde{\alpha}_Z-\alpha_Z)\mathbf{I}+\frac{\xi}{\gamma_2}(\alpha_Z\mathbf{I}-\alpha_A\mathbf{A}\mathbf{A}^\dagger)\right]\tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*\right]=O\left(\frac{\eta}{n}\right).$$

Proof. Using $\tilde{\mathbf{Z}} = \tilde{\eta}\mathbf{u}\tilde{\mathbf{v}}^\top$, we can expand this into three terms. We can take expectations in a similar way via Lemmas 37, 9, 10, 11: Let $c_1 = \tilde{\alpha}_Z - \alpha_Z$. Each term contains a zero expectation:

$$\begin{aligned}\mathbb{E}\left[\tilde{\eta}^2c_1\frac{\eta\|\mathbf{s}\|^2}{\gamma_2}\boldsymbol{\beta}_*^\top\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*\right] &= \frac{\tilde{\eta}^2}{\eta}c_1\left(\mathbb{E}\left[\frac{\eta^2}{\gamma_2}\right]\mathbb{E}[\|\mathbf{s}\|^2]\mathbb{E}[\boldsymbol{\beta}_*^\top\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*]+\sqrt{\text{Var}\left(\frac{\eta^2}{\gamma_2}\right)}O\left(\frac{1}{n}\right)\right) \\ &= \frac{\tilde{\eta}^2}{\eta}c_1\left(\sqrt{\text{Var}\left(\frac{\eta^2}{\gamma_2}\right)}O\left(\frac{1}{n}\right)\right)=O\left(\frac{\eta}{n}\right).\end{aligned}$$

$$\begin{aligned}\mathbb{E}\left[\tilde{\eta}^2\alpha_Z\frac{\eta\xi\|\mathbf{s}\|^2}{\gamma_2^2}\boldsymbol{\beta}_*^\top\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*\right] &= \frac{\alpha_Z\tilde{\eta}^2}{\eta^2}\left(\mathbb{E}\left[\frac{\eta^4}{\gamma_2^2}\right]\mathbb{E}\left[\frac{\xi}{\eta}\right]\mathbb{E}[\|\mathbf{s}\|^2]\mathbb{E}[\boldsymbol{\beta}_*^\top\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*]+\sqrt{\text{Var}\left(\frac{\eta^4}{\gamma_2^2}\right)}O\left(\frac{1}{n}\right)\right) \\ &= \frac{\alpha_Z\tilde{\eta}^2}{\eta^2}\left(\sqrt{\text{Var}\left(\frac{\eta^4}{\gamma_2^2}\right)}O\left(\frac{1}{n}\right)\right)=O\left(\frac{1}{n}\right).\end{aligned}$$

$$\begin{aligned}\mathbb{E}\left[\tilde{\eta}^2\alpha_A\frac{\eta\xi\|\mathbf{s}\|^2}{\gamma_2^2}\boldsymbol{\beta}_*^\top\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*\right] &= \frac{\alpha_Z\tilde{\eta}^2}{\eta^2}\left(\mathbb{E}\left[\frac{\eta^4}{\gamma_2^2}\right]\mathbb{E}\left[\frac{\xi}{\eta}\right]\mathbb{E}[\|\mathbf{s}\|^2]\mathbb{E}[\boldsymbol{\beta}_*^\top\mathbf{A}\mathbf{A}^\dagger\mathbf{u}\mathbf{h}\boldsymbol{\beta}_*]+\sqrt{\text{Var}\left(\frac{\eta^4}{\gamma_2^2}\right)}O\left(\frac{1}{n}\right)\right) \\ &= \frac{\alpha_Z\tilde{\eta}^2}{\eta^2}\left(\sqrt{\text{Var}\left(\frac{\eta^4}{\gamma_2^2}\right)}O\left(\frac{1}{n}\right)\right)=O\left(\frac{1}{n}\right).\end{aligned}$$

Thus the cross term concentrates around zero at a rate of $O(\eta/n)$. □

Lemma 26. *In the same setting as Section 2, we have that*

$$\mathbb{E}[\boldsymbol{\beta}_*^\top\mathbf{Z}(\mathbf{Z}+\mathbf{A})^\dagger(\mathbf{Z}+\mathbf{A})^\dagger\mathbf{Z}\boldsymbol{\beta}_*]=\begin{cases}\frac{\eta^2(\eta^2+\rho^2)}{(\eta^2c+\rho^2)^2}\frac{c^2}{1-c}(\boldsymbol{\beta}_*^\top\mathbf{u})^2+o(1)+O\left(\frac{1}{n}\right) & c < 1 \\ \frac{\eta^2}{\eta^2+\rho^2}\frac{c}{c-1}(\boldsymbol{\beta}_*^\top\mathbf{u})^2+o\left(\frac{1}{\rho^2}\right)+O\left(\frac{1}{\rho^2n}\right) & c > 1\end{cases}.$$

Proof. We start with $c < 1$ and expand this term using Proposition 1:

$$\boldsymbol{\beta}_*^\top\mathbf{Z}(\mathbf{Z}+\mathbf{A})^\dagger(\mathbf{Z}+\mathbf{A})^\dagger\mathbf{Z}\boldsymbol{\beta}_*=\frac{\eta^2\|\mathbf{h}\|^2\xi^2}{\gamma_1^2}(\boldsymbol{\beta}_*^\top\mathbf{u})^2+\frac{\eta^4\|\mathbf{t}\|^4}{\gamma_1^2}(\mathbf{k}^\top\mathbf{A}^\dagger\mathbf{A}^\dagger\mathbf{k})(\boldsymbol{\beta}_*^\top\mathbf{u})^2+\frac{2\eta^3\|\mathbf{t}\|^2\xi}{\gamma_1^2}\mathbf{k}^\top\mathbf{A}^\dagger\mathbf{h}^\top(\boldsymbol{\beta}_*^\top\mathbf{u})^2.$$

We then start plugging in the expectations of these three terms and the ‘‘cumulative’’ variance of the sum according to Lemma 37.

$$\begin{aligned}\mathbb{E}\left[\frac{\eta^2\|\mathbf{h}\|^2\xi^2}{\gamma_1^2}(\boldsymbol{\beta}_*^\top\mathbf{u})^2\right] &= (\boldsymbol{\beta}_*^\top\mathbf{u})^2\mathbb{E}\left[\frac{\eta^4}{\gamma_1^2}\right]\mathbb{E}\left[\frac{\xi^2}{\eta^2}\right]\mathbb{E}[\|\mathbf{h}\|^2]+\sqrt{\text{Var}\left(\frac{\eta^4}{\gamma_1^2}\right)}O\left(\frac{1}{n}\right) \\ &= (\boldsymbol{\beta}_*^\top\mathbf{u})^2\left(\frac{\rho^4\eta^4}{(\eta^2c+\rho^2)^2}+o(1)\right)\left(\frac{1}{\eta^2}+O\left(\frac{1}{\rho^2n}\right)\right)\left(\frac{1}{\rho^2}\frac{c^2}{1-c}+o\left(\frac{1}{\rho^2}\right)\right)+O\left(\frac{1}{n}\right) \\ &= \frac{\eta^2\rho^2}{(\eta^2c+\rho^2)^2}\frac{c^2}{1-c}(\boldsymbol{\beta}_*^\top\mathbf{u})^2+o(1)+O\left(\frac{1}{n}\right).\end{aligned}$$

$$\begin{aligned}
\mathbb{E} \left[\frac{\eta^4 \|\mathbf{t}\|^4}{\gamma_1^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \right] &= (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \mathbb{E} \left[\frac{\eta^4}{\gamma_1^2} \right] (\mathbb{E} [\|\mathbf{t}\|^2])^2 \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}] + \sqrt{\text{Var} \left(\frac{\eta^4}{\gamma_1^2} \right)} O \left(\frac{1}{n} \right) \\
&= (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left(\frac{\rho^4 \eta^4}{(\eta^2 c + \rho^2)^2} + o(1) \right) (1-c)^2 \left(\frac{1}{\rho^4} \frac{c^2}{(1-c)^3} + o \left(\frac{1}{\rho^4} \right) \right) + O \left(\frac{1}{n} \right) \\
&= \frac{\eta^4}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right).
\end{aligned}$$

and

$$\mathbb{E} \left[\frac{\eta^3 \|\mathbf{t}\|^2 \xi}{\gamma_1^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \right] = (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left(\mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \right)^2 \mathbb{E} \left[\frac{\xi}{\eta} \right] \mathbb{E} [\|\mathbf{t}\|^2] \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top] + O \left(\frac{1}{n} \right) = O \left(\frac{1}{n} \right).$$

Now we have the expectations and errors for the three terms. Combining them yields the Lemma statement.

For $c > 1$, we recall that $\mathbf{h}\mathbf{s} = \mathbf{0}$, and Proposition 1 implies

$$\begin{aligned}
\boldsymbol{\beta}_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{Z} \boldsymbol{\beta}_* &= \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_2^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + \frac{\eta^4 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\gamma_2^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + \frac{2\eta^3 \|\mathbf{h}\|^2 \xi}{\gamma_2^2} \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{h} \mathbf{s} \mathbf{u}^\top \boldsymbol{\beta}_* \\
&= \left(\frac{\eta^2 \|\mathbf{h}\|^2 (\xi^2 + \eta^2 \|\mathbf{h}\|^2 \|\mathbf{s}\|^2)}{\gamma_2^2} \right) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \\
&= \left(\frac{\eta^2 \|\mathbf{h}\|^2 \gamma_2}{\gamma_2^2} \right) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \\
&= \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2.
\end{aligned}$$

Hence, we can take expectation:

$$\begin{aligned}
\mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{Z} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{Z} \boldsymbol{\beta}_*] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_2} \right] \mathbb{E} [\|\mathbf{h}\|^2] (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + O \left(\frac{1}{n} \right) \\
&= \frac{\eta^2}{\eta^2 + \rho^2} \frac{c}{c-1} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right).
\end{aligned}$$

□

Lemma 27. *In the same setting as Section 2, we have that,*

$$\mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \boldsymbol{\beta}_*] = \begin{cases} \|\boldsymbol{\beta}_*\|^2 + \frac{\eta^2 (\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 - \frac{2\eta^2 c}{\eta^2 c + \rho^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right) & c < 1 \\ \frac{\|\boldsymbol{\beta}_*\|^2}{c} - \frac{\eta^2}{\eta^2 + \rho^2} \left(\frac{\|\boldsymbol{\beta}_*\|^2}{d} - \frac{(\boldsymbol{\beta}_*^\top \mathbf{u})^2}{c(c-1)} \right) + o(1) + O \left(\frac{1}{n} \right) & c > 1 \end{cases}.$$

Proof. We use similar expansions that follow from Lemma 2.

$$\begin{aligned}
\boldsymbol{\beta}_*^\top \mathbf{A} (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \boldsymbol{\beta}_* &= \|\boldsymbol{\beta}_*\|^2 + \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_1^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + \frac{\eta^4 \|\mathbf{t}\|^4}{\gamma_1^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \\
&\quad + \frac{2\eta^3 \|\mathbf{t}\|^2 \xi}{\gamma_1^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top - \frac{2\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \boldsymbol{\beta}_* - \frac{2\eta \xi}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{h} \boldsymbol{\beta}_*.
\end{aligned}$$

Lemma 26 gives the expectation of the first four terms:

$$\|\boldsymbol{\beta}_*\|^2 + \frac{\eta^2 (\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right).$$

We have done the following expectations in Equations 18, 19:

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{h} \boldsymbol{\beta}_* \right] = O \left(\frac{1}{n} \right), \quad \mathbb{E} \left[\frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \boldsymbol{\beta}_* \right] = \frac{\eta^2 c}{\eta^2 c + \rho^2} + o(1) + O \left(\frac{1}{n} \right).$$

Combining these results yields the lemma statement.

For $c > 1$, with $\mathbf{h}\mathbf{s} = \mathbf{0}$, $\mathbf{s}^\top \mathbf{A}\mathbf{A}^\dagger = \mathbf{0}$, $\mathbf{h}\mathbf{A}\mathbf{A}^\dagger = \mathbf{h}$, we have the following expansion by Lemma 2:

$$\begin{aligned} \beta_*^\top \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}\beta_* &= \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \beta_* + \frac{\eta^2 \|\mathbf{s}\|^2 \xi^2}{\gamma_2^2} \beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* + \frac{\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2}{\gamma_2^2} \beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* \\ &\quad + \frac{\eta^4 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\gamma_2^2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{u}^\top \mathbf{A}\mathbf{A}^\dagger \beta_* + \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_2^2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{u}^\top \mathbf{A}\mathbf{A}^\dagger \beta_* \\ &\quad - \frac{2\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* - \frac{2\eta\xi}{\gamma_2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_* \\ &\quad - \frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 \xi}{\gamma_2^2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_* + \frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 \xi}{\gamma_2^2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_*. \end{aligned}$$

We can combine the coefficients as:

$$\begin{aligned} \frac{\eta^2 \|\mathbf{s}\|^2 \xi^2}{\gamma_2^2} + \frac{\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2}{\gamma_2^2} - \frac{2\eta^2 \|\mathbf{s}\|^2}{\gamma_2} &= \frac{\eta^2 \|\mathbf{s}\|^2 (\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2) - 2\eta^2 \|\mathbf{s}\|^2 \gamma_2}{\gamma_2^2} = -\frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2}, \\ \frac{\eta^4 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\gamma_2^2} + \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_2^2} &= \frac{\eta^2 \|\mathbf{h}\|^2 (\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2)}{\gamma_2^2} = \frac{\eta^2 \|\mathbf{h}\|^2 \gamma_2}{\gamma_2^2} = \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2}. \end{aligned}$$

Then we have that:

$$\begin{aligned} &\beta_*^\top \mathbf{A}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}\beta_* \\ &= \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \beta_* - \frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{u}^\top \mathbf{A}\mathbf{A}^\dagger \beta_* - \frac{2\eta\xi}{\gamma_2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_*. \end{aligned}$$

Recall from Equation 13 that $\mathbb{E}[\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \beta_*] = \|\beta_*\|^2/c$. We then proceed similarly with the other expectations using Lemmas 9, 10, 11, 19:

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_* \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_2} \right] \mathbb{E} [\|\mathbf{s}\|^2] \mathbb{E} [\beta_*^\top \mathbf{h}^\top \mathbf{h} \beta_*] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_2} \right)} O \left(\frac{1}{n} \right) \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 + \rho^2} + o \left(\frac{1}{\rho^2} \right) \right) \left(1 - \frac{1}{c} \right) \left(\frac{\|\beta_*\|^2}{d} \frac{c}{\rho^2(c-1)} + o \left(\frac{1}{d\rho^2} \right) \right) + O \left(\frac{1}{n} \right) \\ &= \frac{\|\beta_*\|^2}{d} \frac{\eta^2}{\eta^2 + \rho^2} + o \left(\frac{1}{d} \right) + O \left(\frac{1}{n} \right). \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} (\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u})^2 \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_2} \right] \mathbb{E} [\|\mathbf{h}\|^2] \mathbb{E} [(\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u})^2] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_2} \right)} O \left(\frac{1}{n} \right) \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 + \rho^2} + o \left(\frac{1}{\rho^2} \right) \right) \left(\frac{1}{\rho^2} \frac{c}{c-1} + o \left(\frac{1}{\rho^2} \right) \right) \left(\frac{1}{c^2} (\beta_*^\top \mathbf{u})^2 + o(1) \right) + O \left(\frac{1}{n} \right) \\ &= \frac{\eta^2}{\eta^2 + \rho^2} \frac{(\beta_*^\top \mathbf{u})^2}{c(c-1)} + o(1) + O \left(\frac{1}{n} \right). \end{aligned}$$

$$\begin{aligned} \mathbb{E} \left[\frac{\eta\xi}{\gamma_2} \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_* \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_2} \right] \mathbb{E} \left[\frac{\xi}{\eta} \right] \mathbb{E} [\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h} \beta_*] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_2} \right)} O \left(\frac{1}{n} \right) \\ &= 0 + O \left(\frac{1}{n} \right). \end{aligned} \tag{17}$$

We combine these results to produce the lemma statement. \square

Lemma 28. *In the same setting as Section 2, we have that*

$$\mathbb{E} [\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A}\beta_*] = \begin{cases} - \left(\frac{\eta^2(\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} - \frac{\eta^2 c}{\eta^2 c + \rho^2} \right) (\beta_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right), & c < 1 \\ - \frac{\eta^2}{\eta^2 + \rho^2} \frac{1}{c-1} (\beta_*^\top \mathbf{u})^2 + o(1) + O \left(\frac{1}{n} \right), & c > 1 \end{cases}$$

Proof. For $c < 1$, we expand it using Proposition 1, Lemma 2. Note that all of the relevant expectations have been evaluated in the proofs of Lemmas 26, 27,

$$\begin{aligned} \beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \beta_* &= \frac{\eta \xi}{\gamma_1} \beta_*^\top \mathbf{u} \mathbf{h} \beta_* + \frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \beta_*^\top \mathbf{u} \mathbf{k}^\top \mathbf{A}^\dagger \beta_* - \frac{2\eta^3 \|\mathbf{t}\|^2 \xi}{\gamma_1^2} (\beta_*^\top \mathbf{u})^2 \mathbf{h} \mathbf{A}^{\dagger\top} \mathbf{k} \\ &\quad - \frac{\eta^4 \|\mathbf{t}\|^4}{\gamma_1^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}) (\beta_*^\top \mathbf{u})^2 - \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_1^2} (\beta_*^\top \mathbf{u})^2. \end{aligned}$$

The expectation of the last three terms is given by Lemma 26. The first two expectations come from Equations 18, 19 respectively. We can plug them in and compute the expectation:

$$\mathbb{E} [\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \beta_*] = - \left(\frac{\eta^2 (\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} - \frac{\eta^2 c}{\eta^2 c + \rho^2} \right) (\beta_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right).$$

For $c > 1$, again with $\mathbf{h} \mathbf{s} = \mathbf{0}$ and $\mathbf{s}^\top \mathbf{A} = \mathbf{0}$, $\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \beta_*$ becomes:

$$\begin{aligned} &\beta_*^\top \frac{\eta \xi}{\gamma_2} \mathbf{u} \mathbf{h} \left(\mathbf{A} \mathbf{A}^\dagger + \frac{\eta \xi}{\gamma_2} \mathbf{s} \mathbf{h} - \frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \mathbf{h}^\top \mathbf{h} - \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{s} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger - \frac{\eta \xi}{\gamma_2} \mathbf{h}^\top \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \right) \beta_* \\ &\quad + \beta_*^\top \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top \left(\mathbf{A} \mathbf{A}^\dagger + \frac{\eta \xi}{\gamma_2} \mathbf{s} \mathbf{h} - \frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \mathbf{h}^\top \mathbf{h} - \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{s} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger - \frac{\eta \xi}{\gamma_2} \mathbf{h}^\top \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \right) \beta_* \\ &= \beta_*^\top \left[\frac{\eta \xi}{\gamma_2} \mathbf{u} \mathbf{h} \mathbf{A} \mathbf{A}^\dagger - \frac{\eta^3 \xi \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2^2} \mathbf{u} \mathbf{h} - \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_2^2} \mathbf{u} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \right] \beta_* \\ &\quad + \beta_*^\top \left[\frac{\eta^3 \|\mathbf{h}\|^2 \|\mathbf{s}\|^2 \xi}{\gamma_2^2} \mathbf{u} \mathbf{h} - \frac{\eta^4 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\gamma_2^2} \mathbf{u} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \right] \beta_* \\ &= \beta_*^\top \left[\frac{\eta \xi}{\gamma_2} \mathbf{u} \mathbf{h} \mathbf{A} \mathbf{A}^\dagger - \frac{\eta^2 \|\mathbf{h}\|^2 \xi^2}{\gamma_2^2} \mathbf{u} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger - \frac{\eta^4 \|\mathbf{h}\|^4 \|\mathbf{s}\|^2}{\gamma_2^2} \mathbf{u} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \right] \beta_* \\ &= (\beta_*^\top \mathbf{u}) \left(\frac{\eta \xi}{\gamma_2} \mathbf{h} \mathbf{A} \mathbf{A}^\dagger \beta_* - \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \beta_* \right) \quad \text{since } \gamma_2 = \eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 + \xi^2. \end{aligned}$$

We need to evaluate two following expectations. Similar to $c < 1$,

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_2} \mathbf{h} \mathbf{A} \mathbf{A}^\dagger \beta_* \right] = O\left(\frac{1}{n}\right).$$

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \beta_* \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_2} \right] \mathbb{E} [\|\mathbf{h}\|^2] \mathbb{E} [\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{u}] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_2} \right)} O\left(\frac{1}{n}\right) \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 + \rho^2} + o\left(\frac{1}{\rho^2}\right) \right) \left(\frac{1}{\rho^2} \frac{c}{c-1} + o\left(\frac{1}{\rho^2}\right) \right) \left(\frac{1}{c} (\beta_*^\top \mathbf{u}) \right) + O\left(\frac{1}{n}\right) \\ &= \frac{\eta^2}{\eta^2 + \rho^2} \frac{(\beta_*^\top \mathbf{u})}{c-1} + o(1) + O\left(\frac{1}{n}\right). \end{aligned}$$

Finally, we have that:

$$\mathbb{E} [\beta_*^\top \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{A} \beta_*] = - \frac{\eta^2}{\eta^2 + \rho^2} \frac{1}{c-1} (\beta_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right).$$

□

Lemma 29. *In the same setting as Section 2, we have that,*

$$\mathbb{E} [\boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \boldsymbol{\varepsilon}] = \begin{cases} \tau_\varepsilon^2 \left(\frac{cd}{\rho^2(1-c)} - \frac{\eta^2}{\rho^2(\eta^2 c + \rho^2)} \frac{c^2}{1-c} \right) + o\left(\frac{\eta}{\rho^2}\right) + O\left(\frac{1}{\rho^2 n}\right), & c < 1 \\ \tau_\varepsilon^2 \left(\frac{d}{\rho^2(c-1)} - \frac{\eta^2}{\rho^2(\eta^2 + \rho^2)} \frac{c}{c-1} \right) + o\left(\frac{\eta}{\rho^2}\right) + O\left(\frac{1}{\rho^2 n}\right), & c > 1 \end{cases}$$

Proof. For $c < 1$, we first expand this term using Theorem 6:

$$\begin{aligned}\varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \varepsilon &= \varepsilon^\top \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger - \frac{\xi}{\gamma_1} \mathbf{p}_1 \mathbf{q}_1^\top \right) \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger - \frac{\xi}{\gamma_1} \mathbf{p}_1 \mathbf{q}_1^\top \right)^\top \varepsilon \\ &= \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \varepsilon + \frac{2\eta}{\xi} \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} \mathbf{t} \varepsilon - \frac{2\xi}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \mathbf{q}_1 \mathbf{p}_1^\top \varepsilon \\ &\quad + \frac{\eta^2}{\xi^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k}) \varepsilon^\top \mathbf{t}^\top \mathbf{t} \varepsilon - \frac{2\eta}{\gamma_1} \varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{q}_1 \mathbf{p}_1^\top \varepsilon + \frac{\xi^2}{\gamma_1^2} \varepsilon^\top \mathbf{p}_1 \mathbf{q}_1^\top \mathbf{q}_1 \mathbf{p}_1^\top \varepsilon\end{aligned}$$

Note that Lemma 21 and the fact that $\mathbf{t} \mathbf{A}^\dagger = \mathbf{0}$ imply that the second term has zero expectation:

$$\mathbb{E}_\varepsilon \left[\frac{2\eta}{\xi} \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} \mathbf{t} \varepsilon \right] = \frac{2\eta \tau_\varepsilon^2}{\xi} \mathbf{t} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} = 0.$$

Similarly, we will later use:

$$\mathbb{E}_\varepsilon [\varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{t} \varepsilon] = \tau_\varepsilon^2 \mathbf{t} \mathbf{A}^\dagger \mathbf{h}^\top = 0, \quad \mathbb{E}_\varepsilon [\varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \varepsilon] = \tau_\varepsilon^2 \text{Tr}(\mathbf{t}^\top \mathbf{k}^\top) = \tau_\varepsilon^2 \text{Tr}(\mathbf{k} \mathbf{t}) = 0.$$

Note that these equalities are exact without taking the expectation over other sources of randomness besides ε .

We now expand the other terms one by one and compute their expectations along the way. We start by eliminating zero expectations and taking expectations w.r.t. ε using Lemma 21.

$$\begin{aligned}\mathbb{E} \left[\frac{\eta^2}{\xi^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k}) \varepsilon^\top \mathbf{t}^\top \mathbf{t} \varepsilon \right] &= \mathbb{E} \left[\frac{\eta^2 \|\mathbf{t}\|^2 \tau_\varepsilon^2}{\xi^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} \right]. \\ \mathbb{E} \left[-\frac{2\xi}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \mathbf{q}_1 \mathbf{p}_1^\top \varepsilon \right] &= \mathbb{E} \left[-\frac{2\xi}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \left(\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{k} + \mathbf{h}^\top \right) \left(\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t} + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^3 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi} \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} \mathbf{t} \varepsilon - \frac{2\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} \mathbf{k}^\top \varepsilon \right. \\ &\quad \left. - \frac{2\eta^2 \|\mathbf{k}\|^2}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{t} \varepsilon - \frac{2\eta \xi}{\gamma_1} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{k}^\top \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^2 \|\mathbf{t}\|^2 \tau_\varepsilon^2}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} - \frac{2\eta \xi \tau_\varepsilon^2}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right]. \\ \mathbb{E} \left[-\frac{2\eta}{\gamma_1} \varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{q}_1 \mathbf{p}_1^\top \varepsilon \right] &= \mathbb{E} \left[-\frac{2\eta}{\gamma_1} \varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \mathbf{A}^\dagger \left(\frac{\eta \|\mathbf{t}\|^2}{\xi} \mathbf{A}^\dagger \mathbf{k} + \mathbf{h}^\top \right) \left(\frac{\eta^2 \|\mathbf{k}\|^2}{\xi} \mathbf{t} + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^4 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k}) \varepsilon^\top \mathbf{t}^\top \mathbf{t} \varepsilon - \frac{2\eta^3 \|\mathbf{k}\|^2}{\gamma_1 \xi} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top) \varepsilon^\top \mathbf{t}^\top \mathbf{t} \varepsilon \right. \\ &\quad \left. - \frac{2\eta^3 \|\mathbf{t}\|^2}{\gamma_1 \xi} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k}) \varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \varepsilon - \frac{2\eta^2}{\gamma_1} (\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top) \varepsilon^\top \mathbf{t}^\top \mathbf{k}^\top \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^4 \|\mathbf{t}\|^4 \|\mathbf{k}\|^2 \tau_\varepsilon^2}{\gamma_1 \xi^2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{k} - \frac{2\eta^3 \|\mathbf{k}\|^2 \|\mathbf{t}\|^2 \tau_\varepsilon^2}{\gamma_1 \xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right].\end{aligned}$$

By the squared norms in Lemmas 6, 7, and Lemma 21,

$$\begin{aligned}
\mathbb{E} \left[\frac{\xi^2}{\gamma_1^2} \boldsymbol{\varepsilon}^\top \mathbf{p}_1 \mathbf{q}_1^\top \mathbf{q}_1 \mathbf{p}_1^\top \boldsymbol{\varepsilon} \right] &= \frac{\xi^2 \tau_\varepsilon^2}{\gamma_1^2} \|\mathbf{p}_1\|^2 \|\mathbf{q}_1\|^2 \\
&= \frac{\xi^2 \tau_\varepsilon^2}{\gamma_1^2} \left(\frac{\eta^2 \|\mathbf{k}\|^2}{\xi^2} \gamma_1 \right) \left(\frac{\eta^2 \|\mathbf{t}\|^4}{\xi^2} \mathbf{k} \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} + \frac{2\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \|\mathbf{h}\|^2 \right) \\
&= \frac{\tau_\varepsilon^2}{\gamma_1} (\eta^2 \|\mathbf{k}\|^2) \left(\frac{\eta^2 \|\mathbf{t}\|^4}{\xi^2} \mathbf{k} \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} + \frac{2\eta \|\mathbf{t}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \|\mathbf{h}\|^2 \right) \\
&= \tau_\varepsilon^2 \left(\frac{\eta^4 \|\mathbf{t}\|^4 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} \mathbf{k} \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} + \frac{2\eta^3 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \frac{\eta^2 \|\mathbf{k}\|^2 \|\mathbf{h}\|^2}{\gamma_1} \right)
\end{aligned}$$

We combine like terms and simplify the coefficients. For the term $\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}$,

$$\begin{aligned}
\tau_\varepsilon^2 \left(\frac{\eta^4 \|\mathbf{t}\|^4 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} - \frac{2\eta^4 \|\mathbf{t}\|^4 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} - \frac{2\eta^2 \|\mathbf{t}\|^2}{\gamma_1} + \frac{\eta^2 \|\mathbf{t}\|^2}{\xi^2} \right) &= \tau_\varepsilon^2 \eta^2 \|\mathbf{t}\|^2 \left(\frac{\eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} - \frac{2\eta^2 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi^2} - \frac{2}{\gamma_1} + \frac{1}{\xi^2} \right) \\
&= \tau_\varepsilon^2 \eta^2 \|\mathbf{t}\|^2 \left(-\frac{\gamma_1 - \xi^2}{\gamma_1 \xi^2} - \frac{2}{\gamma_1} + \frac{1}{\xi^2} \right) \\
&= \tau_\varepsilon^2 \eta^2 \|\mathbf{t}\|^2 \left(-\frac{\gamma_1 - \xi^2}{\gamma_1 \xi^2} - \frac{2\xi^2}{\gamma_1 \xi^2} + \frac{\gamma_1}{\gamma_1 \xi^2} \right) \\
&= -\tau_\varepsilon^2 \frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1}.
\end{aligned}$$

For the term $\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top$,

$$\tau_\varepsilon^2 \left(\frac{2\eta^3 \|\mathbf{t}\|^2 \|\mathbf{k}\|^2}{\gamma_1 \xi} - \frac{2\eta^3 \|\mathbf{k}\|^2 \|\mathbf{t}\|^2}{\gamma_1 \xi} - \frac{2\eta \xi}{\gamma_1} \right) = -\tau_\varepsilon^2 \frac{2\eta \xi}{\gamma_1}.$$

Combining these terms together, we have:

$$\mathbb{E} [\boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger\top} \boldsymbol{\varepsilon}] = \mathbb{E} \left[\boldsymbol{\varepsilon}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \boldsymbol{\varepsilon} - \frac{\eta^2 \|\mathbf{t}\|^2 \tau_\varepsilon^2}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} - \frac{2\eta \xi \tau_\varepsilon^2}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \frac{\eta^2 \|\mathbf{k}\|^2 \|\mathbf{h}\|^2}{\gamma_1} \right].$$

Similarly, using Lemmas 9, 10, 11, 19, 21, we have the following:

$$\mathbb{E} [\boldsymbol{\varepsilon}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \boldsymbol{\varepsilon}] = \tau_\varepsilon^2 \mathbb{E} [\text{Tr}(\mathbf{A}^\dagger \mathbf{A}^{\dagger\top})] = \tau_\varepsilon^2 n \mathbb{E} \left[\frac{1}{\lambda} \right] = \tau_\varepsilon^2 \frac{cd}{\rho^2(1-c)} + o\left(\frac{d}{\rho^2}\right) \quad \text{by Equation 11.}$$

$$\begin{aligned}
\mathbb{E} \left[\frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k} \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} [\|\mathbf{t}\|^2] \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger\top} \mathbf{k}] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_1} \right)} O\left(\frac{1}{n}\right) \\
&= \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right) \right) (1-c) \left(\frac{1}{\rho^4} \frac{c^2}{(1-c)^3} + o\left(\frac{1}{\rho^4}\right) \right) + O\left(\frac{1}{n}\right) \\
&= \frac{\eta^2}{\eta^2 c + \rho^2} \frac{c^2}{\rho^2(1-c)^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right).
\end{aligned}$$

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_1} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right] = \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} \left[\frac{\xi}{\eta} \right] \mathbb{E} [\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_1} \right)} O\left(\frac{1}{n}\right) = O\left(\frac{1}{n}\right).$$

$$\begin{aligned}
\mathbb{E} \left[\frac{\eta^2 \|\mathbf{k}\|^2 \|\mathbf{h}\|^2}{\gamma_1} \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} [\|\mathbf{k}\|^2] \mathbb{E} [\|\mathbf{h}\|^2] + \sqrt{\text{Var} \left(\frac{\eta^2}{\gamma_1} \right)} O\left(\frac{1}{n}\right) \\
&= \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o\left(\frac{1}{\rho^2}\right) \right) \left(\frac{1}{\rho^2} \frac{c^2}{1-c} + o\left(\frac{1}{\rho^2}\right) \right) \left(\frac{1}{\rho^2} \frac{c}{1-c} + o\left(\frac{1}{\rho^2}\right) \right) + O\left(\frac{1}{n}\right) \\
&= \frac{\eta^2}{\eta^2 c + \rho^2} \frac{c^3}{\rho^2(1-c)^2} + o(1) + O\left(\frac{1}{n}\right).
\end{aligned}$$

After simple algebra, the result follows from here.

For $c > 1$, we can expand similarly using Theorem 6,

$$\begin{aligned} \varepsilon^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^\dagger \varepsilon &= \varepsilon^\top \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top - \frac{\xi}{\gamma_2} \mathbf{p}_2 \mathbf{q}_2^\top \right) \left(\mathbf{A}^\dagger + \frac{\eta}{\xi} \mathbf{s} \mathbf{h} \mathbf{A}^\dagger - \frac{\xi}{\gamma_2} \mathbf{q}_2 \mathbf{p}_2^\top \right) \varepsilon \\ &= \varepsilon^\top \mathbf{A}^\dagger \mathbf{A}^\dagger \varepsilon + \frac{2\eta}{\xi} \varepsilon^\top \underbrace{\mathbf{A}^\dagger \mathbf{s} \mathbf{h} \mathbf{A}^\dagger}_0 \varepsilon - \frac{2\xi}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon \\ &\quad + \frac{\eta^2 \|\mathbf{s}\|^2}{\xi^2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{h} \mathbf{A}^\dagger \varepsilon - \frac{2\eta}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon + \frac{\xi^2}{\gamma_2^2} \varepsilon^\top \mathbf{p}_2 \mathbf{q}_2^\top \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon. \end{aligned}$$

We expand the other terms one by one, marking those with zero expectations:

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \|\mathbf{s}\|^2}{\xi^2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{h} \mathbf{A}^\dagger \varepsilon \right] &= \mathbb{E} \left[\frac{\eta^2 \|\mathbf{s}\|^2 \tau_\varepsilon^2}{\xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top \right]. \\ \mathbb{E} \left[-\frac{2\xi}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon \right] &= \mathbb{E} \left[-\frac{2\xi}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \left(\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s} + \mathbf{h}^\top \right) \left(\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h} \mathbf{A}^\dagger + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\xi}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \left(\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h} \mathbf{A}^\dagger + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{h} \mathbf{A}^\dagger \varepsilon - \frac{2\eta\xi}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{k}^\top \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^2 \|\mathbf{s}\|^2 \tau_\varepsilon^2}{\gamma_2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top - \frac{2\eta\xi \tau_\varepsilon^2}{\gamma_2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right]. \\ \mathbb{E} \left[-\frac{2\eta}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon \right] &= \mathbb{E} \left[-\frac{2\eta}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{s}^\top \left(\frac{\eta \|\mathbf{h}\|^2}{\xi} \mathbf{s} + \mathbf{h}^\top \right) \left(\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h} \mathbf{A}^\dagger + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta}{\gamma_2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \left(\frac{\eta \|\mathbf{h}\|^2 \|\mathbf{s}\|^2}{\xi} \right) \left(\frac{\eta^2 \|\mathbf{s}\|^2}{\xi} \mathbf{h} \mathbf{A}^\dagger + \eta \mathbf{k}^\top \right) \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2}{\gamma_2 \xi^2} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{h} \mathbf{A}^\dagger \varepsilon - \frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2 \xi} \varepsilon^\top \mathbf{A}^\dagger \mathbf{h}^\top \mathbf{k}^\top \varepsilon \right] \\ &= \mathbb{E} \left[-\frac{2\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2 \tau_\varepsilon^2}{\gamma_2 \xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top - \frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2 \tau_\varepsilon^2}{\gamma_2 \xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right]. \end{aligned}$$

Using the squared norms from Lemmas 6, 7,

$$\begin{aligned} \mathbb{E} \left[\frac{\xi^2}{\gamma_2^2} \varepsilon^\top \mathbf{p}_2 \mathbf{q}_2^\top \mathbf{q}_2 \mathbf{p}_2^\top \varepsilon \right] &= \mathbb{E} \left[\frac{\xi^2}{\gamma_2^2} \tau_\varepsilon^2 \|\mathbf{p}_2\|^2 \|\mathbf{q}_2\|^2 \right] \\ &= \mathbb{E} \left[\frac{\xi^2 \tau_\varepsilon^2}{\gamma_2^2} \left(\frac{\|\mathbf{h}\|^2}{\xi^2} \gamma_2 \right) \left(\frac{\eta^4 \|\mathbf{s}\|^4}{\xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top + \frac{2\eta^3 \|\mathbf{s}\|^2}{\xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \eta^2 \|\mathbf{k}\|^2 \right) \right] \\ &= \mathbb{E} \left[\tau_\varepsilon^2 \left(\frac{\eta^4 \|\mathbf{h}\|^2 \|\mathbf{s}\|^4}{\gamma_2 \xi^2} \mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top + \frac{2\eta^3 \|\mathbf{h}\|^2 \|\mathbf{s}\|^2}{\gamma_2 \xi} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \frac{\eta^2 \|\mathbf{h}\|^2 \|\mathbf{k}\|^2}{\gamma_2} \right) \right]. \end{aligned}$$

Similarly, we combine the coefficients: For the term $\mathbf{h} \mathbf{A}^\dagger \mathbf{A}^\dagger \mathbf{h}^\top$,

$$\begin{aligned} \tau_\varepsilon^2 \left(\frac{\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2}{\gamma_2 \xi^2} - \frac{2\eta^4 \|\mathbf{s}\|^4 \|\mathbf{h}\|^2}{\gamma_2 \xi^2} - \frac{2\eta^2 \|\mathbf{s}\|^2}{\gamma_2} + \frac{\eta^2 \|\mathbf{s}\|^2}{\xi^2} \right) &= \tau_\varepsilon^2 \eta^2 \|\mathbf{s}\|^2 \left(\frac{\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2 \xi^2} - \frac{2\eta^2 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2 \xi^2} - \frac{2}{\gamma_2} + \frac{1}{\xi^2} \right) \\ &= \tau_\varepsilon^2 \eta^2 \|\mathbf{s}\|^2 \left(-\frac{\gamma_2 - \xi^2}{\gamma_2 \xi^2} - \frac{2}{\gamma_2} + \frac{1}{\xi^2} \right) \\ &= \tau_\varepsilon^2 \eta^2 \|\mathbf{s}\|^2 \left(-\frac{\gamma_2 - \xi^2}{\gamma_2 \xi^2} - \frac{2\xi^2}{\gamma_2 \xi^2} + \frac{\gamma_2}{\gamma_2 \xi^2} \right) \\ &= -\tau_\varepsilon^2 \frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2}. \end{aligned}$$

For the term $\mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top$,

$$\tau_\varepsilon^2 \left(\frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2 \xi} - \frac{2\eta^3 \|\mathbf{s}\|^2 \|\mathbf{h}\|^2}{\gamma_2 \xi} - \frac{2\eta \xi}{\gamma_2} \right) = -\tau_\varepsilon^2 \frac{2\eta \xi}{\gamma_2}.$$

Combining these terms together, we have:

$$\mathbb{E} [\boldsymbol{\varepsilon}^\top (\mathbf{Z} + \mathbf{A})^\dagger (\mathbf{Z} + \mathbf{A})^{\dagger \top} \boldsymbol{\varepsilon}] = \mathbb{E} \left[\boldsymbol{\varepsilon}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger \top} \boldsymbol{\varepsilon} - \frac{\eta^2 \|\mathbf{s}\|^2 \tau_\varepsilon^2}{\gamma_2} \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^\dagger \mathbf{h}^\top - \frac{2\eta \xi \tau_\varepsilon^2}{\gamma_2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top + \frac{\eta^2 \|\mathbf{k}\|^2 \|\mathbf{h}\|^2}{\gamma_2} \right].$$

Similarly, replicating the proof with the $c > 1$ counterparts, we have the following:

$$\mathbb{E} [\boldsymbol{\varepsilon}^\top \mathbf{A}^\dagger \mathbf{A}^{\dagger \top} \boldsymbol{\varepsilon}] = \tau_\varepsilon^2 \mathbb{E} [\text{Tr}(\mathbf{A}^\dagger \mathbf{A}^{\dagger \top})] = \tau_\varepsilon^2 n \mathbb{E} \left[\frac{1}{\lambda} \right] = \tau_\varepsilon^2 \frac{d}{\rho^2(c-1)} + o\left(\frac{d}{\rho^2}\right).$$

$$\mathbb{E} \left[\frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \mathbf{h} \mathbf{A}^{\dagger \top} \mathbf{A}^\dagger \mathbf{h}^\top \right] = \frac{\eta^2}{\eta^2 + \rho^2} \frac{c^2}{\rho^2(c-1)^2} + o(1) + O\left(\frac{1}{n}\right).$$

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_2} \mathbf{k}^\top \mathbf{A}^\dagger \mathbf{h}^\top \right] = O\left(\frac{1}{n}\right).$$

$$\mathbb{E} \left[\frac{\eta^2 \|\mathbf{k}\|^2 \|\mathbf{h}\|^2}{\gamma_2} \right] = \frac{\eta^2}{\eta^2 + \rho^2} \frac{c}{\rho^2(c-1)^2} + o(1) + O\left(\frac{1}{n}\right).$$

After simple algebra, the result follows. \square

E.5.7 TARGET ALIGNMENT: HELPER LEMMAS

Lemma 30. *In the same setting as Section 2, we have that*

$$\mathbb{E} [\boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^{\dagger \top} \mathbf{Z}^\top \boldsymbol{\beta}_*] = \begin{cases} \frac{\eta^2 c}{\rho^2 + \eta^2 c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right) & c < 1 \\ \frac{\eta^2}{\eta^2 + \rho^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right) & c > 1 \end{cases}.$$

Proof. For $c < 1$, from Proposition 1, we get that

$$\boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^{\dagger \top} \mathbf{Z}^\top \boldsymbol{\beta}_* = \frac{\eta \xi}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* + \frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger \top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_*.$$

To begin, we start estimating

$$\mathbb{E} \left[\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right].$$

Using our Spherical Hypercontractivity, we have that $\frac{\xi}{\eta}$ and $\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_*$ satisfy the assumptions for Lemma 36. Then using Lemma 9 we have that

$$\mathbb{E} \left[\frac{\xi}{\eta} \right] = \frac{1}{\eta} \quad \text{and} \quad \text{Var} \left(\frac{1}{\eta} \right) = O\left(\frac{1}{\rho^2 d}\right)$$

and Lemma 11, we have that

$$\mathbb{E} [\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_*] = 0 \quad \text{and} \quad \text{Var} (\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_*) = O\left(\frac{1}{\rho^2 d}\right)$$

Thus, using Lemma 37, we have that

$$\mathbb{E} \left[\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] = 0 + O\left(\frac{1}{\rho^2 d}\right)$$

and using Lemma 36, since all the means are $O(1)$, we have that

$$\text{Var} \left(\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right) = O\left(\max\left(\text{Var}\left(\frac{\xi}{\eta}\right), \text{Var}(\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_*)\right)\right) = O\left(\frac{1}{\rho^2 n}\right).$$

Then Lemma 19 gives mean and variance of $\frac{\eta^2}{\gamma_i}$. Since $\frac{\eta^2}{\gamma_i}$ does not satisfy the higher moment bound, and cannot be directly included in the product, we can include it via the classical bound:

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] = \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} \left[\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] + \sqrt{\text{Var} \left(\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right) \text{Var} \left(\frac{\eta^2}{\gamma_1} \right)} = O \left(\frac{1}{n} \right). \quad (18)$$

For the second term, we begin with

$$\mathbb{E} \left[\|\mathbf{t}\|^2 \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right].$$

Lemma 9 tells us that

$$\mathbb{E}[\|\mathbf{t}\|^2] = 1 - c \quad \text{and} \quad \text{Var}(\|\mathbf{t}\|^2) = O \left(\frac{1}{n} \right)$$

and Lemma 10 tells us

$$\mathbb{E} \left[\boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right] = \frac{1}{\rho^2} \frac{c}{1-c} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + o \left(\frac{1}{\rho^2} \right) \quad \text{and} \quad \text{Var} \left(\boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right) = O \left(\frac{1}{\rho^4 d} \right).$$

Thus using Lemmas 37 and Lemma 36, we get that

$$\mathbb{E} \left[\|\mathbf{t}\|^2 \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right] = (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{c}{\rho^2} + o \left(\frac{1}{\rho^2} \right) + O \left(\frac{1}{n} \right) \quad \text{and} \quad \text{Var} \left(\|\mathbf{t}\|^2 \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right) = O \left(\frac{1}{n} \right)$$

Recalling the mean and variance for $\frac{\eta^2}{\gamma_1}$ from 19, we have that

$$\begin{aligned} \mathbb{E} \left[\frac{\eta^2 \|\mathbf{t}\|^2}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right] &= \mathbb{E} \left[\frac{\eta^2}{\gamma_1} \right] \mathbb{E} \left[\|\mathbf{t}\|^2 \boldsymbol{\beta}_*^\top \mathbf{A}^{\dagger\top} \mathbf{k} \mathbf{u}^\top \boldsymbol{\beta}_* \right] + \sqrt{O \left(\frac{1}{n} \right) \text{Var} \left(\frac{\eta^2}{\gamma_1} \right)} \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 c + \rho^2} + o \left(\frac{1}{\rho^2} \right) \right) \left((\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{c}{\rho^2} + o \left(\frac{1}{\rho^2} \right) + O \left(\frac{1}{n} \right) \right) + O \left(\frac{1}{n} \right) \\ &= (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{\eta^2 c}{\eta^2 c + \rho^2} + o(1) + O \left(\frac{1}{n} \right). \end{aligned} \quad (19)$$

Combining these two terms yields the first result.

Similarly, for $c > 1$, Proposition 1 gives the expansion:

$$\boldsymbol{\beta}_*^\top (\mathbf{Z} + \mathbf{A})^{\dagger\top} \mathbf{Z}^\top \boldsymbol{\beta}_* = \boldsymbol{\beta}_*^\top \left(\frac{\eta \xi}{\gamma_2} \mathbf{u} \mathbf{h} + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{u} \mathbf{s}^\top \right)^\top \boldsymbol{\beta}_* = \frac{\eta \xi}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* + \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \boldsymbol{\beta}_*^\top \mathbf{s} \mathbf{u}^\top \boldsymbol{\beta}_*.$$

For the first term, we begin with

$$\mathbb{E} \left[\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right].$$

Recalling from Lemma 11, we see that

$$\mathbb{E} \left[\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] = 0 \quad \text{and} \quad \text{Var} \left(\boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right) = O \left(\frac{1}{\rho^2 d} \right).$$

Thus again using Lemma 36 and Lemma 37, we see that

$$\mathbb{E} \left[\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] = 0 + O \left(\frac{1}{\rho^2 d} \right) \quad \text{and} \quad \text{Var} \left(\frac{\xi}{\eta} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right) = O \left(\frac{1}{\rho^2 d} \right).$$

Next using the standard covariance bound on the expectation of the product. We see that

$$\mathbb{E} \left[\frac{\eta \xi}{\gamma_1} \boldsymbol{\beta}_*^\top \mathbf{h}^\top \mathbf{u}^\top \boldsymbol{\beta}_* \right] = 0 + O \left(\frac{1}{\rho^2 d} \right) + O \left(\frac{1}{n} \right) = O \left(\frac{1}{n} \right).$$

For the second term, we begin with

$$\mathbb{E} [\|\mathbf{h}\|^2 \beta_*^\top \mathbf{s} \mathbf{u} \beta_*].$$

Recall from Lemma 9 we have that

$$\mathbb{E}[\|\mathbf{h}\|^2] = \frac{1}{\rho^2} \frac{c}{c-1} + o\left(\frac{1}{\rho^2}\right) \quad \text{and} \quad \text{Var}(\|\mathbf{h}\|^2) = O\left(\frac{1}{\rho^4 n}\right)$$

and from Lemma 10

$$\mathbb{E}[\beta_*^\top \mathbf{s} \mathbf{u} \beta_*] = \left(1 - \frac{1}{c}\right) (\beta_*^\top \mathbf{u})^2 \quad \text{and} \quad \text{Var}(\beta_*^\top \mathbf{s} \mathbf{u} \beta_*) = O\left(\frac{1}{d}\right).$$

Thus using Lemma 36 and Lemma 37, we get that

$$\mathbb{E}[\|\mathbf{h}\|^2 \beta_*^\top \mathbf{s} \mathbf{u} \beta_*] = \frac{(\beta_*^\top \mathbf{u})^2}{\rho^2} + o\left(\frac{1}{\rho^2}\right) + O\left(\frac{1}{d}\right) \quad \text{and} \quad \text{Var}(\|\mathbf{h}\|^2 \beta_*^\top \mathbf{s} \mathbf{u} \beta_*) = O\left(\frac{1}{d}\right).$$

Recalling the mean and variance for $\frac{\eta^2}{\gamma_2}$ from Lemma 19 and using the classical covariance bound for the expectation of the product, we get that

$$\begin{aligned} \mathbb{E}\left[\frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \beta_*^\top \mathbf{s} \mathbf{u}^\top \beta_*\right] &= \mathbb{E}\left[\frac{\eta^2}{\gamma_2}\right] \mathbb{E}[\|\mathbf{h}\|^2 \beta_*^\top \mathbf{s} \mathbf{u}^\top \beta_*] + \sqrt{O\left(\frac{1}{n}\right) \text{Var}\left(\frac{\eta^2}{\gamma_2}\right)} \\ &= \left(\frac{\rho^2 \eta^2}{\eta^2 + \rho^2} + o\left(\frac{1}{\rho^2}\right)\right) \left(\frac{(\beta_*^\top \mathbf{u})^2}{\rho^2} + o\left(\frac{1}{\rho^2}\right) + O\left(\frac{1}{d}\right)\right) + O\left(\frac{1}{n}\right) \\ &= \frac{\eta^2}{\eta^2 + \rho^2} (\beta_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right). \end{aligned}$$

Then adding the two together, we get the result for $c > 1$ as well. \square

Lemma 31. *In the same setting as Section 2, we have that, for $c < 1$*

$$\mathbb{E}[\beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_*] = \|\beta_*\|^2 - \frac{\eta^2 c}{\rho^2 + \eta^2 c} (\beta_*^\top \mathbf{u})^2 + o\left(\frac{1}{\rho^2}\right) + O\left(\frac{1}{n}\right).$$

and for $c > 1$

$$\mathbb{E}[\beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_*] = \frac{1}{c} \|\beta_*\|^2 - \frac{\eta^2}{\eta^2 + \rho^2} \left(\frac{\|\beta_*\|^2}{d} + \frac{1}{c} (\beta_*^\top \mathbf{u})^2\right) + o(1) + O\left(\frac{1}{n}\right).$$

Proof. For $c < 1$, using the expectation from Lemma 30, we get

$$\begin{aligned} \mathbb{E}[\beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_*] &= \mathbb{E}\left[\beta_*^\top (\mathbf{I} - \mathbf{Z}(\mathbf{Z} + \mathbf{A})^\dagger)^\top \beta_*\right] \\ &= \|\beta_*\|^2 - \frac{\eta^2 c}{\rho^2 + \eta^2 c} (\beta_*^\top \mathbf{u})^2 + o(1) + O\left(\frac{1}{n}\right). \end{aligned}$$

For $c > 1$, using Lemma 2, we get

$$\beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_* = \beta_*^\top \left(\mathbf{A} \mathbf{A}^\dagger + \frac{\eta \xi}{\gamma_2} \mathbf{h}^\top \mathbf{s}^\top - \frac{\eta^2 \|\mathbf{s}\|^2}{\gamma_2} \mathbf{h}^\top \mathbf{h} - \frac{\eta^2 \|\mathbf{h}\|^2}{\gamma_2} \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{s}^\top - \frac{\eta \xi}{\gamma_2} \mathbf{A} \mathbf{A}^\dagger \mathbf{u} \mathbf{h} \right)^\top \beta_*.$$

We then compute the expectation of each term above. To begin, we have that

$$\mathbb{E}[\beta_*^\top \mathbf{A} \mathbf{A}^\dagger \beta_*] = \frac{1}{c} \|\beta_*\|^2 \quad \text{by Equation 13.}$$

Next, we recall from Lemma 11 that

$$\mathbb{E}[\beta_*^\top \mathbf{h}^\top \mathbf{s}^\top \beta_*] = 0 \quad \text{and} \quad \text{Var}(\beta_*^\top \mathbf{h}^\top \mathbf{s}^\top \beta_*) = O\left(\frac{1}{\rho^2 d}\right).$$

and from Lemma 9 that

$$\mathbb{E}\left[\frac{\xi}{\eta}\right] = \frac{1}{\eta} + o\left(\frac{1}{\rho^2}\right) \quad \text{and} \quad \text{Var}\left(\frac{\xi}{\eta}\right) = O\left(\frac{1}{\rho^2 n}\right)$$

Thus, using Lemmas 36 and Lemma 37, we have that

$$\mathbb{E}\left[\frac{\xi}{\eta}\beta_*^\top \mathbf{h}^\top \mathbf{s}^\top \beta_*\right] = O\left(\frac{1}{\rho^2 n}\right) \quad \text{and} \quad \text{Var}\left(\frac{\xi}{\eta}\beta_*^\top \mathbf{h}^\top \mathbf{s}^\top \beta_*\right) = O\left(\frac{1}{\rho^2 n}\right).$$

Then recalling the mean and variance of η^2/γ_2 from 19, using the standard covariance bound on the difference between the product of the expectation and the expectation of the product, we get that

$$\mathbb{E}\left[\frac{\eta\xi}{\gamma_2}\beta_*^\top \mathbf{h}^\top \mathbf{s}^\top \beta_*\right] = O\left(\frac{1}{n}\right) \quad \text{and} \quad \mathbb{E}\left[\frac{\eta\xi}{\gamma_2}\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\beta_*\right] = O\left(\frac{1}{n}\right).$$

Furthermore, for the next three terms, recall from Lemma 10 that

$$\mathbb{E}[\beta_*^\top \mathbf{h}^\top \mathbf{h}\beta_*] = \frac{\|\beta_*\|^2}{d} \frac{c}{\rho^2(c-1)} + o\left(\frac{1}{\rho^2 d}\right) \quad \text{and} \quad \text{Var}(\beta_*^\top \mathbf{h}^\top \mathbf{h}\beta_*) = O\left(\frac{1}{\rho^2 d^2}\right)$$

and

$$\mathbb{E}[\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top \beta_*] = \frac{c-1}{c^2}(\beta_*^\top \mathbf{u})^2 + o(1) \quad \text{and} \quad \text{Var}(\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top \beta_*) = O\left(\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top \beta_* \frac{1}{d}\right)$$

and from Lemma 11

$$\mathbb{E}[\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\beta_*] = 0 \quad \text{and} \quad \text{Var}(\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\beta_*) = O\left(\frac{1}{\rho^2 d^2}\right).$$

Then recalling from Lemma 9, we have that

$$\mathbb{E}[\|\mathbf{s}\|^2] = 1 - \frac{1}{c} \quad \text{and} \quad \text{Var}(\|\mathbf{s}\|^2) = O\left(\frac{1}{d}\right).$$

Then using Lemma 36 and Lemma 37, we have that for third term

$$\mathbb{E}[\|\mathbf{s}\|^2 \beta_*^\top \mathbf{h}^\top \mathbf{h}\beta_*] = \frac{1}{\rho^2 d} \|\beta_*\|^2 + o\left(\frac{1}{\rho^2 d}\right) + O\left(\frac{1}{d}\right) \quad \text{and} \quad \text{Var}(\|\mathbf{s}\|^2 \beta_*^\top \mathbf{h}^\top \mathbf{h}\beta_*) = O\left(\frac{1}{d}\right)$$

for the fourth term

$$\begin{aligned} \mathbb{E}[\|\mathbf{h}\|^2 \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top] &= \left(\frac{1}{\rho^2} \frac{c}{c-1} + o\left(\frac{1}{\rho^2}\right)\right) \left(\frac{c-1}{c^2}(\beta_*^\top \mathbf{u})^2 + o(1)\right) + O\left(\frac{1}{\rho^2 d}\right) \\ &= \frac{(\beta_*^\top \mathbf{u})^2}{\rho^2 c} + o(1) + O\left(\frac{1}{\rho^2 d}\right) \end{aligned}$$

with variance

$$\text{Var}(\|\mathbf{h}\|^2 \beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top) = O\left(\frac{1}{\rho^2 d}\right).$$

For the first term, we have that

$$\mathbb{E}\left[\frac{\xi}{\eta}\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\beta_*\right] = 0 + O\left(\frac{1}{\rho^2 d}\right) \quad \text{and} \quad \text{Var}\left(\frac{\xi}{\eta}\beta_*^\top \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\beta_*\right) = O\left(\frac{1}{\rho^2 d}\right)$$

Adding the last three terms and using Lemma 34 twice, we get that

$$\mathbb{E}\left[\beta_*^\top \left(\|\mathbf{s}\|^2 \mathbf{h}^\top \mathbf{h} + \|\mathbf{h}\|^2 \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top + \frac{\xi}{\eta} \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\right) \beta_*\right] = \frac{1}{\rho^2 d} \|\beta_*\|^2 + \frac{(\beta_*^\top \mathbf{u})^2}{\rho^2 c} + 0 + o(1) + O\left(\frac{1}{d}\right)$$

With variance

$$\text{Var}\left(\beta_*^\top \left(\|\mathbf{s}\|^2 \mathbf{h}^\top \mathbf{h} + \|\mathbf{h}\|^2 \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top + \frac{\xi}{\eta} \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\right) \beta_*\right) = O\left(\frac{1}{d}\right)$$

Then recalling the mean and variance of η^2/γ_2 from Lemma 19, and using the covariance bound for the expectation of products, we get that

$$\mathbb{E}\left[\frac{\eta^2}{\gamma_2} \beta_*^\top \left(\|\mathbf{s}\|^2 \mathbf{h}^\top \mathbf{h} + \|\mathbf{h}\|^2 \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{s}^\top + \frac{\xi}{\eta} \mathbf{A}\mathbf{A}^\dagger \mathbf{u}\mathbf{h}\right) \beta_*\right] = \frac{\eta^2}{\eta^2 + \rho^2} \left(\frac{\|\beta_*\|^2}{d} + \frac{1}{c}(\beta_*^\top \mathbf{u})^2\right) + o(1) + O\left(\frac{1}{n}\right).$$

Adding all five terms, we get that

$$\mathbb{E}[\beta_*^\top (\mathbf{Z} + \mathbf{A})^\dagger \mathbf{A}^\top \beta_*] = \frac{1}{c} \|\beta_*\|^2 - \frac{\eta^2}{\eta^2 + \rho^2} \left(\frac{\|\beta_*\|^2}{d} + \frac{1}{c}(\beta_*^\top \mathbf{u})^2\right) + o(1) + O\left(\frac{1}{n}\right).$$

□

E.6 STEP 5: UPSCALING AND ASYMPTOTIC RISK FORMULAS

In the previous step we derived downscaled expressions for the four constituent terms of the risk: **Bias, Variance, Data Noise, and Target Alignment**. We stop our abuse of notation and are explicit again about downscaled vs. upscaled.

Bias (downscaled). For $c < 1$, the bias term is

$$\frac{\tilde{\eta}^2}{\tilde{n}} \left(\left[(\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 c + \rho^2} (\alpha_Z - \alpha_A) \right]^2 (\boldsymbol{\beta}_*^\top \mathbf{u})^2 + \tau_{\varepsilon,r}^2 \frac{c}{1-c} \frac{1}{\eta^2 c + \rho^2} \right) + o\left(\frac{1}{\tilde{n}}\right) + o\left(\frac{1}{n}\right).$$

For $c > 1$, the bias term is

$$\begin{aligned} \frac{\tilde{\eta}^2}{\tilde{n}} \left[(\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left((\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right) \right)^2 + \alpha_A^2 \frac{\|\boldsymbol{\beta}_*\|^2}{d} \left(\frac{c-1}{c} \right) \frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} \right. \\ \left. + \tau_{\varepsilon,r}^2 \frac{\eta^2 c + \rho^2}{c-1} \frac{1}{(\eta^2 + \rho^2)^2} \right] + o\left(\frac{1}{\tilde{n}}\right) + o\left(\frac{1}{n}\right) \end{aligned}$$

Variance (downscaled). For $c < 1$, the variance term is

$$\begin{aligned} \frac{\tilde{\rho}^2}{d} \left[\alpha_A^2 \|\boldsymbol{\beta}_*\|^2 + (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \left((\alpha_Z - \alpha_A)^2 \frac{\eta^2 (\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} + 2\alpha_A (\alpha_Z - \alpha_A) \frac{\eta^2 c}{\eta^2 c + \rho^2} \right) \right. \\ \left. + \tau_{\varepsilon,r}^2 \left(\frac{c}{1-c} \frac{d}{\rho^2} - \frac{\eta^2}{\rho^2 (\eta^2 c + \rho^2)} \frac{c^2}{1-c} \right) \right]. \end{aligned}$$

For $c > 1$, the variance term is

$$\frac{\tilde{\rho}^2}{d} \left[\|\boldsymbol{\beta}_*\|^2 \left(\frac{\alpha_A^2}{c} - \frac{\alpha_A^2}{d} \frac{\eta^2}{\eta^2 + \rho^2} \right) + (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{c}{c-1} \frac{\eta^2}{\eta^2 + \rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + \tau_{\varepsilon,r}^2 \left(\frac{d}{\rho^2} \frac{1}{c-1} - \frac{\eta^2}{\rho^2 (\eta^2 + \rho^2)} \frac{c}{c-1} \right) \right].$$

Data noise (downscaled). The data noise term is

$$\frac{\tilde{\alpha}_A^2 \tilde{\rho}^2}{d} \|\boldsymbol{\beta}_*\|^2.$$

Target alignment (downscaled). For $c < 1$, the alignment term is

$$-\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \left(\alpha_A \|\boldsymbol{\beta}_*\|^2 + (\alpha_Z - \alpha_A) (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \frac{\eta^2 c}{\rho^2 + \eta^2 c} \right).$$

For $c > 1$, the alignment term is

$$-\frac{2\tilde{\alpha}_A \tilde{\rho}^2}{d} \left(\frac{\alpha_A}{c} \|\boldsymbol{\beta}_*\|^2 - \frac{\alpha_A}{d} \frac{\eta^2}{\eta^2 + \rho^2} \|\boldsymbol{\beta}_*\|^2 + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\eta^2}{\eta^2 + \rho^2} (\boldsymbol{\beta}_*^\top \mathbf{u})^2 \right).$$

These formulas are expressed in terms of the concentrated building blocks, but still at the ‘‘microscopic’’ scale in which η is $O(\sqrt{d})$, $\rho = \Theta(1)$, and $\tau_{\varepsilon,r}^2 = O(1/d)$.

In this section we return to the macroscopic, or upscaled, version of the problem. Specifically, we multiply each term by d and reparametrize according to

$$\theta^2 = \frac{d}{n} \eta^2, \quad \tilde{\theta}^2 = \frac{d}{\tilde{n}} \tilde{\eta}^2, \quad \tau_\varepsilon^2 = d \tau_{\varepsilon,r}^2,$$

while keeping $\rho, \tilde{\rho}$ fixed. This normalization ensures that the effective spike strength θ , isotropic noise level ρ , and label noise $\tau_{\varepsilon,r}$ are all of order one. In this scaling, the risk is d times larger than in the downscaled representation, and the resulting formulas cleanly separate the contributions of the four terms.

The terms change as follows

Front factors (after multiplying by d).

$$\frac{\tilde{\eta}^2}{\tilde{n}} \xrightarrow{\times d} \tilde{\theta}^2, \quad \frac{\tilde{\rho}^2}{d} \xrightarrow{\times d} \tilde{\rho}^2, \quad \frac{\tilde{\alpha}_A^2 \tilde{\rho}^2}{d} \xrightarrow{\times d} \tilde{\alpha}_A^2 \tilde{\rho}^2, \quad d \tau_{\varepsilon,r}^2 \rightarrow \tau_{\varepsilon}^2. \quad (20)$$

Denominator identities.

$$\eta^2 c + \rho^2 = \theta^2 + \rho^2, \quad \eta^2 + \rho^2 = \frac{\theta^2 + c \rho^2}{c}. \quad (21)$$

Frequently used ratios and their upscaled forms.

$$\frac{\rho^2}{\eta^2 c + \rho^2} = \frac{\rho^2}{\theta^2 + \rho^2}, \quad (22)$$

$$\frac{\eta^2 c}{\eta^2 c + \rho^2} = \frac{\theta^2}{\theta^2 + \rho^2}, \quad (23)$$

$$\frac{\eta^2}{\eta^2 + \rho^2} = \frac{\theta^2}{\theta^2 + c \rho^2}, \quad (24)$$

$$\frac{\rho^2}{\eta^2 + \rho^2} = \frac{c \rho^2}{\theta^2 + c \rho^2}, \quad (25)$$

$$\frac{\eta^2 \rho^2}{(\eta^2 + \rho^2)^2} = \frac{\theta^2 \rho^2}{(\theta^2 + c \rho^2)^2} c, \quad (26)$$

$$\frac{\eta^2(\eta^2 + \rho^2)}{(\eta^2 c + \rho^2)^2} \frac{c^2}{1-c} = \frac{\theta^2(\theta^2 + c \rho^2)}{(\theta^2 + \rho^2)^2} \frac{1}{1-c}. \quad (27)$$

Noise terms with aspect-ratio factors. After multiplying by d and substituting $\tau_{\varepsilon}^2 = d \tau_{\varepsilon,r}^2$:

$$\tau_{\varepsilon,r}^2 \left(\frac{c}{1-c} \frac{d}{\rho^2} - \frac{\eta^2}{\rho^2(\eta^2 c + \rho^2)} \frac{c^2}{1-c} \right) \rightarrow \tau_{\varepsilon}^2 \left(\frac{1}{\rho^2} \frac{c}{1-c} - \frac{\theta^2}{\rho^2(\theta^2 + \rho^2)} \frac{c}{1-c} \right), \quad (28)$$

$$\tau_{\varepsilon,r}^2 \left(\frac{d}{\rho^2} \frac{1}{c-1} - \frac{\eta^2}{\rho^2(\eta^2 + \rho^2)} \frac{c}{c-1} \right) \rightarrow \tau_{\varepsilon}^2 \left(\frac{1}{\rho^2} \frac{1}{c-1} - \frac{\theta^2}{\rho^2(\theta^2 + c \rho^2)} \frac{c}{c-1} \right). \quad (29)$$

Alignment-specific identities.

$$\frac{\eta^2 c}{\rho^2 + \eta^2 c} = \frac{\theta^2}{\rho^2 + \theta^2}, \quad \frac{\eta^2}{\eta^2 + \rho^2} = \frac{\theta^2}{\theta^2 + c \rho^2}. \quad (30)$$

We now state the explicit upscaled limits for each component. As before, we present results separately in the underparametrized regime ($c < 1$) and the overparametrized regime ($c > 1$). Each term has a little $o(1)$ error term.

Bias. For $c < 1$, the bias contribution is

$$\tilde{\theta}^2 \left[\left((\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\theta^2 + \rho^2} (\alpha_Z - \alpha_A) \right)^2 (\beta_*^\top \mathbf{u})^2 + \frac{\tau_{\varepsilon}^2}{d} \frac{c}{1-c} \frac{1}{\theta^2 + \rho^2} \right].$$

For $c > 1$, the bias is

$$\tilde{\theta}^2 \left[(\beta_*^\top \mathbf{u})^2 \left((\tilde{\alpha}_Z - \alpha_Z) + \frac{\rho^2}{\frac{\theta^2}{c} + \rho^2} (\alpha_Z - \frac{\alpha_A}{c}) \right)^2 + \alpha_A^2 \frac{\|\beta_*\|^2}{d} \left(\frac{c-1}{c} \right) \frac{\frac{\theta^2}{c} \rho^2}{(\frac{\theta^2}{c} + \rho^2)^2} + \frac{\tau_{\varepsilon}^2}{d} \frac{1}{c-1} \frac{\theta^2 + \rho^2}{(\frac{\theta^2}{c} + \rho^2)^2} \right].$$

Variance. For $c < 1$, the variance contribution is

$$\tilde{\rho}^2 \left[\alpha_A^2 \|\beta_*\|^2 + (\beta_*^\top \mathbf{u})^2 \left((\alpha_Z - \alpha_A)^2 \frac{\theta^2(\theta^2 + c\rho^2)}{(\theta^2 + \rho^2)^2} \frac{1}{1-c} + 2\alpha_A(\alpha_Z - \alpha_A) \frac{\theta^2}{\theta^2 + \rho^2} \right) + \tau_\varepsilon^2 \left(\frac{1}{\rho^2} \frac{c}{1-c} - \frac{1}{d} \frac{\theta^2}{\rho^2(\theta^2 + \rho^2)} \cdot \frac{c}{1-c} \right) \right].$$

For $c > 1$, the variance is

$$\tilde{\rho}^2 \left[\|\beta_*\|^2 \left(\frac{\alpha_A^2}{c} - \frac{\alpha_A^2}{d} \frac{\theta^2}{\theta^2 + c\rho^2} \right) + (\beta_*^\top \mathbf{u})^2 \frac{c}{c-1} \frac{\theta^2}{\theta^2 + c\rho^2} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + \tau_\varepsilon^2 \left(\frac{1}{\rho^2} \frac{1}{c-1} - \frac{1}{d} \frac{\theta^2}{\rho^2(\theta^2 + c\rho^2)} \cdot \frac{c}{c-1} \right) \right].$$

Data Noise. The data noise term is independent of c :

$$\tilde{\alpha}_A^2 \tilde{\rho}^2 \|\beta_*\|^2.$$

Target Alignment. For $c < 1$, the target alignment contribution is

$$-2\tilde{\alpha}_A \tilde{\rho}^2 \left(\alpha_A \|\beta_*\|^2 + (\alpha_Z - \alpha_A) (\beta_*^\top \mathbf{u})^2 \frac{\theta^2}{\rho^2 + \theta^2} \right).$$

For $c > 1$, the alignment term is

$$-2\tilde{\alpha}_A \tilde{\rho}^2 \left(\frac{\alpha_A}{c} \|\beta_*\|^2 - \frac{\alpha_A}{d} \frac{\theta^2}{\theta^2 + c\rho^2} \|\beta_*\|^2 + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\theta^2}{\theta^2 + c\rho^2} (\beta_*^\top \mathbf{u})^2 \right).$$

Lastly, replacing $\tilde{\rho}, \rho$ with $\tilde{\tau}, \tau$ and using $d/n \rightarrow c$ yield the detailed expressions in Theorem 5, up to simple algebra (rearranging terms and simplifying the fractions).

F PROBABILITY LEMMAS

Proposition 2. If $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ are fixed unit norm vector and $\mathbf{A} \in \mathbb{R}^{d \times n}$ is a Gaussian matrix with i.i.d. $\mathcal{N}(0, 1)$ entries. If $d > n$, then we have that

$$\mathbb{E}[(\mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{v})^2] = \frac{n}{d(d+2)} \left[(\mathbf{u}^\top \mathbf{v})^2 (n+2) + \frac{(1 - (\mathbf{u}^\top \mathbf{v})^2)(d-n)}{d-1} \right] = \frac{1}{c^2} (\mathbf{u}^\top \mathbf{v})^2 + o(1),$$

$$\text{Var}(\mathbf{u}^\top \mathbf{A} \mathbf{A}^\dagger \mathbf{v})^2 = O\left(\frac{1}{d}\right).$$

Proof. Let $\mathbf{P} := \mathbf{A} \mathbf{A}^\dagger$. This is the orthogonal projection matrix onto the column space of \mathbf{A} , denoted $C(\mathbf{A}) = \text{Range}(\mathbf{A})$. The subspace $C(\mathbf{A})$ is an n -dimensional subspace of \mathbb{R}^d . Because the entries A_{ij} are i.i.d. $\mathcal{N}(0, 1)$, the distribution of the random subspace $C(\mathbf{A})$ is isotropic (or rotationally invariant). Consequently, the distribution of the random projection matrix \mathbf{P} is also rotationally invariant. That is, for any fixed $d \times d$ orthogonal matrix \mathbf{Q} , the distribution of $\mathbf{Q} \mathbf{P} \mathbf{Q}^\top$ is the same as the distribution of \mathbf{P} .

We are interested in $\mathbb{E}[(\mathbf{u}^\top \mathbf{P} \mathbf{v})^2]$. Let θ be the angle between \mathbf{u} and \mathbf{v} , such that $\cos(\theta) = \mathbf{u}^\top \mathbf{v}$ (since they are unit vectors). Due to the rotational invariance of the distribution of \mathbf{P} , we can choose an orthonormal basis without loss of generality. Let \mathbf{Q} be an orthogonal matrix such that $\mathbf{u}' = \mathbf{Q} \mathbf{u} = \mathbf{e}_1 = (1, 0, \dots, 0)^\top$ and $\mathbf{v}' = \mathbf{Q} \mathbf{v}$ lies in the span of \mathbf{e}_1 and \mathbf{e}_2 . Specifically, $\mathbf{v}' = \cos(\theta) \mathbf{e}_1 + \sin(\theta) \mathbf{e}_2$. Let $\mathbf{P}' = \mathbf{Q} \mathbf{P} \mathbf{Q}^\top$. \mathbf{P}' has the same distribution as \mathbf{P} . Then,

$$\mathbf{u}^\top \mathbf{P} \mathbf{v} = (\mathbf{Q}^\top \mathbf{u}')^\top \mathbf{P} (\mathbf{Q}^\top \mathbf{v}') = (\mathbf{u}')^\top (\mathbf{Q} \mathbf{P} \mathbf{Q}^\top) \mathbf{v}' = (\mathbf{u}')^\top \mathbf{P}' \mathbf{v}'$$

Substituting $\mathbf{u}' = \mathbf{e}_1$ and $\mathbf{v}' = \cos(\theta) \mathbf{e}_1 + \sin(\theta) \mathbf{e}_2$:

$$\begin{aligned} \mathbf{u}^\top \mathbf{P} \mathbf{v} &= \mathbf{e}_1^\top \mathbf{P}' (\cos(\theta) \mathbf{e}_1 + \sin(\theta) \mathbf{e}_2) \\ &= \cos(\theta) (\mathbf{e}_1^\top \mathbf{P}' \mathbf{e}_1) + \sin(\theta) (\mathbf{e}_1^\top \mathbf{P}' \mathbf{e}_2) \\ &= \cos(\theta) P'_{11} + \sin(\theta) P'_{12} \end{aligned}$$

where P'_{ij} are the elements of \mathbf{P}' . Since \mathbf{P}' has the same distribution as \mathbf{P} , we can drop the prime for calculating expectations involving the elements. Let $\mathbf{X} = \mathbf{u}^\top \mathbf{P} \mathbf{v}$. We then need $\mathbb{E}[X^2]$.

$$\begin{aligned}\mathbb{E}[X^2] &= \mathbb{E}[(\cos(\theta)P_{11} + \sin(\theta)P_{12})^2] \\ &= \mathbb{E}[\cos^2(\theta)P_{11}^2 + \sin^2(\theta)P_{12}^2 + 2\cos(\theta)\sin(\theta)P_{11}P_{12}] \\ &= \cos^2(\theta)\mathbb{E}[P_{11}^2] + \sin^2(\theta)\mathbb{E}[P_{12}^2] + 2\cos(\theta)\sin(\theta)\mathbb{E}[P_{11}P_{12}]\end{aligned}$$

Calculation of Moments. We need to compute $\mathbb{E}[P_{11}^2]$, $\mathbb{E}[P_{12}^2]$, and $\mathbb{E}[P_{11}P_{12}]$.

Consider a reflection matrix \mathbf{R} that maps \mathbf{e}_2 to $-\mathbf{e}_2$ and leaves other basis vectors unchanged (i.e., $\mathbf{R} = \text{diag}(1, -1, 1, \dots, 1)$). Since the distribution of \mathbf{P} is isotropic, it is invariant under reflection. Let $\mathbf{P}^* = \mathbf{R}\mathbf{P}\mathbf{R}^\top = \mathbf{R}\mathbf{P}\mathbf{R}$. \mathbf{P}^* has the same distribution as \mathbf{P} . The components are related:

$$P_{11}^* = (\mathbf{R}\mathbf{P}\mathbf{R})_{11} = R_{11}P_{11}R_{11} = P_{11}$$

and

$$P_{12}^* = (\mathbf{R}\mathbf{P}\mathbf{R})_{12} = R_{11}P_{12}R_{22} = (1)P_{12}(-1) = -P_{12}.$$

Therefore,

$$\mathbb{E}[P_{11}P_{12}] = \mathbb{E}[P_{11}^*P_{12}^*] = \mathbb{E}[P_{11}(-P_{12})] = -\mathbb{E}[P_{11}P_{12}].$$

This implies $2\mathbb{E}[P_{11}P_{12}] = 0$, so $\mathbb{E}[P_{11}P_{12}] = 0$.

The diagonal element $P_{11} = \mathbf{e}_1^\top \mathbf{P} \mathbf{e}_1 = \|\mathbf{P} \mathbf{e}_1\|_2^2$ represents the squared norm of the projection of the fixed unit vector \mathbf{e}_1 onto the random n -dimensional subspace $C(\mathbf{A})$. This variable follows a Beta distribution:

$$P_{11} \sim \text{Beta}\left(\frac{n}{2}, \frac{d-n}{2}\right)$$

The mean and variance of a $\text{Beta}(\alpha, \beta)$ distribution are $\frac{\alpha}{\alpha+\beta}$ and $\frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$, respectively. Here, $\alpha = n/2$ and $\beta = (d-n)/2$, so $\alpha + \beta = d/2$.

$$\mathbb{E}[P_{11}] = \frac{n/2}{d/2} = \frac{n}{d}$$

Next

$$\text{Var}(P_{11}) = \frac{(n/2)((d-n)/2)}{(d/2)^2(d/2+1)} = \frac{n(d-n)/4}{(d^2/4)((d+2)/2)} = \frac{n(d-n) \cdot 8}{4d^2(d+2)} = \frac{2n(d-n)}{d^2(d+2)}$$

Now we find $\mathbb{E}[P_{11}^2]$ using $\mathbb{E}[P_{11}^2] = \text{Var}(P_{11}) + (\mathbb{E}[P_{11}])^2$:

$$\begin{aligned}\mathbb{E}[P_{11}^2] &= \frac{2n(d-n)}{d^2(d+2)} + \left(\frac{n}{d}\right)^2 \\ &= \frac{2n(d-n) + n^2(d+2)}{d^2(d+2)} \\ &= \frac{2nd - 2n^2 + n^2d + 2n^2}{d^2(d+2)} \\ &= \frac{2nd + n^2d}{d^2(d+2)} \\ &= \frac{n(n+2)}{d(d+2)}.\end{aligned}$$

We use the property that \mathbf{P} is a projection matrix, so $\mathbf{P}^2 = \mathbf{P}$. The trace is $\text{Tr}(\mathbf{P}) = n$. Also $\text{Tr}(\mathbf{P}^2) = \text{Tr}(\mathbf{P}) = n$. We can write $\text{Tr}(\mathbf{P}^2) = \text{Tr}(\mathbf{P}\mathbf{P}^\top)$ since \mathbf{P} is symmetric.

$$\text{Tr}(\mathbf{P}^2) = \sum_{i=1}^d \sum_{j=1}^d (P_{ij})^2$$

Taking the expectation:

$$\mathbb{E}[\text{Tr}(\mathbf{P}^2)] = \mathbb{E} \left[\sum_{i,j} P_{ij}^2 \right] = \sum_{i,j} \mathbb{E}[P_{ij}^2] = n$$

By rotational symmetry, $\mathbb{E}[P_{ii}^2]$ is the same for all i , and $\mathbb{E}[P_{ij}^2]$ is the same for all $i \neq j$.

$$\sum_{i=1}^d \mathbb{E}[P_{ii}^2] + \sum_{i \neq j} \mathbb{E}[P_{ij}^2] = n.$$

There are d diagonal terms and $d(d-1)$ off-diagonal terms.

$$d \mathbb{E}[P_{11}^2] + d(d-1) \mathbb{E}[P_{12}^2] = n$$

Substitute the value for $\mathbb{E}[P_{11}^2]$ (assuming $d > 1$):

$$\begin{aligned} d \left(\frac{n(n+2)}{d(d+2)} \right) + d(d-1) \mathbb{E}[P_{12}^2] &= n \\ \frac{n(n+2)}{d+2} + d(d-1) \mathbb{E}[P_{12}^2] &= n \\ d(d-1) \mathbb{E}[P_{12}^2] &= n - \frac{n(n+2)}{d+2} = \frac{n(d+2) - n(n+2)}{d+2} = \frac{nd + 2n - n^2 - 2n}{d+2} = \frac{n(d-n)}{d+2} \\ \mathbb{E}[P_{12}^2] &= \frac{n(d-n)}{d(d-1)(d+2)} \end{aligned}$$

Substitute the moments back into the expression for $\mathbb{E}[X^2]$:

$$\mathbb{E}[X^2] = \cos^2(\theta) \mathbb{E}[P_{11}^2] + \sin^2(\theta) \mathbb{E}[P_{12}^2] + 2 \cos(\theta) \sin(\theta) \cdot 0$$

Using $\cos(\theta) = \mathbf{u}^\top \mathbf{v}$, $\cos^2(\theta) = (\mathbf{u}^\top \mathbf{v})^2$, and $\sin^2(\theta) = 1 - \cos^2(\theta) = 1 - (\mathbf{u}^\top \mathbf{v})^2$:

$$\begin{aligned} \mathbb{E}[(\mathbf{u}^\top \mathbf{A} \mathbf{A}^\top \mathbf{v})^2] &= (\mathbf{u}^\top \mathbf{v})^2 \left(\frac{n(n+2)}{d(d+2)} \right) + (1 - (\mathbf{u}^\top \mathbf{v})^2) \left(\frac{n(d-n)}{d(d-1)(d+2)} \right) \\ &= \frac{n}{d(d+2)} \left[(\mathbf{u}^\top \mathbf{v})^2 (n+2) + \frac{(1 - (\mathbf{u}^\top \mathbf{v})^2)(d-n)}{d-1} \right] \\ &= \frac{1}{c^2} (\mathbf{u}^\top \mathbf{v})^2 + O\left(\frac{1}{d}\right). \end{aligned}$$

Calculation of Variance. Recall that reflection $\mathbf{R} = \text{diag}(1, -1, 1, \dots, 1)$ implies $\mathbf{P} \stackrel{d}{=} \mathbf{R} \mathbf{P} \mathbf{R}$ (equal in distribution) and thus $\mathbb{E}[P_{11} P_{12}] = 0$, and in general any mixed moment with an odd power of P_{12} vanishes. Therefore, we have the following expansion:

$$\mathbb{E}[X^4] = \cos^4 \theta \mathbb{E}[P_{11}^4] + 6 \cos^2 \theta \sin^2 \theta \mathbb{E}[P_{11}^2 P_{12}^2] + \sin^4 \theta \mathbb{E}[P_{12}^4]. \quad (31)$$

We start with $\mathbb{E}[P_{11}^4]$. Since $P_{11} \sim \text{Beta}(\alpha, \beta)$ with $\alpha = \frac{n}{2}$, $\beta = \frac{d-n}{2}$. We need the higher moments for the Beta distribution: for $m \geq 1$,

$$\mathbb{E}[P_{11}^m] = \frac{\alpha^{(m)}}{(\alpha + \beta)^{(m)}} = \frac{\left(\frac{n}{2}\right)^{(m)}}{\left(\frac{d}{2}\right)^{(m)}, \quad x^{(m)} := x(x+1) \cdots (x+m-1).$$

In particular, we have the following third and fourth moments:

$$\mathbb{E}[P_{11}^3] = \frac{\left(\frac{n}{2}\right)^{(3)}}{\left(\frac{d}{2}\right)^{(3)}} = \frac{1}{c^3} + O\left(\frac{1}{d}\right), \quad \mathbb{E}[P_{11}^4] = \frac{\left(\frac{n}{2}\right)^{(4)}}{\left(\frac{d}{2}\right)^{(4)}} = \frac{1}{c^4} + O\left(\frac{1}{d}\right).$$

We now move on to $\mathbb{E}[P_{11}^2 P_{12}^2]$. From idempotency, $(P^2)_{11} = P_{11}$ gives the row identity $P_{11} = \sum_{k=1}^d P_{1k}^2$. Multiplying by P_{11}^2 and taking expectations, we have that

$$\mathbb{E}[P_{11}^3] = \mathbb{E}[P_{11}^4] + \sum_{k=2}^d \mathbb{E}[P_{11}^2 P_{1k}^2] = \mathbb{E}[P_{11}^4] + (d-1) \mathbb{E}[P_{11}^2 P_{12}^2].$$

$$\mathbb{E}[P_{11}^2 P_{12}^2] = \frac{\mathbb{E}[P_{11}^3] - \mathbb{E}[P_{11}^4]}{d-1} = \frac{1}{d-1} \left(\frac{\binom{n}{2}^{(3)}}{\binom{d}{2}^{(3)}} - \frac{\binom{n}{2}^{(4)}}{\binom{d}{2}^{(4)}} \right) = \frac{1}{d-1} \left(\frac{1}{c^3} - \frac{1}{c^4} + O\left(\frac{1}{d}\right) \right) = O\left(\frac{1}{d}\right).$$

We still need to evaluate or upper bound $\mathbb{E}[P_{12}^4]$. From $P_{11} = \sum_{k=1}^d P_{1k}^2$ we have $\sum_{k=2}^d P_{1k}^2 = P_{11} - P_{11}^2$. By Cauchy–Schwarz,

$$\sum_{k=2}^d P_{1k}^4 = \left(\sum_{k=2}^d P_{1k}^2 \right)^2 = (P_{11} - P_{11}^2)^2.$$

Taking expectations, we get:

$$(d-1)\mathbb{E}[P_{12}^4] \leq \mathbb{E}[(P_{11} - P_{11}^2)^2] = \mathbb{E}[P_{11}^2] - 2\mathbb{E}[P_{11}^3] + \mathbb{E}[P_{11}^4].$$

$$\mathbb{E}[P_{12}^4] \leq \frac{1}{d-1} \left(\frac{1}{c^2} - \frac{2}{c^3} + \frac{1}{c^4} \right) + O\left(\frac{1}{d^2}\right) = O\left(\frac{1}{d}\right).$$

We can now plug these expectation bounds into Equation 31:

$$\begin{aligned} \mathbb{E}[X^4] &= \cos^4 \theta \frac{\binom{n}{2}^{(4)}}{\binom{d}{2}^{(4)}} + O\left(\frac{1}{d}\right) 6 \cos^2 \theta \sin^2 \theta + O\left(\frac{1}{d}\right) \sin^4 \theta \\ &= \frac{1}{c^4} (\mathbf{u}^\top \mathbf{v})^4 + O\left(\frac{1}{d}\right). \end{aligned}$$

Recall from the prior proof that:

$$\mathbb{E}[X^2] = \cos^2 \theta \frac{n(n+2)}{d(d+2)} + \sin^2 \theta \frac{n(d-n)}{d(d-1)(d+2)} = \frac{1}{c^2} (\mathbf{u}^\top \mathbf{v})^2 + O\left(\frac{1}{d}\right).$$

Finally, we have that the variance is of order:

$$\text{Var}(X^2) = \mathbb{E}[X^4] - (\mathbb{E}[X^2])^2 = O\left(\frac{1}{d}\right).$$

□

Lemma 32. *Let $a \neq 0$ be a constant and suppose that $\zeta = a + o(f(n))$ as $n \rightarrow \infty$. Then,*

$$\frac{1}{\zeta} = \frac{1}{a} + o(f(n)).$$

Proof. Write $\zeta = a + r_n$ with $r_n = o(f(n))$. Then

$$\frac{1}{\zeta} = \frac{1}{a + r_n} = \frac{1}{a} \cdot \frac{1}{1 + \frac{r_n}{a}}.$$

Using the expansion

$$\frac{1}{1+u} = 1 - u + O(u^2) \quad \text{as } u \rightarrow 0,$$

with $u = r_n/a$, we obtain

$$\frac{1}{\zeta} = \frac{1}{a} \left(1 - \frac{r_n}{a} + O((r_n/a)^2) \right) = \frac{1}{a} - \frac{r_n}{a^2} + O(r_n^2).$$

Since $r_n = o(f(n))$ and $f(n) \rightarrow 0$, we have $r_n^2 = o(f(n))$. Therefore

$$\frac{1}{\zeta} = \frac{1}{a} + o(f(n)),$$

which is the desired expansion. □

Lemma 33 (Variance of a reciprocal). *Let X be a random variable satisfying*

$$\mathbb{E}[X] = a > 0 \quad \text{and} \quad \text{Var}(X) = \sigma^2 = o(1),$$

and assume that X is bounded away from zero with high probability. That is, there exists $C \in (0, a)$ such that

$$\Pr[X \geq C] = 1 - o(1)$$

If there exists an M such that

$$\mathbb{E}[X^{-8}] \leq M \quad \text{and} \quad \mathbb{E}[(X - \mathbb{E}[X])^4] = O(\sigma^4)$$

Then

$$\text{Var}\left(\frac{1}{X}\right) = \frac{1}{a^4} \text{Var}(X) + o(\text{Var}(X)),$$

so in particular, $\text{Var}(1/X) = o(1)$.

Proof. Let $Y := X - a$. Then

$$\mathbb{E}[Y] = 0, \quad \mathbb{E}[Y^2] = \sigma^2, \quad \mathbb{E}[Y^4] = O(\sigma^4).$$

By Taylor's theorem with Lagrange remainder for $f(x) = 1/x$, there exists $\theta = \theta(X) \in (0, 1)$ such that

$$\frac{1}{X} = \frac{1}{a} - \frac{Y}{a^2} + Z, \quad Z := \frac{Y^2}{(a + \theta Y)^3} \geq 0.$$

Write $\Delta := \frac{1}{X} - \frac{1}{a} = -\frac{Y}{a^2} + Z$. Then

$$\text{Var}\left(\frac{1}{X}\right) = \mathbb{E}[\Delta^2] - (\mathbb{E}[\Delta])^2.$$

We will show

$$\mathbb{E}[\Delta^2] = \frac{\sigma^2}{a^4} + o(\sigma^2) \quad \text{and} \quad (\mathbb{E}[\Delta])^2 = o(\sigma^2).$$

Let $G := \{X \geq C\}$ and $B := \{X < C\}$. Since $C < a$ and $\mathbb{E}[Y^2] = \sigma^2$, Chebyshev gives the quantitative bound

$$\Pr[B] = \Pr[|Y| \geq a - C] \leq \frac{\mathbb{E}[Y^2]}{(a - C)^2} = \frac{\sigma^2}{(a - C)^2} = O(\sigma^2) = o(1).$$

Second moment $\mathbb{E}[\Delta^2]$. We split over G and B .

On G . Since $a + \theta Y = \theta X + (1 - \theta)a \geq C$, we have

$$|Z| \leq \frac{Y^2}{C^3}, \quad Z^2 \leq \frac{Y^4}{C^6}.$$

Therefore

$$\mathbb{E}\left[\left(-\frac{Y}{a^2} + Z\right)^2 \mathbf{1}_G\right] = \frac{1}{a^4} \mathbb{E}[Y^2 \mathbf{1}_G] - \frac{2}{a^2} \mathbb{E}[YZ \mathbf{1}_G] + \mathbb{E}[Z^2 \mathbf{1}_G].$$

We bound each term as follows.

$$\mathbb{E}[Z^2 \mathbf{1}_G] \leq \frac{1}{C^6} \mathbb{E}[Y^4] = O(\sigma^4),$$

and, using $\mathbf{1}_G \leq 1$ and Lyapunov/monotonicity of L^p norms,

$$\mathbb{E}[|YZ| \mathbf{1}_G] \leq \frac{1}{C^3} \mathbb{E}[|Y|^3] \leq \frac{1}{C^3} (\mathbb{E}[Y^4])^{3/4} = O(\sigma^3) = o(\sigma^2).$$

Moreover,

$$\mathbb{E}[Y^2 \mathbf{1}_G] = \sigma^2 - \mathbb{E}[Y^2 \mathbf{1}_B], \quad \mathbb{E}[Y^2 \mathbf{1}_B] \leq (\mathbb{E}[Y^4])^{1/2} \Pr[B]^{1/2} = O(\sigma^2) \Pr[B]^{1/2} = o(\sigma^2).$$

Hence

$$\mathbb{E} \left[\left(-\frac{Y}{a^2} + Z \right)^2 \mathbf{1}_G \right] = \frac{\sigma^2}{a^4} + o(\sigma^2).$$

On B . Using the algebraic identity

$$\left(\frac{1}{X} - \frac{1}{a} \right)^2 = \frac{Y^2}{a^2 X^2},$$

Cauchy–Schwarz and Hölder (with exponents 2, 2) give

$$\mathbb{E} [\Delta^2 \mathbf{1}_B] = \frac{1}{a^2} \mathbb{E} \left[\frac{Y^2}{X^2} \mathbf{1}_B \right] \leq \frac{1}{a^2} (\mathbb{E} [Y^4])^{1/2} (\mathbb{E} [X^{-4} \mathbf{1}_B])^{1/2} \leq \frac{1}{a^2} O(\sigma^2) (\mathbb{E} [X^{-8}])^{1/4} \Pr [B]^{1/4}.$$

Under the lemma’s assumption $\mathbb{E} [X^{-8}] \leq M$, we get

$$\mathbb{E} [\Delta^2 \mathbf{1}_B] = O(\sigma^2) \Pr [B]^{1/4} = o(\sigma^2).$$

Combining the G and B parts,

$$\mathbb{E} [\Delta^2] = \frac{\sigma^2}{a^4} + o(\sigma^2).$$

Mean correction $(\mathbb{E} [\Delta])^2$. Since $\mathbb{E} [Y] = 0$, we have

$$\mathbb{E} [\Delta] = \mathbb{E} [Z] = \mathbb{E} [Z \mathbf{1}_G] + \mathbb{E} [Z \mathbf{1}_B].$$

On G , $Z \leq Y^2/C^3$, so

$$\mathbb{E} [Z \mathbf{1}_G] \leq \frac{1}{C^3} \mathbb{E} [Y^2 \mathbf{1}_G] \leq \frac{1}{C^3} \sigma^2.$$

On B , The inequality

$$Z = \frac{Y^2}{a + \theta Y} \leq \frac{X^2}{Y^3}$$

holds on set B because on this set as $X < a$, meaning the point $a + \theta Y$ lies between X and a , so $a + \theta Y > X$. Thus, using Cauchy–Schwarz and Hölder,

$$\mathbb{E} [Z \mathbf{1}_B] \leq \mathbb{E} \left[\frac{Y^2}{X^3} \mathbf{1}_B \right] \leq (\mathbb{E} [Y^4])^{1/2} (\mathbb{E} [X^{-6} \mathbf{1}_B])^{1/2} \leq O(\sigma^2) (\mathbb{E} [X^{-12}])^{1/4} \Pr [B]^{1/4} = o(\sigma^2).$$

Thus $|\mathbb{E} [\Delta]| = O(\sigma^2)$ and therefore

$$(\mathbb{E} [\Delta])^2 = O(\sigma^4) = o(\sigma^2).$$

Putting the two steps together,

$$\text{Var} \left(\frac{1}{X} \right) = \mathbb{E} [\Delta^2] - (\mathbb{E} [\Delta])^2 = \frac{\sigma^2}{a^4} + o(\sigma^2) = \frac{1}{a^4} \text{Var}(X) + o(\text{Var}(X)).$$

□

Lemma 34 (Variance of a sum). *Let A and B be any random variables with finite variances $V(A) = \text{Var}(A)$ and $V(B) = \text{Var}(B)$. Then,*

$$\text{Var}(A + B) \leq \left(\sqrt{V(A)} + \sqrt{V(B)} \right)^2.$$

Proof. Recall that

$$\text{Var}(A + B) = \text{Var}(A) + \text{Var}(B) + 2 \text{Cov}(A, B).$$

By the Cauchy–Schwarz inequality, we have

$$|\text{Cov}(A, B)| \leq \sqrt{V(A)V(B)}.$$

Thus,

$$\text{Var}(A + B) \leq V(A) + V(B) + 2\sqrt{V(A)V(B)} = \left(\sqrt{V(A)} + \sqrt{V(B)} \right)^2.$$

□

Lemma 35 (Variance of one product). *Let A, B be real random variables with means $a = \mathbb{E}[A]$, $b = \mathbb{E}[B]$ and finite variances. Assume*

$$\mathbb{E}[(A - a)^4] \leq K_A \text{Var}(A)^2, \quad \mathbb{E}[(B - b)^4] \leq K_B \text{Var}(B)^2.$$

Then, with $C_4 := (K_A K_B)^{1/4}$,

$$\sqrt{\text{Var}(AB)} \leq |a| \sqrt{\text{Var}(B)} + |b| \sqrt{\text{Var}(A)} + C_4 \sqrt{\text{Var}(A) \text{Var}(B)}.$$

Moreover, as $\text{Var}(A), \text{Var}(B) \rightarrow 0$,

$$\text{Var}(AB) = O(a^2 \text{Var}(B)) + O(b^2 \text{Var}(A)) + o(\text{Var}(A) + \text{Var}(B)).$$

It directly follows that if all the means are $O(1)$,

$$\begin{aligned} \text{Var}(AB) &= O(\text{Var}(B)) + O(\text{Var}(A)). \\ \text{Var}(ABC) &= O(\text{Var}(C)) + O(\text{Var}(B)) + O(\text{Var}(A)) \quad \text{and so on by induction.} \end{aligned}$$

Proof. Write

$$AB - ab = a\tilde{B} + b\tilde{A} + \tilde{A}\tilde{B}.$$

Using $\text{Var}(U + V) = \text{Var}(U) + \text{Var}(V) + 2\text{Cov}(U, V)$ and $|\text{Cov}(U, V)| \leq \sqrt{\text{Var}(U)\text{Var}(V)}$, we get

$$\begin{aligned} \text{Var}(AB) &= \text{Var}(a\tilde{B} + b\tilde{A} + \tilde{A}\tilde{B}) \\ &\leq \left(|a| \sqrt{\text{Var}(\tilde{B})} + |b| \sqrt{\text{Var}(\tilde{A})} + \sqrt{\text{Var}(\tilde{A}\tilde{B})} \right)^2. \end{aligned}$$

Since $\text{Var}(\tilde{A}) = \text{Var}(A)$ and $\text{Var}(\tilde{B}) = \text{Var}(B)$, it remains to bound $\text{Var}(\tilde{A}\tilde{B})$. By Cauchy–Schwarz (Hölder with $p = q = 2$),

$$\text{Var}(\tilde{A}\tilde{B}) \leq \mathbb{E}[\tilde{A}^2 \tilde{B}^2] \leq \left(\mathbb{E}[\tilde{A}^4] \right)^{1/2} \left(\mathbb{E}[\tilde{B}^4] \right)^{1/2}.$$

Since we assume fourth–moment control $\mathbb{E}[\tilde{A}^4] \leq K_A \text{Var}(A)^2$ and $\mathbb{E}[\tilde{B}^4] \leq K_B \text{Var}(B)^2$, then

$$\sqrt{\text{Var}(\tilde{A}\tilde{B})} \leq (K_A K_B)^{1/4} \sqrt{\text{Var}(A) \text{Var}(B)}.$$

Hence

$$\text{Var}(AB) \leq \left(|a| \sqrt{\text{Var}(B)} + |b| \sqrt{\text{Var}(A)} + C_4 \sqrt{\text{Var}(A) \text{Var}(B)} \right)^2, \quad C_4 := (K_A K_B)^{1/4}.$$

For the moreover part, using the exact variance–covariance expansion,

$$\text{Var}(AB) = a^2 \text{Var}(B) + b^2 \text{Var}(A) + 2ab \text{Cov}(A, B) + \text{Var}(\tilde{A}\tilde{B}) + 2a \text{Cov}(\tilde{B}, \tilde{A}\tilde{B}) + 2b \text{Cov}(\tilde{A}, \tilde{A}\tilde{B}),$$

we bound the three remainder terms using Cauchy–Schwarz and the fourth–moment control:

$$\begin{aligned} \text{Var}(\tilde{A}\tilde{B}) &\leq \mathbb{E}[\tilde{A}^2 \tilde{B}^2] \leq \left(\mathbb{E}[\tilde{A}^4] \right)^{1/2} \left(\mathbb{E}[\tilde{B}^4] \right)^{1/2} \leq C_4^2 \text{Var}(A) \text{Var}(B), \\ |\text{Cov}(\tilde{B}, \tilde{A}\tilde{B})| &\leq \sqrt{\text{Var}(\tilde{B})} \sqrt{\text{Var}(\tilde{A}\tilde{B})} \leq C_4 \text{Var}(B) \sqrt{\text{Var}(A)}, \\ |\text{Cov}(\tilde{A}, \tilde{A}\tilde{B})| &\leq \sqrt{\text{Var}(\tilde{A})} \sqrt{\text{Var}(\tilde{A}\tilde{B})} \leq C_4 \text{Var}(A) \sqrt{\text{Var}(B)}. \end{aligned}$$

As $\text{Var}(A), \text{Var}(B) \rightarrow 0$, each of these is $o(\text{Var}(A) + \text{Var}(B))$.

For the covariance term, Cauchy–Schwarz and the inequality $2uv \leq \varepsilon u^2 + \varepsilon^{-1}v^2$ (for any $\varepsilon > 0$) with $u := |a|\sqrt{\text{Var}(B)}$, $v := |b|\sqrt{\text{Var}(A)}$ give

$$|2ab \text{Cov}(A, B)| \leq 2|ab|\sqrt{\text{Var}(A)\text{Var}(B)} \leq \varepsilon a^2 \text{Var}(B) + \varepsilon^{-1}b^2 \text{Var}(A).$$

Therefore,

$$\text{Var}(AB) \leq (1 + \varepsilon) a^2 \text{Var}(B) + (1 + \varepsilon^{-1}) b^2 \text{Var}(A) + o(\text{Var}(A) + \text{Var}(B)).$$

Choosing, e.g., $\varepsilon = 1$ yields

$$\text{Var}(AB) = O(a^2 \text{Var}(B)) + O(b^2 \text{Var}(A)) + o(\text{Var}(A) + \text{Var}(B)),$$

which proves the moreover statement. \square

Lemma 36 (Variance of general product). *Let $m \geq 2$ and let X_1, \dots, X_m be real random variables with nonzero means $\mu_i := \mathbb{E}[X_i] \neq 0$ and variances $f_i(n) := \text{Var}(X_i) \rightarrow 0$ as $n \rightarrow \infty$. Assume that for some integer $M \geq m$ (it is enough to take $M = m$),*

$$\mathbb{E}[|X_i - \mu_i|^{2M}] = O(\text{Var}(X_i)^M) \quad \text{for each } i = 1, \dots, m. \quad (32)$$

Then

$$\text{Var}\left(\prod_{i=1}^m X_i\right) = O\left(\left(\sum_{i=1}^m \sqrt{f_i(n)}\right)^2\right) = O\left(\max_{1 \leq i \leq m} f_i(n)\right).$$

Proof. Write $\Delta_i := X_i - \mu_i$ so that $\mathbb{E}[\Delta_i] = 0$ and $\|\Delta_i\|_{L_2} = \sigma_i$. By assumption Equation 32 with $M \geq m$ and monotonicity of L_p norms,

$$\|\Delta_i\|_{L_{2k}} = O\left(\sqrt{f_i(n)}\right) \quad \text{for every } 1 \leq k \leq m, i = 1, \dots, m.$$

Expand the product multilinearly:

$$\prod_{i=1}^m X_i - \prod_{i=1}^m \mu_i = \sum_{\emptyset \neq S \subseteq [m]} \left(\prod_{j \in S^c} \mu_j\right) \left(\prod_{i \in S} \Delta_i\right).$$

Taking L_2 norms and using the triangle inequality,

$$\left\|\prod_{i=1}^m X_i - \prod_{i=1}^m \mu_i\right\|_{L_2} \leq \sum_{\emptyset \neq S \subseteq [m]} \left(\prod_{j \in S^c} |\mu_j|\right) \left\|\prod_{i \in S} \Delta_i\right\|_{L_2}.$$

For a fixed nonempty S with $|S| = k$, apply Hölder with exponents all equal to $2k$:

$$\left\|\prod_{i \in S} \Delta_i\right\|_{L_2} \leq \prod_{i \in S} \|\Delta_i\|_{L_{2k}} = O\left(\prod_{i \in S} \sqrt{f_i}\right),$$

where we used $\|\Delta_i\|_{L_{2k}} = O(\sqrt{f_i})$ for $k \leq m$.

Let $c_i := \sqrt{f_i(n)}$. Summing over subsets S shows

$$\left\|\prod_{i=1}^m X_i - \prod_{i=1}^m \mu_i\right\|_{L_2} \leq A \left(\prod_{i=1}^m (1 + c_i) - 1\right) \leq A(e^\Xi - 1),$$

where $\Xi := \sum_{i=1}^m c_i$ and A is a constant depending only on m , $\{\mu_i\}$, and the moment constants (not on n). Hence

$$\text{Var}\left(\prod_{i=1}^m X_i\right) \leq \left\|\prod_{i=1}^m X_i - \prod_{i=1}^m \mu_i\right\|_{L_2}^2 = O(\Xi^2) = O\left(\left(\sum_{i=1}^m \sqrt{f_i(n)}\right)^2\right).$$

Since m is fixed, $(\sum_{i=1}^m \sqrt{f_i})^2 \leq m^2 \max_i f_i$, giving the claimed bound. \square

Corollary 1 (Higher moments of the centered product). *Fix $p \geq 1$. Under the hypotheses of Lemma 36, then*

$$\left\| \prod_{i=1}^m X_i - \prod_{i=1}^m \mathbb{E}[X_i] \right\|_{L_{2p}} \leq C_{p,m} \sum_{\emptyset \neq S \subseteq [m]} \left(\prod_{j \in S^c} |\mathbb{E}[X_j]| \right) \prod_{i \in S} \sqrt{f_i} = o(1),$$

and hence $\mathbb{E} \left| \prod_{i=1}^m X_i - \prod_{i=1}^m \mathbb{E}[X_i] \right|^{2p} = o(1)$.

Lemma 37 (Expectation of Product vs. Product of Expectations). *Fix $k \geq 2$. Let X_1, \dots, X_k be random variables. Assume:*

1. *Uniformly bounded means:* $\sup_{n,i} |\mathbb{E}[X_i]| \leq M < \infty$.
2. *Vanishing variances:* $\text{Var}(X_i) = f_i(n)$ with $f_i(n) \rightarrow 0$ as $n \rightarrow \infty$ for each i .
3. *Moment control up to order k :* For each i and every $p \in \{2, \dots, k\}$,

$$\mathbb{E} [|X_i - \mathbb{E}[X_i]|^p] \leq C_p \text{Var}(X_i)^{p/2},$$

with constants C_p .

Then for finite k , we have:

$$\left| \mathbb{E} \left[\prod_{i=1}^k X_i \right] - \prod_{i=1}^k \mathbb{E}[X_i] \right| = O \left(\left(\sum_{i=1}^k \sqrt{f_i(n)} \right)^2 \right) = O \left(\max_{1 \leq i \leq k} f_i(n) \right).$$

Proof. Set $\Delta_i := X_i - \mathbb{E}[X_i]$, so $\mathbb{E}\Delta_i = 0$, $\text{Var}(X_i) = \text{Var}(\Delta_i) = f_i(n)$, and by assumption

$$\|\Delta_i\|_{L_p} := (\mathbb{E} |\Delta_i|^p)^{1/p} \leq C_p^{1/p} f_i(n)^{1/2}, \quad p = 2, \dots, k.$$

Using the multilinearity of expectation,

$$\prod_{i=1}^k X_i = \prod_{i=1}^k (\mathbb{E}[X_i] + \Delta_i) = \sum_{S \subseteq [k]} \left(\prod_{i \in S} \Delta_i \right) \left(\prod_{j \notin S} \mathbb{E}[X_j] \right),$$

Thus,

$$\prod_{i=1}^k X_i - \prod_{i=1}^k \mathbb{E}[X_i] = \sum_{\emptyset \neq S \subseteq [k]} \left[\prod_{i \in S} \Delta_i \right] \prod_{j \notin S} \mathbb{E}[X_j].$$

Then taking the expectation and noting that $\prod_{j \notin S} \mathbb{E}[X_j]$ is a constant, we get

$$\mathbb{E} \left[\prod_{i=1}^k X_i \right] - \prod_{i=1}^k \mathbb{E}[X_i] = \sum_{\emptyset \neq S \subseteq [k]} \mathbb{E} \left[\prod_{i \in S} \Delta_i \right] \prod_{j \notin S} \mathbb{E}[X_j].$$

If $S = \{\ell\}$ then $\mathbb{E} \left[\prod_{i \in S} \Delta_i \right] = \mathbb{E}[\Delta_\ell] = 0$. Hence every singleton term vanishes exactly, and the sum begins at $|S| = 2$. From the bounded means assumption,

$$\left| \prod_{j \notin S} \mathbb{E}[X_j] \right| \leq M^{k-|S|}, \quad \forall S \subseteq [k].$$

Fix a nonempty subset S with $|S| = m \geq 2$. By generalized Hölder with all exponents equal to m (so $\sum_{i \in S} \frac{1}{m} = 1$),

$$\left| \mathbb{E} \left[\prod_{i \in S} \Delta_i \right] \right| \leq \prod_{i \in S} \|\Delta_i\|_{L_m} \leq \prod_{i \in S} \left(C_m^{1/m} f_i(n)^{1/2} \right) = C_m \prod_{i \in S} \sqrt{f_i(n)}.$$

Therefore, for every S with $|S| = m \geq 2$,

$$\left| \mathbb{E} \left[\prod_{i \in S} \Delta_i \right] \prod_{j \notin S} \mathbb{E}[X_j] \right| \leq M^{k-m} C_m \prod_{i \in S} \sqrt{f_i(n)}.$$

Let $c_i := \sqrt{f_i(n)} \geq 0$. Denote by

$$e_m(c_1, \dots, c_k) := \sum_{\substack{S \subseteq [k] \\ |S|=m}} \prod_{i \in S} c_i$$

the m -th elementary symmetric polynomial. Summing the bound from, we get

$$\left| \mathbb{E} \left[\prod_{i=1}^k X_i \right] - \prod_{i=1}^k \mathbb{E} X_i \right| \leq \sum_{m=2}^k M^{k-m} C_m e_m(c_1, \dots, c_k).$$

Let $M_\star := \max_{2 \leq m \leq k} M^{k-m} C_m$. Since $e_m \geq 0$ for $c_i \geq 0$,

$$\sum_{m=2}^k M^{k-m} C_m e_m \leq M_\star \sum_{m=2}^k e_m(c_1, \dots, c_k).$$

Recall the identity

$$\prod_{i=1}^k (1 + c_i) = \sum_{m=0}^k e_m(c_1, \dots, c_k) = 1 + \sum_{m=1}^k e_m(c_1, \dots, c_k),$$

so that $\sum_{m=2}^k e_m = \prod_{i=1}^k (1 + c_i) - 1 - \sum_{i=1}^k c_i$. Hence

$$\left| \mathbb{E} \left[\prod_{i=1}^k X_i \right] - \prod_{i=1}^k \mathbb{E} X_i \right| \leq M_\star \left(\prod_{i=1}^k (1 + c_i) - 1 - \sum_{i=1}^k c_i \right).$$

Let $\Xi := \sum_{i=1}^k c_i \rightarrow 0$ as $n \rightarrow \infty$. Since $\log(1 + u) \leq u$ for $u \geq 0$,

$$\prod_{i=1}^k (1 + c_i) = \exp \left(\sum_{i=1}^k \log(1 + c_i) \right) \leq \exp(\Xi).$$

Thus, the difference is at most $M_\star(e^\Xi - 1 - \Xi)$. By Taylor's theorem, $e^\Xi = 1 + \Xi + \frac{1}{2}\Xi^2 e^\xi$ for some $\xi \in [0, \Xi]$, so $e^\Xi - 1 - \Xi = \frac{1}{2}\Xi^2 e^\xi \leq \frac{1}{2}\Xi^2 e^\Xi$ (since $\xi \leq \Xi$ and $e^\xi \leq e^\Xi$). Therefore,

$$\left| \mathbb{E} \left[\prod_{i=1}^k X_i \right] - \prod_{i=1}^k \mathbb{E} X_i \right| \leq \frac{M_\star}{2} \Xi^2 e^\Xi = O(\Xi^2),$$

as $\Xi \rightarrow 0$ and $e^\Xi \rightarrow 1$. Since $\Xi = O\left(\sum_{i=1}^k \sqrt{f_i(n)}\right)$, we get the result. \square

Lemma 38 (Moment preservation under monomial \leftrightarrow Hermite change of basis). *Fix $M \in \mathbb{N}$ and degree $r \in \mathbb{N}$. Let*

$$\mathcal{M} := \{x^\gamma : \gamma \in \mathbb{N}^M, |\gamma| \leq r\}, \quad \mathcal{H} := \{\mathbf{H}_\alpha : \alpha \in \mathbb{N}^M, |\alpha| \leq r\},$$

with $\mathbf{H}_\alpha(x) = \prod_{j=1}^M H_{\alpha_j}(x_j)$ the probabilists' Hermite basis. For any (random) coefficients $\{a_\gamma\}_{|\gamma| \leq r}$ define the random polynomial $P(x) = \sum_{|\gamma| \leq r} a_\gamma x^\gamma$. Then there is a deterministic, invertible matrix $T = T(M, r)$ such that the Hermite coefficients $c = \{c_\alpha\}_{|\alpha| \leq r}$ in $P(x) = \sum_{|\alpha| \leq r} c_\alpha \mathbf{H}_\alpha(x)$ satisfy

$$c = T a.$$

Consequently, for any $p \geq 1$,

$$\|c_\alpha\|_{L_p} \leq \sum_{|\gamma| \leq r} |T_{\alpha\gamma}| \|a_\gamma\|_{L_p} \quad \text{for all } \alpha,$$

so if each $a_\gamma \in L_p$ then each $c_\alpha \in L_p$. Moreover, since T is invertible, the converse also holds: if each $c_\alpha \in L_p$ then each $a_\gamma \in L_p$.

Proof. In one dimension, each monomial admits a finite Hermite expansion $x^m = \sum_{j=0}^{\lfloor m/2 \rfloor} t_{m,j} H_{m-2j}(x)$ with deterministic coefficients $t_{m,j}$; in several dimensions, take tensor products to obtain $x^\gamma = \sum_{|\alpha| \leq |\gamma|} T_{\alpha\gamma} \mathbf{H}_\alpha(x)$. Ordering multi-indices by total degree yields a block upper-triangular, deterministic, invertible matrix $T = T(M, r)$. Linearity gives $c = T a$. For $p \geq 1$, Minkowski's inequality yields $\|c_\alpha\|_{L_p} = \left\| \sum_\gamma T_{\alpha\gamma} a_\gamma \right\|_{L_p} \leq \sum_\gamma |T_{\alpha\gamma}| \|a_\gamma\|_{L_p}$, so finiteness of all $\|a_\gamma\|_{L_p}$ implies finiteness of all $\|c_\alpha\|_{L_p}$. Invertibility gives the converse using $a = T^{-1}c$ and the same argument with T^{-1} . \square

G PROOF OF SPECIFIC CASES AND OVERFITTING

G.1 PROOF OF THEOREM 1.

Proof. We set $\alpha_Z = \alpha_A = \tilde{\alpha}_Z = \tilde{\alpha}_A = \alpha$, $\tilde{\theta} = \theta$, $\tilde{\tau} = \tau$ in the above Theorem 5 and note that it greatly simplifies each term. Algebra shows that for $c < 1$

$$\begin{aligned} \text{Bias} &= \tau_\varepsilon^2 \frac{c}{1-c} \frac{\theta^2}{d(\theta^2 + \tau^2)}, & \text{Variance} &= \alpha^2 \tau^2 \|\beta_*\|^2 + \tau_\varepsilon^2 \frac{c}{1-c} \left[1 - \frac{\theta^2}{d(\theta^2 + \tau^2)} \right], \\ \text{Data Noise} &= \alpha^2 \tau^2 \|\beta_*\|^2, & \text{Target Alignment} &= -2\alpha^2 \tau^2 \|\beta_*\|^2, \end{aligned}$$

While for $c > 1$, we can first send $d, n \rightarrow \infty$ and many terms become asymptotically 0. In the end, we get that:

$$\begin{aligned} \text{Bias} &= \alpha^2 \theta^2 (\beta_*^\top \mathbf{u})^2 \left(1 - \frac{1}{c} \right)^2 \left(\frac{\tau^2 c}{\theta^2 + \tau^2 c} \right)^2, & \text{Data Noise} &= \alpha^2 \tau^2 \|\beta_*\|^2, \\ \text{Variance} &= \alpha^2 \tau^2 \|\beta_*\|^2 \frac{1}{c} + \alpha^2 \tau^2 (\beta_*^\top \mathbf{u})^2 \frac{\theta^2}{\theta^2 + \tau^2 c} \left(1 - \frac{1}{c} \right) + \tau_\varepsilon^2 \frac{1}{c-1}. \\ \text{Target Alignment} &= -2\alpha^2 \tau^2 \left(\left(1 - \frac{1}{c} \right) \frac{\theta^2}{\theta^2 + \tau^2 c} (\beta_*^\top \mathbf{u})^2 + \|\beta_*\|^2 \frac{1}{c} \right), \end{aligned}$$

Adding these terms together, we see with simple algebra that many terms cancel or can be combined, establishing the stated formula. \square

G.2 PROOF OF THEOREM 2.

Proof. We set $\alpha_Z = \tilde{\alpha}_Z$, $\alpha_A = \tilde{\alpha}_A$, $\tilde{\theta} = \theta$, $\tilde{\tau} = \tau$, and send $d, n \rightarrow \infty$ in Theorem 5. Recall that $\Delta_c = \alpha_Z - \frac{\alpha_A}{c}$ and $\Delta_1 = \alpha_Z - \alpha_A$. Then some algebra shows that for $c < 1$,

$$\begin{aligned} \text{Bias} &= \theta^2 (\beta_*^\top \mathbf{u})^2 \Delta_1^2 \left(\frac{\tau^2}{\theta^2 + \tau^2} \right)^2, & \text{Data Noise} &= \alpha_A^2 \tau^2 \|\beta_*\|^2, \\ \text{Target Alignment} &= -2\alpha_A^2 \tau^2 \|\beta_*\|^2 - 2\alpha_A \tau^2 (\beta_*^\top \mathbf{u})^2 \Delta_1 \frac{\theta^2}{\theta^2 + \tau^2}, \\ \text{Variance} &= \alpha_A^2 \tau^2 \|\beta_*\|^2 + \tau_\varepsilon^2 \frac{c}{1-c} + \tau^2 (\beta_*^\top \mathbf{u})^2 \left[\frac{1}{1-c} \frac{\theta^4 + \theta^2 \tau^2 c}{(\theta^2 + \tau^2)^2} \Delta_1^2 + 2\alpha_A \Delta_1 \frac{\theta^2}{\theta^2 + \tau^2} \right]. \end{aligned}$$

For $c > 1$, we have that

$$\begin{aligned} \text{Bias} &= \theta^2 (\beta_*^\top \mathbf{u})^2 \Delta_c^2 \left(\frac{\tau^2 c}{\theta^2 + \tau^2 c} \right)^2, & \text{Data Noise} &= \alpha_A^2 \tau^2 \|\beta_*\|^2, \\ \text{Target Alignment} &= -2\alpha_A^2 \tau^2 \frac{\|\beta_*\|^2}{c} - 2\alpha_A \tau^2 (\beta_*^\top \mathbf{u})^2 \Delta_c \frac{\theta^2}{\theta^2 + \tau^2 c}, \\ \text{Variance} &= \alpha_A^2 \tau^2 \frac{\|\beta_*\|^2}{c} + \tau_\varepsilon^2 \frac{1}{c-1} + \tau^2 (\beta_*^\top \mathbf{u})^2 \frac{c}{1-c} \frac{\theta^2}{\theta^2 + \tau^2 c} \Delta_c^2. \end{aligned}$$

We proceed by adding these terms together and the results follow from algebra. \square

G.3 PROOF OF THEOREM 3.

Proof. We set $\tilde{\theta} = \theta$ and $\tilde{\tau} = \tau$ in Theorem 5 and have the regime of equal operator norm $\theta^2 = \gamma \tau^2$. Since we are interested in the limit $c \rightarrow \infty$, we only consider the overparameterized case $c > 1$. We first take the limit $d, n \rightarrow \infty$ and have that:

$$\begin{aligned} \text{Bias} &= \tau^2 (\beta_*^\top \mathbf{u})^2 \left(\sqrt{\gamma} (\tilde{\alpha}_Z - \alpha_Z) + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{c\sqrt{\gamma}}{\gamma + c} \right)^2, & \text{Data Noise} &= \tilde{\alpha}_A^2 \tau^2 \|\beta_*\|^2, \\ \text{Target Alignment} &= -2\tilde{\alpha}_A \tau^2 \left(\left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\gamma}{\gamma + c} (\beta_*^\top \mathbf{u})^2 + \alpha_A \frac{\|\beta_*\|^2}{c} \right), \end{aligned}$$

$$\text{Variance} = \tau^2 \alpha_A^2 \frac{\|\beta_*\|^2}{c} + \tau^2 (\beta_*^\top \mathbf{u})^2 \frac{c}{(c-1)} \frac{\gamma}{\gamma+c} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + \tau_\varepsilon^2 \left(\frac{1}{c-1} \right).$$

The rest follows from simple calculus: if $\tilde{\alpha}_Z \neq \alpha_Z$, $\gamma = \omega_c(1)$, and $\beta_*^\top \mathbf{u} \neq 0$, the bias will diverge and other terms are controlled, yielding catastrophic. If $\tilde{\alpha}_Z = \alpha_Z$, $\omega_c(1) \leq \gamma \leq o_c(c^2)$, and $\beta_*^\top \mathbf{u} \neq 0$, a similar thing happens. In other cases, all of these terms are controlled and become finite values in the limit $\lim_{c \rightarrow \infty} \mathcal{R}_c - \tau_\varepsilon^2$, giving us tempered overfitting.

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = \begin{cases} \tilde{\alpha}_A^2 \tau^2 \|\beta_*\|^2 & \beta \perp u \\ \tau^2 [\gamma \tilde{\alpha}_Z^2 (\beta_*^\top \mathbf{u})^2 + \tilde{\alpha}_A^2 \|\beta_*\|^2] & \beta \not\perp u, \gamma = \Theta_c(1) \\ \infty & \alpha_Z \neq \tilde{\alpha}_Z, \beta_* \not\perp u, \gamma = \omega(1) \\ \infty & \alpha_Z = \tilde{\alpha}_Z, \beta_* \not\perp u, \omega(1) \leq \gamma \leq o(c^2) \\ \tau^2 \left[\left(\frac{\phi}{(\phi+1)^2} \alpha_Z^2 - 2\tilde{\alpha}_A \alpha_Z \right) (\beta_*^\top \mathbf{u})^2 + \alpha_A^2 \|\beta_*\|^2 \right] & \alpha_Z = \tilde{\alpha}_Z, \beta_* \not\perp u, \gamma = \phi c^2 \\ \tau^2 [(\alpha_Z^2 - 2\tilde{\alpha}_A \alpha_Z) (\beta_*^\top \mathbf{u})^2 + \alpha_A^2 \|\beta_*\|^2] & \alpha_Z = \tilde{\alpha}_Z, \beta_* \not\perp u, \gamma = \omega(c^2) \end{cases}$$

□

G.4 PROOF OF THEOREM 4.

Proof. We start with the first part and assume that $\alpha_Z \neq \tilde{\alpha}_Z$. Similarly, we have that $\tilde{\theta} = \theta$ and $\tilde{\tau} = \tau$ in Theorem 5. To achieve equal Frobenius norm, we set $\theta^2 = d\tau^2$ and send $d, n \rightarrow \infty$ so several terms would vanish.

In particular, for $c < 1$, we have that

$$\text{Bias} = \theta^2 (\beta_*^\top \mathbf{u})^2 \left(\tilde{\alpha}_Z - \alpha_Z + (\alpha_Z - \alpha_A) \frac{\tau^2}{\theta^2 + \tau^2} \right)^2 = \tau^2 (\beta_*^\top \mathbf{u})^2 \left(\sqrt{d}(\tilde{\alpha}_Z - \alpha_Z) + (\alpha_Z - \alpha_A) \frac{\sqrt{d}}{d+1} \right)^2,$$

It is clear that this term becomes ∞ since the term inside the parentheses scales with d . Note that the variance and data noise are non-negative, and target alignment is controlled. We have that $\mathcal{R}_c = \infty$ for $c \in (0, 1)$.

For $c > 1$, the same logic follows, and we also note that:

$$\text{Bias} = \theta^2 (\beta_*^\top \mathbf{u})^2 \left(\tilde{\alpha}_Z - \alpha_Z + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\tau^2 c}{\theta^2 + \tau^2 c} \right)^2 = \tau^2 (\beta_*^\top \mathbf{u})^2 \left(\sqrt{d}(\tilde{\alpha}_Z - \alpha_Z) + \left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{\sqrt{dc}}{d+c} \right)^2,$$

which scales with d with other terms controlled. Hence, $\mathcal{R}_c = \infty$ for all $c \neq 1$.

Now assume that $\alpha_Z = \tilde{\alpha}_Z$. Since we are interested in $c \rightarrow \infty$, we only consider $c > 1$. First, from algebra and taking the limit for d, n , we have that:

$$\text{Bias} = \tau^2 (\beta_*^\top \mathbf{u})^2 \left(\left(\alpha_Z - \frac{\alpha_A}{c} \right) \frac{c\sqrt{d}}{d+c} \right)^2 \rightarrow 0, \quad \text{Data Noise} = \tilde{\alpha}_A^2 \tau^2 \|\beta_*\|^2,$$

$$\text{Target Alignment} = -2\tilde{\alpha}_A \tau^2 \left(\left(\alpha_Z - \frac{\alpha_A}{c} \right) (\beta_*^\top \mathbf{u})^2 + \alpha_A \frac{\|\beta_*\|^2}{c} \right),$$

$$\text{Variance} = \tau^2 \alpha_A^2 \frac{\|\beta_*\|^2}{c} + \tau^2 (\beta_*^\top \mathbf{u})^2 \frac{c}{(c-1)} \left(\alpha_Z - \frac{\alpha_A}{c} \right)^2 + \tau_\varepsilon^2 \left(\frac{1}{c-1} \right).$$

We now take $c \rightarrow \infty$ and many terms vanish in this limit, yielding:

$$\lim_{c \rightarrow \infty} \mathcal{R}_c = -2\tilde{\alpha}_A \alpha_Z \tau^2 (\beta_*^\top \mathbf{u})^2 + \tau^2 (\beta_*^\top \mathbf{u})^2 \alpha_Z^2 + \tilde{\alpha}_A^2 \tau^2 \|\beta_*\|^2 = \tau^2 [(\beta_*^\top \mathbf{u})^2 (\alpha_Z^2 - 2\tilde{\alpha}_A \alpha_Z) + \|\beta_*\|^2 \tilde{\alpha}_A^2].$$

□

Proposition 3 (Non–existence of a canceling scale parameter). *Let $\alpha_A, \alpha_Z > 0$ be fixed scalars, let $\mathbf{u}, \beta_* \in \mathbb{R}^d$ be fixed vectors, and set*

$$a := \|\beta_*\|^2 > 0, \quad b := (\beta_*^\top \mathbf{u})^2 \in [0, a].$$

For every positive real number ϕ define

$$f(\phi) = \alpha_A^2 a + \left(\alpha_Z^2 \left(1 + \frac{1}{\phi} \right) - 2\alpha_Z \alpha_A \right) b.$$

Then

$$f(\phi) > 0 \quad \text{for all } \phi > 0.$$

Consequently the equation $f(\phi) = 0$ has no solution with $\phi \in (0, \infty)$.

Proof. If $b = 0$ (i.e. β_* is orthogonal to \mathbf{u}) we have $f(\phi) = \alpha_A^2 a > 0$, so no positive ϕ can cancel the expression. Hence assume $b > 0$.

Writing $r := b/a \in (0, 1]$ we obtain

$$f(\phi) = a \left[\alpha_A^2 + \alpha_Z(\alpha_Z - 2\alpha_A) r + \frac{\alpha_Z^2 r}{\phi} \right]. \quad (*)$$

Since $r \leq 1$,

$$\alpha_A^2 + \alpha_Z(\alpha_Z - 2\alpha_A) r \geq \alpha_A^2 + \alpha_Z(\alpha_Z - 2\alpha_A) = (\alpha_A - \alpha_Z)^2 \geq 0.$$

Thus the square bracket in (*) is the sum of a non–negative term and a strictly positive term. □