
IN-CONTEXT BENIGN OVERFITTING: A FEATURE-SELECTION MODEL IN IN-CONTEXT LINEAR REGRESSION

Anonymous authors

Paper under double-blind review

ABSTRACT

In in-context learning (ICL), a frozen pre-trained model solves tasks by conditioning on a prompt of a few input–output examples, without gradient updates. If the task was present in pretraining but the particular prompt sequence was not, the resulting in-distribution generalization is retrieval-based ICL. Learning-based ICL instead reflects out-of-distribution generalization: the model succeeds on prompts generated by a novel task. Empirically, both forms improve with scale. By analogy to benign overfitting in supervised learning, we call this in-context benign overfitting: larger models more faithfully memorize the pretraining tasks (improving retrieval ICL) while also generalizing better to novel tasks (improving learning ICL). We prove that this phenomenon already arises in a minimal in-context linear-regression feature-selection model. In contrast, standard in-context linear-regression models exhibit a retrieval–learning tradeoff, where the emergence of learning-based ICL coincides with degraded retrieval-based performance.

1 INTRODUCTION

In-context learning (ICL) has established itself as a cornerstone capability of large language models, where the model solves tasks at inference time by conditioning on a prompt of input–output demonstrations, without any updates to model weights (Brown et al., 2020). While the prompt contains only a handful of examples, the model also draws on information acquired during pretraining. As a result, there are two qualitatively different settings in which a model can succeed: it can perform well because the underlying task was present during pretraining, or it can perform well on prompts generated by a task that is novel relative to pretraining. Pan et al. (2023) introduced a useful distinction between these two modes of ICL, which they call *retrieval* and *learning*. Concretely, let $\mathcal{D}_{\text{task}}^{\text{train}} = \{\theta_1, \theta_2, \dots, \theta_M\}$ denote the set of tasks present during pretraining. For a task θ , a prompt consists of T demonstrations and a query, $X_\theta = [(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_T, y_T), \mathbf{x}_q]$, where labels $y_t = f_\theta(\mathbf{x}_t)$ are generated through a fixed mapping f that is shared across demonstrations inside the prompt. This mapping is parameterized by the task θ and applies to all $t \in [T]$ as well as the query \mathbf{x}_q . The two modes are then defined as follows:

- **Retrieval mode (ID generalization).** We sample $\theta \sim \mathcal{D}_{\text{task}}^{\text{train}}$, so while the prompt X_θ is unseen during pretraining, the underlying task is drawn from the pretraining task set. We evaluate whether the model correctly predicts the query label y_q . Success in this mode is consistent with the model having *memorized* the pretraining tasks and *retrieving* the correct task mechanism from its parameters. This mode, therefore, probes task-level in-distribution (ID) generalization.
- **Learning mode (OOD generalization).** We sample $\theta \sim \mathcal{D}_{\text{task}}$ from a distribution that differs from $\mathcal{D}_{\text{task}}^{\text{train}}$. We again evaluate prediction on y_q . Since the task is novel relative to pretraining, the model cannot rely on task memorization alone and must *learn* from demonstrations. This mode probes task-level out-of-distribution (OOD) generalization.

Hereon, we will refer to these as ID and OOD generalization, respectively.

Understanding these two modes has motivated a surge of work in synthetic, fully controlled settings. These studies train transformers from scratch on sequences generated by a finite task set, for example in-context linear regression (Garg et al., 2022; Raventos et al., 2023; Carroll et al., 2025) or Markov-chains (Park et al., 2025), and then evaluate ID and OOD generalization as discussed above. This line

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

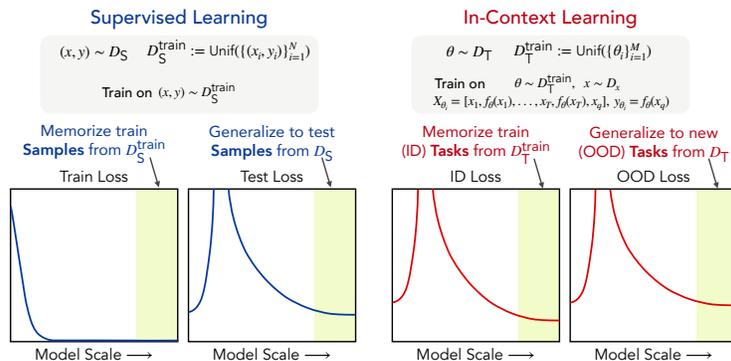


Figure 1: **Benign Overfitting in Supervised vs. In-Context Learning.** Left illustrates benign overfitting commonly observed in standard supervised learning setups, where large models far beyond the interpolation threshold (green band) memorize the train *samples* while simultaneously generalizing to new ones. We introduce and formalize an in-context analogue of this phenomenon (right): large models (in green) memorize *tasks* seen during (pre)training while simultaneously generalizing to novel out-of-distribution tasks.

of work has uncovered several qualitative phenomena, including task-diversity thresholds (Raventos et al., 2023; Park et al., 2025) for the emergence of *good* OOD generalization and transient dynamics (Singh et al., 2023; Carroll et al., 2025) in how the two modes evolve over the course of training. We defer a detailed discussion and pointers to the relevant papers to Appendix A.

One question, however, remains conceptually unsettled: *how should ID and OOD generalization evolve with model scale?* In the era of scaling laws, the prevailing intuition is that larger models should excel in both modes. Indeed, this has been observed experimentally: Wei et al. (2023) (Figs. 2 and 3) show that in large language models, both modes tend to improve with scale: larger models more accurately solve tasks encountered during pretraining and also generalize better to novel tasks at inference time.

Contributions. Drawing an analogy to the literature on benign overfitting in supervised learning (Bartlett et al., 2020; Hastie et al., 2020; Belkin et al., 2020), we frame the scaling behavior of ID and OOD generalization in ICL as what we call *in-context benign overfitting*.

In-Context Benign Overfitting

Definition 1. *In-context benign overfitting is the property whereby large (overparameterized) models memorize ID ICL tasks, achieving low risk on prompts from $D_{\text{task}}^{\text{train}}$, while simultaneously generalizing to novel OOD ICL tasks (from D_{task}).*

This lifts the classical notion of benign overfitting from the sample level to the task level: whereas standard benign overfitting concerns memorization of noisy training *samples*, in-context benign overfitting concerns memorization of pretraining *tasks*. See Fig. 1 for an illustration.

We ground our framework in the minimal setup of in-context linear regression (Garg et al., 2022), employing a simplified linear attention model consistent with recent theoretical works (Zhang et al., 2023; Lu et al., 2025; Wu et al., 2023; Zhang et al., 2025). To explicitly model the effect of scale, we incorporate a feature subsampling mechanism inspired by the double descent literature (Hastie et al., 2020; Belkin et al., 2020). Using data generated from this model, we study the two canonical predictors extensively studied in classical statistical learning: the minimum-norm interpolator (in the overparameterized regime) and the least-squares solution (in the underparameterized regime). We analyze these predictors in the proportional high-dimensional asymptotic limit, where the number of (pre)training tasks, number of (pre)training prompts, number of in-context demonstrations, number of learnable parameters, and the ambient feature dimension grow strictly in proportion (Lu et al., 2025). In this regime, we derive sharp analytical characterizations of the ID and OOD risks, demonstrating that our theoretical predictions accurately track the empirical risk curves and faithfully capture the phenomenon of in-context benign overfitting.

2 PROBLEM SETUP

Data Model. We study in-context linear regression first introduced by Garg et al. (2022). Input sequences are interleaved pairs $X = [\mathbf{x}_1, y_1, \dots, \mathbf{x}_T, y_T, \mathbf{x}_{T+1}]$, where features $\mathbf{x}_t \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{x}})$ and labels $y_t = \mathbf{w}^\top \mathbf{x}_t + \epsilon_t$ for a task vector $\mathbf{w} \sim \mathcal{D}_{\text{task}} = \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{w}})$ and noise $\epsilon_t \sim \mathcal{N}(0, \sigma_n^2)$. We refer to $\mathbf{x}_q := \mathbf{x}_{T+1}$ as the query vector, $y_q := y_{T+1}$ as the query label, and the number of demonstrations T as the context-length. Given the T demonstrations $(\mathbf{x}_t, y_t)_{t=1}^T$, the goal is predicting y_q .

Training Data. We construct a training set by sampling a finite pool of M task vectors $\{\mathbf{w}_i\}_{i=1}^M \sim \mathcal{D}_{\text{task}}$ independently. The training task distribution is uniform over this pool, i.e. $\mathcal{D}_{\text{task}}^{\text{train}} = \text{unif}\{\mathbf{w}_1, \dots, \mathbf{w}_M\}$. We generate N training sequences X_i by sampling $\mathbf{w}_i \sim \mathcal{D}_{\text{task}}^{\text{train}}$ and constructing query labels $y_q^i = \mathbf{w}_i^\top \mathbf{x}_q + \epsilon_q^i$, for all $i \in [N]$.

Learning Model. To solve the ICL linear regression problem, we employ a simplified model inspired by linear attention. In its vanilla form, this model has been used extensively in prior work as a tractable framework for theoretically analyzing ICL (Zhang et al., 2025; Lu et al., 2025; Wu et al., 2023; Zhang et al., 2024; Frei & Vardi, 2024).

Our key modeling contribution is a feature-selection formulation inspired by the double descent literature in supervised learning (Hastie et al., 2020; Belkin et al., 2020). We vary model capacity by restricting the learner to a subset of original features $\mathcal{S} \subseteq [d]$ of size $p = |\mathcal{S}|$, sampled uniformly without replacement. The model first constructs a feature embedding $\mathbf{h}_{\mathcal{S}_p}(X) = \text{vec}(\mathbf{x}_q[\mathcal{S}], \hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}]^\top) \in \mathbb{R}^{p^2}$ where $\hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}] = \frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t[\mathcal{S}]$ is the *averaging estimator* restricted to \mathcal{S} . The prediction is then linear in these features $f_{\theta, \mathcal{S}}(X) = \theta^\top \mathbf{h}_{\mathcal{S}_p}(X)$, with trainable parameters $\theta \in \mathbb{R}^{p^2}$.

Given training sequences $\{X_i\}_{i=1}^N$, we obtain θ by minimizing the empirical risk $\mathcal{L}_{\text{emp}}(\theta) := \frac{1}{N} \|\mathbf{H}_{\mathcal{S}_p} \theta - \mathbf{y}\|^2$. The solution $\hat{\theta}_{\mathcal{S}_p}$ follows the problem geometry:

- **Underparameterized regime:** For full column rank $\mathbf{H}_{\mathcal{S}_p}$, $\hat{\theta}_{\mathcal{S}_p}$ is the unique least-squares solution $\hat{\theta}_{\mathcal{S}_p} := \arg \min_{\theta} \|\mathbf{y} - \mathbf{H}_{\mathcal{S}_p} \theta\|^2$.
- **Overparameterized regime:** For full row rank $\mathbf{H}_{\mathcal{S}_p}$, gradient descent initialized at zero converges to the minimum-norm interpolator (Hastie et al., 2020; Bartlett et al., 2020), i.e. $\hat{\theta}_{\mathcal{S}_p} := \arg \min_{\theta} \|\theta\|$ s.t. $\mathbf{H}_{\mathcal{S}_p} \theta = \mathbf{y}$.

This characterization allows us to study generalization of $\hat{\theta}_{\mathcal{S}_p}$ across both regimes as the model size p varies. We evaluate the trained predictor $\hat{\theta}_{\mathcal{S}_p}$ on two types of generalization:

- In-distribution (ID):** performance on training tasks, i.e., $\mathbf{w} \sim \mathcal{D}_{\text{task}}^{\text{train}}$;
- Out-of-distribution (OOD):** performance on novel tasks, i.e., $\mathbf{w} \sim \mathcal{D}_{\text{task}}$.

To define these metrics, we first introduce the *task-specific mean-square-error loss*, which measures the prediction error of $\hat{\theta}_{\mathcal{S}_p}$ on sequences generated from a fixed task vector \mathbf{w} , $\mathcal{L}(\hat{\theta}_{\mathcal{S}_p}; \mathbf{w}) := \mathbb{E}_{X, \epsilon} [(\hat{\theta}_{\mathcal{S}_p}^\top \mathbf{h}_{\mathcal{S}_p}(X) - y_q)^2]$, where the expectation is over the feature vectors and noise in the sequence X , i.e., $\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{x}_q \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{x}})$ and $\epsilon_1, \dots, \epsilon_q \sim \mathcal{N}(0, \sigma_n^2)$. The ID and OOD risks are then defined by averaging over the respective task distributions: $\mathcal{L}_{\text{ID}}(\hat{\theta}_{\mathcal{S}_p}) := \mathbb{E}_{\mathcal{S}_p} \mathbb{E}_{\mathbf{w} \sim \mathcal{D}_{\text{task}}^{\text{train}}} \mathcal{L}(\hat{\theta}_{\mathcal{S}_p}; \mathbf{w})$, $\mathcal{L}_{\text{OOD}}(\hat{\theta}_{\mathcal{S}_p}) := \mathbb{E}_{\mathcal{S}_p} \mathbb{E}_{\mathbf{w} \sim \mathcal{D}_{\text{task}}} \mathcal{L}(\hat{\theta}_{\mathcal{S}_p}; \mathbf{w})$. In both cases, the outer expectation $\mathbb{E}_{\mathcal{S}_p}$ averages over the random feature subset \mathcal{S} , ensuring that we evaluate performance independently of any particular realization of it.

3 SHARP ASYMPTOTICS OF ID AND OOD RISKS

We characterize the ID and OOD risks in the joint high-dimensional limit where the ambient dimension d , context length T , model size p , number of training tasks M , and number of training sequences N all grow to infinity as follows: $N, T, M, d, p \rightarrow \infty$ with

$$\frac{N}{d^2} \rightarrow \nu, \quad \frac{M}{d} \rightarrow \mu, \quad \frac{T}{d} \rightarrow \tau, \quad \frac{p}{d} \rightarrow \rho. \quad (1)$$

Throughout this section, we assume isotropic features and task vectors, i.e., $\Sigma_{\mathbf{x}} = \sigma_x^2 \mathbb{I}_d$ and $\Sigma_{\mathbf{w}} = \sigma_w^2 \mathbb{I}_d$. As we will show this setting already suffices to establish the in-context benign overfitting phenomenon. Further, following Lu et al. (2025), we fix $\sigma_x = 1/\sqrt{d}$ and $\sigma_w = 1$.

Analysis Framework. Our derivations leverage a Gaussian equivalence principle (Mei & Montanari, 2019; Goldt et al., 2020). To find the asymptotic risks, we analyze a Linear Gaussian Equivalent Problem (LGP) which matches the first and second-order statistics of our original model (see App. B). The predictive power of this equivalence is rooted in the broader phenomenon of Gaussian universality (Goldt et al., 2020). While a formal proof of equivalence for our problem is left for future work, we analyze the LGP using the Convex Gaussian Min-Max Theorem (CGMT) and verify experimentally (Sec. 4) that its theoretical risk curves accurately track the original model’s performance.

Next, we present our main results. See Appendix B for a detailed overview of the derivations.

Theorem 1 (Overparameterized Regime, $\nu < \rho^2$). *Under the asymptotic scaling (Eq. (1)), as $d \rightarrow \infty$, the risks of minimum-norm estimator $\hat{\theta}_{S_p}$ for the LGP (Def. 3) converge in probability to:*

$$\mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) \xrightarrow{P} \left(1 - \frac{\nu}{\rho^2}\right) \rho(1+ac)(1-c(1+a+ac)) + (1+a+ac)(\sigma_n^2 + \rho c + 1 - \rho), \quad \mathcal{L}_{\text{ID}}(\hat{\theta}_{S_p}) \xrightarrow{P} \frac{\nu}{\alpha^2} \bar{\beta}_*^2,$$

where $c := \frac{\sigma_n^2 + 1}{\tau}$ and the scalars a and $\bar{\beta}_*$ are the unique positive solutions to Eq. (6) (see App. B).

Theorem 2 (Underparameterized Regime, $\nu > \rho^2$). *Under the asymptotic scaling (Eq. (1)), as $d \rightarrow \infty$, the risks for the least-squares estimator $\hat{\theta}_{S_p}$ for the LGP (Def. 3) converge in probability to:*

$$\mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) \xrightarrow{P} m_c(\kappa_\infty^2(1+c) - 2c^2\rho) + (1+c)\rho c^2 m'_c + \sigma_n^2 + 1 - \rho + \rho c, \quad \mathcal{L}_{\text{ID}}(\hat{\theta}_{S_p}) \xrightarrow{P} \frac{\nu}{\rho^2} \kappa_\infty^2,$$

where κ_∞ is a scalar defined in Eq. (9) (see App. B), $m_c := m_\gamma(z_c)$ is the Stieltjes transform of the Marchenko-Pastur law with parameter $\gamma = \rho/\mu$ evaluated at $z_c = -c$, and m'_c denotes the derivative of $m_\gamma(\cdot)$ at z_c .

4 EXPERIMENTAL RESULTS

We present experimental results demonstrating in-context benign overfitting and validating our theoretical predictions. We follow the setup in Sec. 2 unless otherwise noted; see App. I for details.

Fig. 2 compares our derived asymptotic risk curves (Theorems 1 and 2) against empirical loss values. We plot both ID and OOD risk as a function of the relative model scale ρ^2/ν (specifically, we fix N and vary p). We observe that the theoretical predictions closely match the empirical observations. The vertical asymptote at $\rho^2/\nu = 1$ marks the phase transition between the underparameterized ($\rho^2/\nu < 1$) and overparameterized ($\rho^2/\nu > 1$) regimes. As model scale increases, both ID and OOD loss exhibit double descent. Importantly, increasing the model scale causes the ID loss to converge to a significantly lower value compared to the underparameterized regime, while the OOD loss converges to a value comparable to the underparameterized minimum. This demonstrates that the model benignly overfits the training tasks—memorizing ID tasks while simultaneously generalizing to OOD tasks, capturing the essence of in-context benign overfitting. See App. I for experiments illustrating effect of task diversity (M , Fig. 4 in App.), and experiments with non-isotropic Σ_x, Σ_w (Fig. 5, 6 in App.).

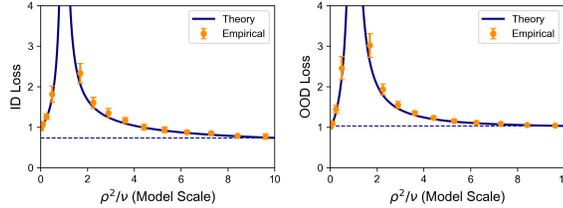


Figure 2: In-Context Benign Overfitting. Comparison of our theoretical risk curves (solid lines) vs. empirical observations (points) for ID Loss (left) and OOD Loss (right) as a function of model scale ρ^2/ν (with fixed ν). The horizontal dashed lines denote the asymptotic value in the overparameterized regime. The vertical asymptote at $\rho^2/\nu = 1$ marks the interpolation threshold. In the overparameterized regime ($\rho^2/\nu > 1$), the model exhibits in-context benign overfitting: it effectively memorizes pre-training tasks (low ID risk) while simultaneously yielding a small risk on OOD tasks.

5 CONCLUSION

We introduced *in-context benign overfitting*—a task-level analogue of classical benign overfitting, and provided sharp asymptotic characterizations showing that it arises in a minimal feature-selection model of in-context linear regression. Moving forward, rigorously establishing the Gaussian equivalence principle in our setting and extending the analysis to richer architectures (e.g., softmax attention) and other task families (e.g., Markov chains) remain natural next steps. More broadly, our results suggest a reassuring message: scaling up need not force a tradeoff between solving known tasks and generalizing to new ones, even though prior theoretical models of ICL might suggest otherwise (Lu et al., 2025; Raventos et al., 2023; Park et al., 2025) (See Fig. 3 in App. for discussion).

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

REFERENCES

- Ehsan Abbasi, Fariborz Salehi, and Babak Hassibi. Universality in learning from linear measurements. In *Advances in Neural Information Processing Systems*, pp. 12372–12382, 2019.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=LziniAXEI9>.
- Danil Akhtiamov, Reza Ghane, and Babak Hassibi. One-bit quantization for random features models. *arXiv preprint arXiv:2510.16250*, 2025.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? Investigations with linear models. In *Int. Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- 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, and Ji Xu. Two models of double descent for weak features. *SIAM Journal on Mathematics of Data Science*, 2(4):1167–1180, January 2020. ISSN 2577-0187. doi: 10.1137/20m1336072. URL <http://dx.doi.org/10.1137/20M1336072>.
- David Bosch, Ashkan Panahi, and Babak Hassibi. Precise asymptotic analysis of deep random feature models. In *The Thirty Sixth Annual Conference on Learning Theory*, pp. 4132–4179. PMLR, 2023.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf.
- Liam Carroll, Jesse Hoogland, Matthew Farrugia-Roberts, and Daniel Mufet. Dynamics of transient structure in in-context linear regression transformers, 2025. URL <https://arxiv.org/abs/2501.17745>.
- Xiangyu Chang, Yingcong Li, Samet Oymak, and Christos Thrampoulidis. Provable benefits of overparameterization in model compression: From double descent to pruning neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6974–6983, 2021.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Training dynamics of multi-head softmax attention for in-context learning: Emergence, convergence, and optimality, 2024a. URL <https://arxiv.org/abs/2402.19442>.
- Siyu Chen, Heejune Sheen, Tianhao Wang, and Zhuoran Yang. Unveiling induction heads: Provable training dynamics and feature learning in transformers. *arXiv preprint arXiv:2409.10559*, 2024b.
- Yatin Dandi, Ludovic Stephan, Florent Krzakala, Bruno Loureiro, and Lenka Zdeborová. Universality laws for gaussian mixtures in generalized linear models. *Advances in Neural Information Processing Systems*, 36:54754–54768, 2023.
- Zeyu Deng, Abla Kammoun, and Christos Thrampoulidis. A model of double descent for high-dimensional binary linear classification. *arXiv preprint arXiv:1911.05822*, 2019.
- Puneesh Deora, Bhavya Vasudeva, Tina Behnia, and Christos Thrampoulidis. In-context occam’s razor: How transformers prefer simpler hypotheses on the fly. In *Second Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=ZSMnX3LBva>.

-
- 270 Ezra Edelman, Nikolaos Tsilivis, Benjamin L. Edelman, Eran Malach, and Surbhi Goel. The evolution
271 of statistical induction heads: In-context learning markov chains. In *The Thirty-eighth Annual*
272 *Conference on Neural Information Processing Systems*, 2024. URL [https://openreview.](https://openreview.net/forum?id=qaRT6QTIqJ)
273 [net/forum?id=qaRT6QTIqJ](https://openreview.net/forum?id=qaRT6QTIqJ).
274
- 275 Spencer Frei and Gal Vardi. Trained transformer classifiers generalize and exhibit benign overfitting
276 in-context. *arXiv preprint arXiv:2410.01774*, 2024.
- 277 Deqing Fu, Tian qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn to achieve second-order
278 convergence rates for in-context linear regression. In *The Thirty-eighth Annual Conference on*
279 *Neural Information Processing Systems*, 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=L8h6cozcbn)
280 [id=L8h6cozcbn](https://openreview.net/forum?id=L8h6cozcbn).
281
- 282 Shivam Garg, Dimitris Tsipras, Percy Liang, and Gregory Valiant. What can transformers learn
283 in-context? a case study of simple function classes. In Alice H. Oh, Alekh Agarwal, Danielle
284 Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
285 URL <https://openreview.net/forum?id=flNZJ2eOet>.
- 286 Federica Gerace, Florent Krzakala, Bruno Loureiro, Ludovic Stephan, and Lenka Zdeborová. Gaus-
287 sian universality of perceptrons with random labels. *Physical Review E*, 109(3):034305, 2024.
288
- 289 Reza Ghane, Danil Akhiamov, and Babak Hassibi. Universality in transfer learning for linear models.
290 *Advances in Neural Information Processing Systems*, 37:125729–125779, 2024.
- 291 Reza Ghane, Anthony Bao, Danil Akhiamov, and Babak Hassibi. Gaussian universality for diffusion
292 models. *IEEE Signal Processing Letters*, 33:116–120, 2025.
- 293 Sebastian Goldt, Galen Reeves, Marc Mézard, Florent Krzakala, and Lenka Zdeborová. The gaussian
294 equivalence of generative models for learning with two-layer neural networks. *arXiv preprint*
295 *arXiv:2006.14709*, 2020.
296
- 297 Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J Tibshirani. Surprises in high-
298 dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.
299
- 300 Trevor Hastie, Andrea Montanari, Saharon Rosset, and Ryan J. Tibshirani. Surprises in high-
301 dimensional ridgeless least squares interpolation, 2020. URL [https://arxiv.org/abs/](https://arxiv.org/abs/1903.08560)
302 [1903.08560](https://arxiv.org/abs/1903.08560).
- 303 Hong Hu and Yue M Lu. Universality laws for high-dimensional learning with random features.
304 *arXiv preprint arXiv:2009.07669*, 2020.
305
- 306 Hong Hu and Yue M Lu. Universality laws for high-dimensional learning with random features.
307 *IEEE Transactions on Information Theory*, 69(3):1932–1964, 2022.
- 308 Tengyuan Liang and Pragya Sur. A precise high-dimensional asymptotic theory for boosting and
309 min-11-norm interpolated classifiers. *arXiv preprint arXiv:2002.01586*, 2020.
310
- 311 Bruno Loureiro, Gabriele Sicuro, Cédric Gerbelot, Alessandro Pocco, Florent Krzakala, and Lenka
312 Zdeborová. Learning gaussian mixtures with generalized linear models: Precise asymptotics in
313 high-dimensions. *Advances in Neural Information Processing Systems*, 34:10144–10157, 2021.
- 314 Yue M. Lu, Mary I. Letey, Jacob A. Zavatore-Veth, Anindita Maiti, and Cengiz Pehlevan. Asymptotic
315 theory of in-context learning by linear attention, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2405.11751)
316 [2405.11751](https://arxiv.org/abs/2405.11751).
317
- 318 Song Mei and Andrea Montanari. The generalization error of random features regression: Precise
319 asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.
- 320 Theodor Misiakiewicz and Basil Saeed. A non-asymptotic theory of kernel ridge regression: deter-
321 ministic equivalents, test error, and gcv estimator. *arXiv preprint arXiv:2403.08938*, 2024.
322
- 323 Andrea Montanari and Phan-Minh Nguyen. Universality of the elastic net error. In *2017 IEEE*
International Symposium on Information Theory (ISIT), pp. 2338–2342. IEEE, 2017.

324 Andrea Montanari and Basil N Saeed. Universality of empirical risk minimization. In *Conference on*
325 *Learning Theory*, pp. 4310–4312. PMLR, 2022.

326

327 Andrea Montanari, Feng Ruan, Youngtak Sohn, and Jun Yan. The generalization error of max-margin
328 linear classifiers: High-dimensional asymptotics in the overparametrized regime. *arXiv preprint*
329 *arXiv:1911.01544*, 2019.

330 Andrea Montanari, Feng Ruan, Basil Saeed, and Youngtak Sohn. Universality of max-margin
331 classifiers. *arXiv preprint arXiv:2310.00176*, 2023.

332

333 Samet Oymak and Joel A Tropp. Universality laws for randomized dimension reduction, with
334 applications. *Information and Inference: A Journal of the IMA*, 7(3):337–446, 2018.

335

336 Jane Pan, Tianyu Gao, Howard Chen, and Danqi Chen. What in-context learning “learns” in-context:
337 Disentangling task recognition and task learning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki
338 Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 8298–
339 8319, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/
340 v1/2023.findings-acl.527. URL <https://aclanthology.org/2023.findings-acl.527/>.

341

342 Ashkan Panahi and Babak Hassibi. A universal analysis of large-scale regularized least squares
343 solutions. *Advances in Neural Information Processing Systems*, 30, 2017.

344

345 Core Francisco Park, Ekdeep Singh Lubana, and Hidenori Tanaka. Algorithmic phases of in-context
346 learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL
<https://openreview.net/forum?id=XgH1wfHSX8>.

347

348 Luca Pesce, Florent Krzakala, Bruno Loureiro, and Ludovic Stephan. Are gaussian data all you
349 need? the extents and limits of universality in high-dimensional generalized linear estimation. In
350 *International Conference on Machine Learning*, pp. 27680–27708. PMLR, 2023.

351

352 Nived Rajaraman, Marco Bondaschi, Ashok Vardhan Makkuva, Kannan Ramchandran, and Michael
353 Gastpar. Transformers on markov data: Constant depth suffices. In *The Thirty-eighth Annual*
354 *Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=5uG9tp3v2q>.

355

356 Allan Raventos, Mansheej Paul, Feng Chen, and Surya Ganguli. Pretraining task diversity and the
357 emergence of non-bayesian in-context learning for regression. In *Thirty-seventh Conference on*
358 *Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=BtAz4a5xDg>.

359

360 Aaditya K. Singh, Stephanie C. Y. Chan, Ted Moskovitz, Erin Grant, Andrew M. Saxe, and Felix
361 Hill. The transient nature of emergent in-context learning in transformers, 2023. URL <https://arxiv.org/abs/2311.08360>.

362

363 Mihailo Stojnic. A framework to characterize performance of lasso algorithms. *arXiv preprint*
364 *arXiv:1303.7291*, 2013.

365

366 Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Society, 2023.

367

368 Christos Thrampoulidis, Samet Oymak, and Babak Hassibi. Regularized linear regression: A
369 precise analysis of the estimation error. In Peter Grünwald, Elad Hazan, and Satyen Kale (eds.),
370 *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine*
371 *Learning Research*, pp. 1683–1709, Paris, France, 03–06 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v40/Thrampoulidis15.html>.

372

373 Christos Thrampoulidis, Ehsan Abbasi, and Babak Hassibi. Precise error analysis of regularized
374 m -estimators in high dimensions. *IEEE Transactions on Information Theory*, 64(8):5592–5628,
375 2018. doi: 10.1109/TIT.2018.2840720.

376

377 Johannes von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev,
Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent.
arXiv preprint arXiv:2212.07677, 2022.

378 Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao
379 Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning
380 differently, 2023. URL <https://arxiv.org/abs/2303.03846>.
381

382 Garrett G Wen, Hong Hu, Yue M Lu, Zhou Fan, and Theodor Misiakiewicz. When does gaussian
383 equivalence fail and how to fix it: Non-universal behavior of random features with quadratic
384 scaling. *arXiv preprint arXiv:2512.03325*, 2025.

385 Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett.
386 How many pretraining tasks are needed for in-context learning of linear regression? *arXiv preprint*
387 *arXiv:2310.08391*, 2023.

388 Ruiqi Zhang, Spencer Frei, and Peter L. Bartlett. Trained transformers learn linear models in-context,
389 2023.
390

391 Ruiqi Zhang, Jingfeng Wu, and Peter Bartlett. In-context learning of a linear transformer block:
392 Benefits of the mlp component and one-step gd initialization. *Advances in Neural Information*
393 *Processing Systems*, 37:18310–18361, 2024.

394 Yedi Zhang, Aaditya K. Singh, Peter E. Latham, and Andrew Saxe. Training dynamics of in-context
395 learning in linear attention, 2025. URL <https://arxiv.org/abs/2501.16265>.
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

432	CONTENTS	
433		
434		
435		
436	1 Introduction	1
437		
438	2 Problem Setup	3
439		
440		
441	3 Sharp Asymptotics of ID and OOD risks	3
442		
443	4 Experimental Results	4
444		
445		
446	5 Conclusion	4
447		
448	A Related Work	10
449		
450		
451	B Overview of Theory	10
452	B.1 Linear Gaussian Equivalent Model	10
453	B.1.1 Overparameterized regime	11
454	B.1.2 Underparameterized regime	14
455		
456		
457	C Gaussian Equivalent Data Model	14
458		
459		
460	D Forming the Auxiliary Optimization Problem	15
461	D.1 Solving the Auxiliary Optimization.	16
462	D.2 Simplifying Per-Task Loss	17
463		
464		
465	E Proofs for Min Norm AO	18
466	E.1 Wishart Matrix Asymptotics	18
467	E.2 General Notation	18
468	E.3 Asymptotic limit of $S(\tau, \beta)$	19
469	E.4 Asymptotic limit of the quadratic forms	20
470	E.5 First Order Optimality Conditions of $D(\bar{u}, \bar{\beta})$	21
471		
472		
473		
474		
475	F Proof for LS AO	24
476		
477		
478	G Asymptotic Risk	24
479	G.1 Asymptotic Risk of Minimum-Norm LGP	24
480	G.2 Asymptotic Risk of LS LGP	33
481		
482		
483	H Helper Lemmas	35
484		
485		
	I Additional Results and Details of Experimental Settings	40

A RELATED WORK

To deconstruct the mechanics of in-context learning, a growing body of literature has turned to controlled synthetic environments where Transformers are pre-trained from scratch on canonical function classes. Linear regression has served as a primary testbed for these investigations (Garg et al., 2022; Raventos et al., 2023; Akyürek et al., 2023; von Oswald et al., 2022; Ahn et al., 2023; Zhang et al., 2025), alongside sequence modeling tasks defined by Markov chains (Park et al., 2025; Edelman et al., 2024; Rajaraman et al., 2024; Deora et al., 2025). A central theme in this line of inquiry is algorithmic discovery: determining whether trained Transformers implement known estimation procedures, such as gradient descent variants for regression (von Oswald et al., 2022; Ahn et al., 2023; Fu et al., 2024) or induction-head mechanisms for Markov processes (Edelman et al., 2024; Rajaraman et al., 2024; Chen et al., 2024b). Additionally, a growing body of work also investigates the training dynamics of in-context learning for linear regression, specifically examining the optimization dynamics for both one-layer linear attention (Zhang et al., 2025; 2023; 2024) and softmax attention (Chen et al., 2024a).

Retrieval and Learning Modes of ICL. Most relevant to our work are studies investigating the dual nature of ICL: the retrieval and learning modes (Pan et al., 2023). Empirical works involving transformers trained on finite task sets have identified task diversity thresholds—critical points where the model transitions from retrieval (good ID generalization) to learning (good OOD tasks) (Raventos et al., 2023; Park et al., 2025). While Wu et al. (2023) provided an initial theoretical quantification of this transition in linear attention for in-context linear regression, our analysis is most closely related to and partly inspired by Lu et al. (2025). They derive precise risk asymptotics in a proportional limit similar to ours (Eq. (1)), but with a crucial distinction: their setup is restricted to $\rho = 1$ (utilizing all features). This difference is fundamental; as we demonstrate in Section 4 (Fig. 3), fixing $\rho = 1$ couples model scale with the ambient dimension. Under this constraint, the in-context benign overfitting phenomenon disappears, highlighting the necessity of the feature-selection mechanism to decouple model capacity from data dimensionality.

Finally, we note that Frei & Vardi (2024) also coin the term *in-context benign overfitting* though in a fundamentally different context. Their work focuses on in-context binary classification and interprets *overfitting* through the lens of label noise: they show that models can fit prompts with flipped labels (noise) while still generalizing to the query. This is strictly analogous to classical benign overfitting over samples. In contrast, our framework applies the concept to the task level, describing models that overfit (memorize) the pre-training *task* distribution while generalizing to novel distributions.

B OVERVIEW OF THEORY

B.1 LINEAR GAUSSIAN EQUIVALENT MODEL

Our asymptotic analysis leverages powerful machinery from high-dimensional statistics that provides sharp characterizations for solutions of *Linear Gaussian Problems* (LGPs).

Definition 2 (Linear Gaussian MSE Problem). *An LGP* $(\theta_*, \Sigma_g, \Sigma_\epsilon)$ with $\Sigma_g, \Sigma_\epsilon$ positive definite, is a linear regression problem with squared loss on N samples $\{(\mathbf{g}_i, y_i)\}_{i=1}^N$ generated as follows: features $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_{\tilde{p}})$ and labels $y_i = \mathbf{g}_i^\top \Sigma_g^{-1/2} \theta_* + \epsilon_i$, where $\epsilon = (\epsilon_1, \dots, \epsilon_N)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma_\epsilon)$. Let \tilde{p} denote the dimension of the regressor $\theta_* \in \mathbb{R}^{\tilde{p}}$, and let $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_N]^\top \in \mathbb{R}^{N \times \tilde{p}}$ be the feature matrix. The estimator is defined as

$$\hat{\theta} = \Sigma_g^{-1/2} \hat{\mathbf{a}} + \theta_*,$$

with the error vector $\hat{\mathbf{a}}$ given as follows.

- If $\tilde{p} < N$: the least-squares solution

$$\hat{\mathbf{a}} = \mathbf{a}_{\text{LS}} := \arg \min_{\mathbf{a}} \|\epsilon - \mathbf{G}\mathbf{a}\|^2.$$

- If $\tilde{p} > N$: the minimum-norm interpolator

$$\hat{\mathbf{a}} = \mathbf{a}_{\text{MN}} := \arg \min_{\mathbf{a}} \|\mathbf{a}\| \quad \text{s.t.} \quad \mathbf{G}\mathbf{a} = \epsilon.$$

LGPs are central to our analysis for two reasons: (i) Well-established tools yield sharp asymptotic characterizations for the generalization error of squared-loss minimizers in LGPs; (ii) A *Gaussian equivalence principle* allows us to transfer these LGP results to our original (non-Gaussian) problem setup in Section 2. We will empirically verify the accuracy of these predictions.

The LGP solves least-squares or minimum-norm optimization as in our original setup. However, there are two key advantages in the LGP setting. First, the transition threshold between regimes is precisely $\tilde{p} = N$: as \tilde{p} and N grow proportionally, the Gaussian feature matrix \mathbf{G} is almost surely full column rank if and only if $\tilde{p} < N$. Second, and most importantly, the optimization problems (LS or min-norm) for LGP involve only an isotropic Gaussian matrix and independent Gaussian noise, which admits sharp asymptotic characterizations. Instead, the original problems involve non-Gaussian features $\mathbf{h}_{S_p}(X_i)$ and correlated labels $y_q^i = \mathbf{w}_i^\top \mathbf{x}_q + \epsilon$.

We now explicitly relate LGPs to our original setup of interest detailed in Section 2.

Definition 3 (Linear Gaussian Equivalent Problem). *Consider N features $\mathbf{h}_{S_p}(X_i) \in \mathbb{R}^{p^2}$ as in original model with labels $y_q^i = \mathbf{w}_i^\top \mathbf{x}_q + \epsilon$ for $\mathbf{w}_i \sim \mathcal{D}_{task}^{train}$ and the corresponding least-squares estimates. The equivalent LGP $(\boldsymbol{\theta}_*, \boldsymbol{\Sigma}_g, \boldsymbol{\Sigma}_\epsilon)$ is given by setting $\boldsymbol{\Sigma}_g \leftarrow \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}}$, $\boldsymbol{\theta}_* \leftarrow \frac{1}{N} \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}}^{-1} \boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}}$ and $\boldsymbol{\Sigma}_\epsilon \leftarrow \text{diag}(\sigma_{n_1}, \dots, \sigma_{n_N})$ with*

$$\boldsymbol{\Sigma}_{\mathbf{H}_{S_p}} := \frac{1}{N} \mathbb{E}_{X, \epsilon} \left[\mathbf{H}_{S_p}^\top \mathbf{H}_{S_p} \right] \in \mathbb{R}^{p^2 \times p^2}, \quad (2a)$$

$$\boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}} = \mathbb{E}_{X, \epsilon} [\mathbf{H}_{S_p}^\top \mathbf{y}] \quad (2b)$$

$$\sigma_{n_i}^2 := \sigma_x^2 \|\mathbf{w}_i\|^2 + \sigma_n^2 - \boldsymbol{\theta}_*^\top \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_* \quad i \in [N]. \quad (2c)$$

See App. C for explicit formulas for the expectations above.

Equivalence is in the strong sense that the asymptotic *empirical distribution* of the LGP estimator $\hat{\boldsymbol{\theta}}$ matches that of the original estimator $\hat{\boldsymbol{\theta}}_S$ (Goldt et al., 2020; Hu & Lu, 2020). Thus, it suffices to characterize the former, since knowing the latter suffices to compute ID and OOD limiting risks.

The LGP parameters are chosen so that the first and second-order statistics of the LGP samples $\{(\mathbf{g}_i, y_i)\}_{i \in [N]}$ match the first and second-order statistics (with respect to X, ϵ) of the original data $\{(\mathbf{h}_S(X_i), y_q^i)\}_{i \in [N]}$. The principle that Gaussian systems matching a non-Gaussian system to first and second order are predictive of properties of the latter is broadly known as *universality*. Its roots lie in classical random matrix theory results on eigendistributions of random ensembles (Tao, 2023). The specific form relevant here, where the “system” refers to a random optimization problem and “properties” to quantities such as its optimal cost or empirical distribution of its solution, has been developed within the high-dimensional statistics literature, from early applications in compressed sensing (e.g. (Montanari & Nguyen, 2017; Oymak & Tropp, 2018; Panahi & Hassibi, 2017; Abbasi et al., 2019)) to more recent machine learning settings (e.g. (Hastie et al., 2019; Mei & Montanari, 2019; Goldt et al., 2020; Bosch et al., 2023; Montanari & Saeed, 2022; Hu & Lu, 2022; Ghane et al., 2024; Akhtiamov et al., 2025)). The terminology *Gaussian equivalence principle* was introduced by Mei & Montanari (2019); Goldt et al. (2020), who applied it to study linear regression with random features. Since then, numerous extensions have applied and confirmed the gaussian universality principle in various contexts, e.g., (Ghane et al., 2024; Dandi et al., 2023; Gerace et al., 2024; Pesce et al., 2023; Misiakiewicz & Saeed, 2024; Ghane et al., 2025). It is also understood that universality does not hold always and very recent work has extended the methodology to handle such cases (Wen et al., 2025). Here, we identify another instance where the Gaussian equivalence principle applies. Figure 2 empirically demonstrates this. Formalizing this equivalence rigorously is left for future work, as our focus here is on deriving predictions and establishing in-context benign overfitting; a formal proof could likely build on recent formal universality equivalence proofs such as (Hu & Lu, 2022; Montanari et al., 2023).

B.1.1 OVERPARAMETERIZED REGIME

We begin with the overparameterized regime $\rho^2 > \nu$ (i.e., $p^2 > N$), analyzing the following minimum-norm linear Gaussian equivalent optimization:

$$\min_{\mathbf{a}} \frac{1}{2} \|\boldsymbol{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2 \quad \text{s.t.} \quad \mathbf{G} \mathbf{a} = \mathbf{D}_w \mathbf{g}', \quad (\text{Primal})$$

where $\Sigma_{\mathbf{H}_{S_p}}$ and θ_* are as in Definition 3; $\mathbf{D}_w \in \mathbb{R}^{N \times N}$ is a diagonal matrix with entries $(\mathbf{D}_w)_{ii} = \sigma_{n_i}$ with σ_{n_i} as given also in Definition 3; $\mathbf{G} \in \mathbb{R}^{N \times p^2}$ has i.i.d. standard Gaussian entries; $\mathbf{g}' \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_N)$ is independent of \mathbf{G} . Note that $\Sigma_{\mathbf{H}_{S_p}}$, θ_* , and \mathbf{D}_w depend on the pretraining task vectors $\{\mathbf{w}_1, \dots, \mathbf{w}_M\}$ sampled from $\mathcal{D}_{\text{task}}^{\text{train}}$, while \mathbf{G} and \mathbf{g}' are independent of these.

We will derive an asymptotic characterization for the empirical distribution of the solution to (Primal). To this end, we apply the convex Gaussian min-max theorem (CGMT) (Stojnic, 2013; Thrampoulidis et al., 2015). The application of the CGMT to analyze minimum-norm linear Gaussian problems (Definition 2) appears in various recent works, e.g., (Deng et al., 2019; Montanari et al., 2019; Loureiro et al., 2021; Chang et al., 2021; Liang & Sur, 2020). While the CGMT framework provides a general recipe, its application still requires problem-specific analysis, which we carry out here since, to the best of our knowledge, no plug-and-play result exists for our specific setting. We outline the key calculations below and defer remaining derivations to the appendix.

Our starting point is what is the following Auxiliary Optimization (AO):

$$\min_{\mathbf{a}} \max_{\lambda} \|\mathbf{a}\| \lambda^\top \mathbf{h} + \|\lambda\| \mathbf{a}^\top \mathbf{g} - \lambda^\top \mathbf{D}_w \mathbf{g}' + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \theta_*\|^2 \quad (3)$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_{p^2})$ and $\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_N)$ are independent of each other and of all other quantities. The CGMT guarantees that the optimal cost of Eq. (3) converges to the same asymptotic limit as that of (Primal), which sets the stage for establishing that the empirical distribution of the AO solution converges to the same limit as that of the primal (Montanari et al., 2019; Chang et al., 2021).

Solving the AO. Following standard CGMT machinery, we reduce the high-dimensional AO in Eq. (3) to a two-dimensional min-max optimization:

$$\min_{u>0} \max_{\beta \geq 0} \frac{\beta u}{2} + \beta \frac{\sum_{i=1}^n \sigma_{n_i}^2}{2u} + \frac{\|\theta_*\|^2}{2} - \frac{1}{2} (\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \theta_*)^\top \mathbf{M}^{-1} (\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \theta_*), \quad (\text{AO})$$

where we define the shorthand

$$\mathbf{M} := \mathbf{M}(\beta, u) := \frac{\beta N}{u} \mathbb{I}_{p^2} + \Sigma_{\mathbf{H}_{S_p}}^{-1}.$$

Direct differentiation yields that for any fixed $u > 0$, the objective in (AO) is *strictly* concave on $\beta > 0$ (except on a measure-zero set). Also, for any fixed $\beta > 0$, it is strictly convex in $u > 0$. This implies that the saddle point (u_*, β_*) is unique whenever it satisfies $\beta_* > 0$.

Moreover, the same reduction that yields (AO) also shows the following closed form for \mathbf{a}_* , the minimizer of Eq. (3):

$$\mathbf{a}_* = \mathbf{a}_*(\beta_*, u_*) = \mathbf{M}(\beta_*, u_*)^{-1} (\beta_* \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \theta_*), \quad (4)$$

where (β_*, u_*) is the unique saddle point of (AO) (when $\beta_* > 0$). In particular, conditional on this saddle point, \mathbf{a}_* is Gaussian, with mean and covariance determined by (β_*, τ_*) . As we show next, in the high-dimensional limit, (β_*, τ_*) converge to deterministic values, and thus \mathbf{a}_* is asymptotically Gaussian with mean and covariance determined by the limiting saddle point.

Deterministic reduction. In the regime specified in Eq. (1), the AO objective converges to a deterministic function of the scalar parameters (u, β) . To state this limit compactly, we introduce deterministic quantities defined as the in-probability limits of the following random terms:

$$c_\infty := \text{plim}_{d \rightarrow \infty} \frac{1}{d^2} \sum_{i=1}^N \sigma_{n_i}^2, \quad (5a)$$

$$s_\infty(\bar{u}, \bar{\beta}) := \text{plim}_{d \rightarrow \infty} \frac{1}{d^2} \text{tr}(\bar{\mathbf{M}}^{-1}), \quad (5b)$$

$$v_\infty(\bar{u}, \bar{\beta}) := \text{plim}_{d \rightarrow \infty} \frac{1}{d} \bar{\theta}_*^\top \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \bar{\mathbf{M}}^{-1} \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \bar{\theta}_*, \quad (5c)$$

where $\bar{\beta}, \bar{u}, \bar{\Sigma}_{\mathbf{H}_{S_p}}, \bar{\boldsymbol{\theta}}_*, \bar{\mathbf{M}}$ are the normalized versions of $\beta, u, \Sigma_{\mathbf{H}_{S_p}}, \boldsymbol{\theta}_*, \mathbf{M}$, chosen so that these quantities remain $O(1)$ (see App. D.1). We prove these convergences (with respect to the randomness in \mathbf{g} and the task vectors $\{\mathbf{w}_i\}_{i \in [N]}$) in App. E. With these definitions, the AO objective admits the following deterministic limit

$$D(\bar{u}, \bar{\beta}) := \frac{\bar{\beta}\bar{u}}{2} + \frac{\bar{\beta}}{2\bar{u}}c_\infty - \frac{\bar{\beta}^2}{2}s_\infty(\bar{u}, \bar{\beta}) - \frac{1}{2}v_\infty(\bar{u}, \bar{\beta}),$$

Moreover, since D is strictly convex-concave on $\{\bar{u} > 0, \bar{\beta} > 0\}$, the limiting saddle point $(\bar{\beta}_*, \bar{u}_*)$ is unique and characterized by first-order optimality conditions, which, after algebraic simplification, yield:

$$\begin{aligned} \bar{u}_* &= \frac{\bar{\beta}_* \nu}{a}, \quad \text{where} \quad \left(1 - \frac{\nu}{\rho^2}\right) a = m_a, \\ \bar{\beta}_*^2 &= \frac{a^2(\sigma_n^2 + 1 - \rho + \rho c) + \rho(1 - a^2 c^2)m_a - \rho a z_a^2 m'_a}{\rho^2 \left(1 - \frac{\nu}{\rho^2} - \frac{m'_a}{a^2}\right)}. \end{aligned} \quad (6)$$

Here, $m_a := m_\gamma(z_a)$ is the Stieltjes transform of the Marchenko-Pastur law with $\gamma := p/M = \rho/\mu$ (using the scaling defined in Eq. (1)) evaluated at point $z_a := -1/a - c$. See Definition 4 in the Appendix. In addition, $m'_a := m'_\gamma(z_a)$ denotes the derivative of $m_\gamma(z)$ evaluated at z_a , and we have used the shorthand $c := \frac{\sigma_n^2 + 1}{\tau}$.

The appearance of the Marchenko-Pastur Stieltjes transform follows from the spectral structure of $\Sigma_{\mathbf{H}_{S_p}}$: Its explicit form (Eq. (10) in the App.) shows that $\Sigma_{\mathbf{H}_{S_p}}$ can be expressed in terms of blocks of the form $\mathbb{I}_p + \zeta \frac{1}{M} \sum_{i=1}^M \mathbf{w}_i[\mathbf{S}]\mathbf{w}_i[\mathbf{S}]^\top$ (for a scalar ζ independent of $\{\mathbf{w}_i\}$). Because $\mathbf{w}_i[\mathbf{S}]$ are Gaussian, the matrix $\frac{1}{M} \sum_{i=1}^M \mathbf{w}_i[\mathbf{S}]\mathbf{w}_i[\mathbf{S}]^\top$ is Wishart, whose empirical spectral distribution converges to the Marchenko-Pastur law in the proportional limit. Consequently, trace terms such as $\frac{1}{p} \text{tr}(\mathbb{I}_p + \zeta \frac{1}{M} \sum_{i=1}^M \mathbf{w}_i[\mathbf{S}]\mathbf{w}_i[\mathbf{S}]^\top)^{-1}$ converge to quantities given in terms of the Stieljes transform.

Evaluation of risk. We are now ready to evaluate the asymptotic limit of the ID and OOD risks. Fix a task vector \mathbf{w} . We first express the per-task loss (Eq. (7)) in terms of the error vector \mathbf{a} ; see Lemma 1 for the proof:

$$\begin{aligned} \mathcal{L}(\mathbf{a}; \mathbf{w}) &= \underbrace{\mathbf{a}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a}}_{\text{Term-1(a)}} + \underbrace{\frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top \bar{\Sigma}(\mathbf{w}) \bar{\boldsymbol{\theta}}_*}_{\text{Term-1(b)}} \\ &+ 2 \underbrace{\frac{1}{\sqrt{d}} \mathbf{a}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\boldsymbol{\theta}}_*}_{\text{Term-1(c)}} + \frac{\|\mathbf{w}\|^2}{d} + \sigma_n^2 \\ &- 2 \underbrace{\left(\frac{1}{\sqrt{d}} \mathbf{a}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top (\mathbf{w} \otimes \mathbf{w}) \right)}_{\text{Term-2}}. \end{aligned} \quad (7)$$

Here $\Sigma(\mathbf{w}) := \mathbb{E}_{X, \epsilon}[\mathbf{h}_{S_p}(X)\mathbf{h}_{S_p}(X)^\top]$ and $\bar{\Sigma}(\mathbf{w})$ denotes its normalized version; likewise, $\bar{\Sigma}_{\mathbf{H}_{S_p}}$ is the normalized version of $\Sigma_{\mathbf{H}_{S_p}}$. All normalizations are chosen so that each term in Eq. (7) is $O(1)$ under the scaling in Eq. (1).

Then, to compute the limits for the ID and OOD losses, we take the expectation of Eq. (7) over $\mathcal{D}_{\text{task}}^{\text{train}}$ and $\mathcal{D}_{\text{task}}$, respectively. We then substitute the AO's prediction for the error vector \mathbf{a} (derived in Eq. (4)), normalized and evaluated at the deterministic optimal scalars $(\bar{\beta}_*, \bar{u}_*)$:

$$\mathbf{a} = \frac{1}{d} \bar{\mathbf{M}}(\bar{\beta}_*, \bar{u}_*)^{-1} \left(\bar{\beta}_* \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_* \right), \quad (8)$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_{p^2})$ is independent of the environment variables $\bar{\Sigma}_{\mathbf{H}_{S_p}}, \bar{\boldsymbol{\theta}}_*, \{\sigma_{n_i}\}_{i=1}^N$. Finally, we derive closed-form limits for each component of Eq. (7) using Eq. (8). Detailed computations for the OOD loss limits (Terms 1(a)-(c) and Term 2) are provided in Lemmas 4 and 5, and the corresponding ID loss limits are derived in Lemmas 6 and 7.

With these we arrive at our first main result.

Theorem 3 (Asymptotic Risk of Minimum-Norm LGP). *Consider the minimum norm estimator $\hat{\theta}_{S_p}$ for the linear Gaussian equivalent problem in Definition 3 under the overparameterized regime ($\nu < \rho^2$). Under the asymptotic scaling defined in Eq. (1), as $d \rightarrow \infty$, the ID and OOD risks converge in probability to the following deterministic limits:*

$$\begin{aligned}\mathcal{L}_{OOD}(\hat{\theta}_{S_p}) &\xrightarrow{P} \left(1 - \frac{\nu}{\rho^2}\right) \rho(1+ac)(1-c(1+a+ac)) \\ &\quad + (1+a+ac)(\sigma_n^2 + \rho c + 1 - \rho), \\ \mathcal{L}_{ID}(\hat{\theta}_{S_p}) &\xrightarrow{P} \frac{\nu}{a^2} \bar{\beta}_*^2,\end{aligned}$$

where $c := \frac{\sigma_n^2 + 1}{\tau}$ and the scalars a and $\bar{\beta}_*$ are the unique positive solutions to Eq. (6).

B.1.2 UNDERPARAMETERIZED REGIME

In the underparameterized regime (where $\nu > \rho^2$), the estimator $\hat{\theta}_{S_p}$ is the unique least-squares solution. We analyze this setting using the same framework applied in Appendix B.1.1 for the overparameterized case. Since the derivation is analogous (in fact, considerably simpler due to the strong convexity of the objective and simpler analysis of the AO) we defer the intermediate steps to App. F and present the final limits directly. The asymptotic behavior in this regime is governed by a scalar κ_∞ , which arises as the high-dimensional limit of the dual variable in the auxiliary optimization (recall Eq. (5)):

$$\kappa_\infty := \sqrt{\frac{c_\infty}{\nu \left(\frac{\nu}{\rho^2} - 1\right)}}. \quad (9)$$

Theorem 4 (Asymptotic Risk of LS LGP). *Consider the least-squares estimator $\hat{\theta}_{S_p}$ for the linear Gaussian equivalent problem in the underparameterized regime ($\nu > \rho^2$). Under the asymptotic scaling in Eq. (1), as $d \rightarrow \infty$, the ID and OOD risks converge in probability to:*

$$\begin{aligned}\mathcal{L}_{OOD}(\hat{\theta}_{S_p}) &\xrightarrow{P} m_c(\kappa_\infty^2(1+c) - 2c^2\rho) + (1+c)\rho c^2 m'_c \\ &\quad + \sigma_n^2 + 1 - \rho + \rho c, \\ \mathcal{L}_{ID}(\hat{\theta}_{S_p}) &\xrightarrow{P} \frac{\nu}{\rho^2} \kappa_\infty^2,\end{aligned}$$

where κ_∞ is defined in Eq. (9), $c := \frac{\sigma_n^2 + 1}{\tau}$, and $m_c := m_\gamma(z_c)$ is the Stieltjes transform of the Marchenko-Pastur law with parameter $\gamma = \rho/\mu$ evaluated at $z_c = -c$. In addition, m'_c denotes the derivative of $m_\gamma(\cdot)$ at z_c .

C GAUSSIAN EQUIVALENT DATA MODEL

Here, we discuss the moment-matching process for the first and second moments involving features \mathbf{g} and labels y in Definition 3 with those from the original data model discussed in Section 2. Recall that for the Gaussian equivalent data model, the features are i.i.d $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{H}_{S_p}})$ for all $i \in [N]$, where $\Sigma_{\mathbf{H}_{S_p}}$ matches the second moment of the original feature matrix \mathbf{H}_{S_p} by definition as $\Sigma_{\mathbf{H}_{S_p}} := \frac{1}{N} \mathbb{E}[\mathbf{H}_{S_p}^\top \mathbf{H}_{S_p}]$. Below, we show that the first moments also match, i.e., $\mathbb{E}[\mathbf{g}] = \mathbb{E}[\mathbf{h}_{S_p}(X)]$.

We compute the first and second moments of feature matrix \mathbf{H}_{S_p} by finding the expected value of each feature row. By definition, we have

$$\mathbb{E}[\mathbf{h}_{S_p}(X)] = \mathbb{E}[\text{vec}(\mathbf{x}_q[S] \hat{\mathbf{w}}_{\text{avg}}[S])^\top] = \mathbb{E}[\hat{\mathbf{w}}_{\text{avg}}[S]] \otimes \mathbb{E}[\mathbf{x}_q[S]] = \mathbf{0},$$

where the second equality uses the independence of the query vector from the rest of the context, and the third equality uses $\mathbf{x}_q \sim \mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbb{I}_d)$.

Next, we compute the second moment of \mathbf{H}_{S_p} . We have

$$\begin{aligned}\Sigma_{\mathbf{H}_{S_p}} &:= \frac{1}{n} \mathbb{E} \left[\mathbf{H}_{S_p}^\top \mathbf{H}_{S_p} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{h}_i \mathbf{h}_i^\top] \\ &= \frac{\sigma_x^4}{TM} \sum_{i=1}^M \left[(\sigma_x^2 \|\mathbf{w}_i\|^2 + \sigma_n^2) \mathbb{I}_{p^2} + \sigma_x^2 (T+1) (\mathbf{w}_i[\mathcal{S}] \mathbf{w}_i[\mathcal{S}]^\top) \otimes \mathbb{I}_p \right].\end{aligned}\quad (10)$$

Here, the first equality uses the independence of the rows of \mathbf{H}_{S_p} , and the second equality follows by using Lemma 16.

Recall that the labels in the Gaussian equivalent problem are generated as $y_i = \mathbf{g}_i^\top \boldsymbol{\theta}_* + \epsilon_i$, where we select $\boldsymbol{\theta}_*$ and noise ϵ_i such that the moments, $\mathbb{E}[y]$, $\boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}} := \mathbb{E}[\mathbf{H}_{S_p}^\top \mathbf{y}]$ and $\mathbb{E}[\mathbf{y} \mathbf{y}^\top]$ also match with that of the original problem.

First, in the original model, $\mathbb{E}[y] = 0$. In the Gaussian equivalent data model, $\mathbb{E}[y] = \mathbb{E}[\mathbf{g}]^\top \boldsymbol{\theta}_* + \mathbb{E}[\epsilon] = \mathbb{E}[\epsilon]$. To match, we require $\mathbb{E}[\epsilon] = 0$.

Next, as the labels y are independent, $\mathbb{E}[\mathbf{y} \mathbf{y}^\top]$ is a diagonal matrix, and we only need to look at the entries $\mathbb{E}[y^2]$.

In the original model, we have

$$\mathbb{E}[y^2] = \mathbb{E}[(\mathbf{w}^\top \mathbf{x} + \epsilon)^2] = \sigma_x^2 \|\mathbf{w}\|^2 + \sigma_n^2.$$

For the Gaussian equivalent data model,

$$\mathbb{E}[y^2] = \boldsymbol{\theta}_*^\top \Sigma_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_* + \mathbb{E}[\epsilon^2].$$

Therefore, for the Gaussian equivalent data model $\epsilon_i \sim \mathcal{N}(0, \sigma_{n_i}^2)$, where

$$\sigma_{n_i}^2 = \sigma_x^2 \|\mathbf{w}_i\|^2 + \sigma_n^2 - \boldsymbol{\theta}_*^\top \Sigma_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_* =: d_i^2.\quad (11)$$

Here, \mathbf{w}_i corresponds to the task vector for sample $i \in [n]$.

Recall that $\mathbf{G} \in \mathbb{R}^{N \times p^2}$ has i.i.d. entries $G_{ij} \sim \mathcal{N}(0, 1)$, so that $\mathbf{G} \Sigma_{\mathbf{H}_{S_p}}^{1/2}$ has i.i.d. rows distributed as $\mathcal{N}(0, \Sigma_{\mathbf{H}_{S_p}})$. Therefore, the remaining second moment-matching condition gives

$$\boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}} = \mathbb{E}[\Sigma_{\mathbf{H}_{S_p}}^{1/2} \mathbf{G}^\top \mathbf{G} \Sigma_{\mathbf{H}_{S_p}}^{1/2} \boldsymbol{\theta}_*] = \Sigma_{\mathbf{H}_{S_p}}^{1/2} \mathbb{E}[\mathbf{G}^\top \mathbf{G}] \Sigma_{\mathbf{H}_{S_p}}^{1/2} \boldsymbol{\theta}_* = N \Sigma_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_*,$$

which implies

$$\boldsymbol{\theta}_* = \frac{1}{N} \Sigma_{\mathbf{H}_{S_p}}^{-1} \boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}}.\quad (12)$$

Here, using Lemma 16, we have

$$\begin{aligned}\boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}} &= \mathbb{E} \left[\sum_{i=1}^n \mathbf{h}_i y_i \right] = \sum_{i=1}^M \frac{n}{M} \mathbb{E}[\mathbf{h}_i y_i] \\ &= \frac{n}{M} \sum_{i=1}^M \sigma_x^4 \mathbf{w}_i[\mathcal{S}] \otimes \mathbf{w}_i[\mathcal{S}].\end{aligned}\quad (13)$$

D FORMING THE AUXILIARY OPTIMIZATION PROBLEM

Min norm using Gaussian equivalent. First, recall the minimum norm defined using the Gaussian equivalent data. We have

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta}\|^2 \quad \text{subject to} \quad \mathbf{G}(\boldsymbol{\theta} - \boldsymbol{\theta}_*) = \mathbf{D}_w \mathbf{g}',\quad (14)$$

where \mathbf{D}_w is a diagonal matrix with the i^{th} entry equal to σ_{n_i} for sample $i \in [n]$. Lastly, we use $\mathbf{g}' \sim \mathcal{N}(\mathbf{0}, \mathbb{I})$ to denote noise that is independent of entries in $\mathbf{G} \in \mathbb{R}^{N \times p^2}$. Here, the rows of \mathbf{G} , $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{H}_{S_p}})$ are i.i.d. sampled for all $i \in [N]$.

810 Changing variable

$$811 \quad \Sigma_{\mathbf{H}_{S_p}}^{1/2} (\boldsymbol{\theta} - \boldsymbol{\theta}_*) =: \mathbf{a}, \quad (15)$$

812 and substituting $\mathbf{G} \leftarrow \mathbf{G}\Sigma_{\mathbf{H}_{S_p}}^{1/2}$, where \mathbf{G} now refers to entries from unit variance isotropic Gaussian,
813 Eq. (14) is now

$$814 \quad \min_{\mathbf{a}} \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2 \quad \text{subject to} \quad \mathbf{G}\mathbf{a} = \mathbf{D}_w \mathbf{g}'.$$

815 Using Lagrangian formulation, the solution to the above problem is the same as the one to the
816 following unconstrained min-max problem.

$$817 \quad \min_{\mathbf{a}} \max_{\lambda} \lambda^\top \mathbf{G}\mathbf{a} - \lambda^\top \mathbf{D}_w \mathbf{g}' + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2.$$

818 Applying the CGMT [Thrapoulidis et al. \(2015; 2018\)](#), the corresponding Auxiliary Optimization
819 (AO) is

$$820 \quad \min_{\mathbf{a}} \max_{\lambda} \|\mathbf{a}\| \lambda^\top \mathbf{h} + \|\mathbf{a}\| \mathbf{a}^\top \mathbf{g} - \lambda^\top \mathbf{D}_w \mathbf{g}' + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2$$

821 Let $\beta = \|\mathbf{a}\|$, and $\hat{\boldsymbol{\lambda}}$, be the corresponding unit norm direction, then the above can be written as

$$822 \quad \min_{\mathbf{a}} \max_{\hat{\boldsymbol{\lambda}}, \beta \geq 0} \beta \hat{\boldsymbol{\lambda}}^\top (\|\mathbf{a}\| \mathbf{h} - \mathbf{D}_w \mathbf{g}') + \beta \mathbf{a}^\top \mathbf{g} + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2$$

$$823 \quad = \min_{\mathbf{a}} \max_{\beta \geq 0} \beta (\|\mathbf{a}\| \|\mathbf{h} - \mathbf{D}_w \mathbf{g}'\| + \beta \mathbf{a}^\top \mathbf{g} + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2),$$

824 where the final step uses the fact that the objective maximizes over $\hat{\boldsymbol{\lambda}}$ when $\hat{\boldsymbol{\lambda}} = \frac{\|\mathbf{a}\| \mathbf{h} - \mathbf{D}_w \mathbf{g}'}{\|\|\mathbf{a}\| \mathbf{h} - \mathbf{D}_w \mathbf{g}'\|}$.

825 The objective is convex in \mathbf{a} and linear, hence concave, in β . As shown in previous works, e.g.
826 ([Chang et al., 2021](#)), the constraint sets can be restricted to compact subsets, hence using Sion's
827 minimax theorem we can swap the order of the optimization to min-max to $\max_{\beta} \min_{\mathbf{a}}$:

$$828 \quad \max_{\beta \geq 0} \min_{\mathbf{a}} \beta (\|\mathbf{a}\| \|\mathbf{h} - \mathbf{D}_w \mathbf{g}'\| + \beta \mathbf{a}^\top \mathbf{g} + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2)$$

829 Next, using the variational characterization of the sqrt function (again, see for example ([Chang et al.,
830 2021](#), Sec. B.4)), the above optimization is

$$831 \quad \max_{\beta \geq 0} \min_{\mathbf{a}, u > 0} \frac{\beta u}{2} + \beta \frac{(\|\mathbf{a}\| \|\mathbf{h} - \mathbf{D}_w \mathbf{g}'\|)^2}{2u} + \beta \mathbf{a}^\top \mathbf{g} + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2 \quad (16)$$

832 Next, as $n \rightarrow \infty$, we have $\frac{\|\mathbf{h}\|^2}{n} \rightarrow 1$ and $\frac{\mathbf{h}^\top \mathbf{D}_w \mathbf{g}'}{n} \rightarrow 0$, and $\frac{\|\mathbf{D}_w \mathbf{g}'\|^2}{n} - \frac{1}{n} \sum_{i=1}^n d_i^2 \rightarrow 0$ in
833 probability. Here and onwards, we denote $d_i = \sigma_{n_i}$. Replacing these random quantities in Eq. (16)
834 yields the following asymptotically equivalent optimization wpa1:

$$835 \quad \max_{\beta \geq 0} \min_{\mathbf{a}, u > 0} \frac{\beta u}{2} + \beta \frac{n \|\mathbf{a}\|^2 + \sum_{i=1}^n d_i^2}{2u} + \beta \mathbf{a}^\top \mathbf{g} + \frac{1}{2} \|\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*\|^2 \quad (17)$$

836 Notice, we have a quadratic in \mathbf{a} . Minimizing the above w.r.t. \mathbf{a} , the optimal is at

$$837 \quad \mathbf{a}_* := \mathbf{a}_*(\beta, u) = \mathbf{M}^{-1}(\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \boldsymbol{\theta}_*), \quad \mathbf{M} := \frac{\beta n}{u} \mathbb{I} + \Sigma_{\mathbf{H}_{S_p}}^{-1}. \quad (18)$$

838 Plugging this back in Eq. (17), it simplifies to

$$839 \quad \max_{\beta \geq 0} \min_{u > 0} \underbrace{\frac{\beta u}{2} + \beta \frac{\sum_{i=1}^n d_i^2}{2u} - \frac{1}{2} (\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \boldsymbol{\theta}_*)^\top \mathbf{M}^{-1} (\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \boldsymbol{\theta}_*) + \frac{\|\boldsymbol{\theta}_*\|^2}{2}}_{R(u, \beta)} \quad (\text{AO})$$

840 D.1 SOLVING THE AUXILIARY OPTIMIZATION.

841 The subsequent analysis will use the saddle point $(\bar{u}_*, \bar{\beta}_*)$ of the deterministic limit objective
842 defined below. Formally relating $(\bar{u}_*, \bar{\beta}_*)$ to the finite-dimensional saddle point requires a uniform
843 convergence argument for the AO objective for which we refer the reader to [Chang et al. \(2021, Sec.
844 B.4\)](#). We treat this as a standard technical condition and proceed.

Choice of scaling. We next identify the natural normalization under the asymptotic regime. For our setup, note that $\|\boldsymbol{\theta}_*\| = \Theta(d^{3/2})$, $\lambda_i = \Theta(1/d^3)$ for all $i \in [p^2]$, where λ_i denote the Eigen values of $\boldsymbol{\Sigma}_{\mathbf{H}_{S_p}}$. Moreover, $\bar{\lambda}_i = \Theta(d^3)$, where $\bar{\lambda}_i$ denote the Eigen values of $\bar{\mathbf{M}}$. Using this, and Eqs. (1), (10), (12), (13) and (18), we define the following rescaled quantities:

$$\bar{u} := u/d, \quad \bar{\beta} := \beta/d^2, \quad (19)$$

$$\bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}} := d^3 \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}} = \frac{1}{\tau} \left(\sigma_n^2 + \frac{1}{M} \sum_{i=1}^M \frac{\|\mathbf{w}_i\|^2}{d} \right) \mathbb{I}_{p^2} + \frac{1}{M} \sum_{i=1}^M (\mathbf{w}_i[\mathcal{S}]\mathbf{w}_i[\mathcal{S}]^\top) \otimes \mathbb{I}_p, \quad (20)$$

$$\bar{\mathbf{M}} := \frac{M}{d^3} = \frac{\bar{\beta}\nu}{\bar{u}} \mathbb{I} + \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1} \quad (21)$$

$$\bar{\boldsymbol{\theta}}_* := \frac{\boldsymbol{\theta}_*}{d} = \frac{1}{dn} \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}}^{-1} \boldsymbol{\sigma}_{\mathbf{H}_{S_p}, \mathbf{y}} = \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1} \left(\frac{1}{M} \sum_{i=1}^M \mathbf{w}_i[\mathcal{S}] \otimes \mathbf{w}_i[\mathcal{S}] \right). \quad (22)$$

Next, we normalize the auxiliary objective (AO) by d^3 , and get the following objective

$$\begin{aligned} D_d(u, \beta) &:= \frac{1}{d^3} R(u, \beta) \\ &= \bar{\beta}\bar{u} + \frac{\bar{\beta}}{2\bar{u}} \frac{1}{d^2} \sum_{i=1}^n d_i^2 - \frac{1}{2d^2} \bar{\beta}^2 \mathbf{g}^\top \bar{\mathbf{M}}^{-1} \mathbf{g} - \frac{1}{2d} \bar{\boldsymbol{\theta}}_*^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\mathbf{M}}^{-1} \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_* + \frac{1}{d^{7/2}} \mathbf{g}^\top \bar{\mathbf{M}}^{-1} \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_* + \frac{\|\bar{\boldsymbol{\theta}}_*\|^2}{2d}. \end{aligned}$$

Using Lemma 14, it follows that

$$\lim_{d \rightarrow \infty} \frac{1}{2d^2} \mathbf{g}^\top \bar{\mathbf{M}}^{-1} \mathbf{g} = \lim_{d \rightarrow \infty} \frac{1}{d^2} \text{tr}(\bar{\mathbf{M}}^{-1}) =: s_\infty(\bar{u}, \bar{\beta}), \quad \lim_{d \rightarrow \infty} \frac{1}{d^{7/2}} \mathbf{g}^\top \bar{\mathbf{M}}^{-1} \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_* = 0. \quad (23)$$

Further, define the deterministic limits

$$c_\infty := \lim_{d \rightarrow \infty} \frac{1}{d^2} \sum_{i=1}^n d_i^2, \quad v_\infty(\bar{u}, \bar{\beta}) := \frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\mathbf{M}}^{-1} \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*. \quad (24)$$

Using Eq. (23), Eq. (24) we have the following deterministic limit of $D_d(\tau, \beta)$

$$D(\bar{u}, \bar{\beta}) = \frac{\bar{\beta}\bar{u}}{2} + \frac{\bar{\beta}}{2\bar{u}} c_\infty - \frac{\bar{\beta}^2}{2} s_\infty(\bar{u}, \bar{\beta}) - \frac{1}{2} v_\infty(\bar{u}, \bar{\beta}). \quad (25)$$

Here, we first compute $\lim_{d \rightarrow \infty} D_d(u, \beta)$, and then, since $\lim_{d \rightarrow \infty} \frac{\|\bar{\boldsymbol{\theta}}_*\|^2}{2d} = 0.5$, we omit it as it is independent of $\bar{\beta}, \bar{u}$.

D.2 SIMPLIFYING PER-TASK LOSS

Lemma 1. *The per-task loss evaluated at fixed \mathbf{w} with respect to the error vector \mathbf{a} is given by.*

$$\begin{aligned} \mathcal{L}(\mathbf{a}; \mathbf{w}) &= \underbrace{\mathbf{a}^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\Sigma}}(\mathbf{w}) \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{a}}_{\text{Term-1(a)}} + \underbrace{\frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top \bar{\boldsymbol{\Sigma}}(\mathbf{w}) \bar{\boldsymbol{\theta}}_*}_{\text{Term-1(b)}} + 2 \underbrace{\frac{1}{\sqrt{d}} \mathbf{a}^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\Sigma}}(\mathbf{w}) \bar{\boldsymbol{\theta}}_*}_{\text{Term-1(c)}} \\ &\quad - 2 \underbrace{\left(\frac{1}{d^{1/2}} \mathbf{a}^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top (\mathbf{w} \otimes \mathbf{w}) \right)}_{\text{Term-2}} + \frac{\|\mathbf{w}\|^2}{d} + \sigma_n^2. \end{aligned}$$

Proof.

$$\begin{aligned} \mathcal{L}(\hat{\boldsymbol{\theta}}_{S_p}; \mathbf{w}) &= \mathbb{E}_{X, \epsilon} [\hat{\boldsymbol{\theta}}_{S_p}^\top \mathbf{h}_{S_p}(X) - y_q]^2 \\ &= \hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X, \epsilon} [\mathbf{h}_{S_p}(X) \mathbf{h}_{S_p}(X)^\top] \hat{\boldsymbol{\theta}}_{S_p} + \mathbb{E}_{X, \epsilon} [y_q^2] - 2 \hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X, y} [\mathbf{h}_{S_p}(X) y_q] \\ &= \underbrace{\hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X, \epsilon} [\mathbf{h}_{S_p}(X) \mathbf{h}_{S_p}(X)^\top] \hat{\boldsymbol{\theta}}_{S_p}}_{\text{Term-1}} + \sigma_x^2 \|\mathbf{w}\|^2 + \sigma_n^2 - 2 \underbrace{\hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X, y} [\mathbf{h}_{S_p}(X) y_q]}_{\text{Term-2}}. \quad (26) \end{aligned}$$

Let $\Sigma(\mathbf{w}) := \mathbb{E}_{X,\epsilon}[\mathbf{h}_{S_p}(X)\mathbf{h}_{S_p}(X)^\top]$. Then, using Lemma 16 and Eq. (1), define

$$\bar{\Sigma}(\mathbf{w}) := d^3 \Sigma(\mathbf{w}) = \frac{1}{\tau} \left(\frac{\|\mathbf{w}\|^2}{d} + \sigma_n^2 \right) \mathbb{I}_{p^2} + \left(1 + \frac{1}{\tau d} \right) (\mathbf{w}[\mathcal{S}]\mathbf{w}[\mathcal{S}]^\top) \otimes \mathbb{I}_p.$$

To simplify Term-1 and Term-2, we use $\hat{\boldsymbol{\theta}}_{S_p} = \Sigma_{H_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*$ from Eq. (15), and Lemma 16.

First, we have Term-1:

$$\begin{aligned} (\Sigma_{H_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*)^\top \Sigma(\mathbf{w}) (\Sigma_{H_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*) &= \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \Sigma(\mathbf{w}) \Sigma_{H_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*^\top \Sigma(\mathbf{w}) \boldsymbol{\theta}_* + 2 \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \Sigma(\mathbf{w}) \boldsymbol{\theta}_* \\ &= \underbrace{\mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \Sigma_{H_{S_p}}^{-1/2} \mathbf{a}}_{\text{Term-1(a)}} + \underbrace{\frac{1}{d} \boldsymbol{\theta}_*^\top \bar{\Sigma}(\mathbf{w}) \boldsymbol{\theta}_*}_{\text{Term-1(b)}} + 2 \underbrace{\frac{1}{\sqrt{d}} \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \boldsymbol{\theta}_*}_{\text{Term-1(c)}}. \end{aligned} \quad (27)$$

Next, we have Term-2:

$$\begin{aligned} \hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X,y}[\mathbf{h}_{S_p}(X)y_q] &= (\Sigma_{H_{S_p}}^{-1/2} \mathbf{a} + \boldsymbol{\theta}_*)^\top (\sigma_x^4 \mathbf{w}[\mathcal{S}] \otimes \mathbf{w}[\mathcal{S}]) \\ &= \sigma_x^4 \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \mathbf{w} \otimes \mathbf{w} + \sigma_x^4 \boldsymbol{\theta}_*^\top \mathbf{w} \otimes \mathbf{w} \\ &= \frac{1}{d^2} \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} \mathbf{w} \otimes \mathbf{w} + \frac{1}{d^2} \boldsymbol{\theta}_*^\top \mathbf{w} \otimes \mathbf{w} \\ &= \frac{1}{d^{1/2}} \mathbf{a}^\top \Sigma_{H_{S_p}}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \boldsymbol{\theta}_*^\top (\mathbf{w} \otimes \mathbf{w}). \end{aligned} \quad (28)$$

Substituting Eqs. (27) and (28) in Eq. (26) finishes the proof. \square

E PROOFS FOR MIN NORM AO

E.1 WISHART MATRIX ASYMPTOTICS

Definition 4. Let μ be a probability measure on \mathbb{R} . The Stieltjes transform of μ , denoted by $m_\mu(z)$, is a function defined for all complex numbers $z \in \mathbb{C} \setminus \text{supp}(\mu)$ by:

$$m_\mu(z) = \int_{\mathbb{R}} \frac{1}{x - z} d\mu(x).$$

Remark 1 (Stieltjes transform of the Marchenko-Pastur law). In the asymptotic regime where $p, M \rightarrow \infty$ with $p/M \rightarrow \gamma$, the limiting spectral distribution of the normalized Wishart matrix \mathbf{W}_p is the Marchenko-Pastur law μ_γ . Its Stieltjes transform $m_\gamma(z)$ is determined by the quadratic equation:

$$z\gamma m_\gamma(z)^2 + (z + \gamma - 1)m_\gamma(z) + 1 = 0.$$

For $z \in \mathbb{C} \setminus \text{supp}(\mu_\gamma)$:

$$m_\gamma(z) = \frac{-(z + \gamma - 1) - \sqrt{(z + \gamma - 1)^2 - 4z\gamma}}{2z\gamma}.$$

E.2 GENERAL NOTATION

We use the following notation throughout the appendix. Let \mathbf{G} be a $M \times p$ random matrix with independent and identically distributed entries $G_{ij} \sim \mathcal{N}(0, 1)$. Let \mathbf{g}_i denote row i of \mathbf{G} . Let $\mathbf{W}_p = \frac{1}{M} \mathbf{G}^\top \mathbf{G}$ denote the normalized Wishart matrix. Let $\lambda_1, \dots, \lambda_p$ denote its eigenvalues. Let $\mathbf{g} \sim \mathcal{N}(0, \mathbb{I}_{p^2})$ be independent of \mathbf{W}_p , and let $\mathbf{w} \sim \mathcal{N}(0, \mathbb{I}_d)$ be independent of $(\mathbf{g}, \mathbf{W}_p)$.

Let \otimes denote the Kronecker product, and $\text{tr}(\cdot)$ denote the trace. We state some of its properties below. For conformable matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$, vectors \mathbf{a}, \mathbf{b} ,

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1}, \quad (29)$$

$$\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A}) \text{tr}(\mathbf{B}), \quad (30)$$

$$(\mathbf{A} \otimes \mathbb{I}) \text{vec}(\mathbf{C}) = \text{vec}(\mathbf{AC}), \quad (31)$$

$$(\mathbf{A} \otimes \mathbb{I})(\mathbf{B} \otimes \mathbb{I}) = (\mathbf{AB}) \otimes \mathbb{I}, \quad (32)$$

$$\text{vec}(\mathbf{ab}^\top) \text{vec}(\mathbf{ab}^\top)^\top = \mathbf{bb}^\top \otimes \mathbf{aa}^\top. \quad (33)$$

Also,

$$\text{vec}(\mathbf{A})^\top \text{vec}(\mathbf{B}) = \text{tr}(\mathbf{A}^\top \mathbf{B}). \quad (34)$$

Let $\Pi_{p \times d}$ denote the permutation matrix with all entries 0 or 1, and each row has exactly one 1 and most one 1 per column. Let $a, b, c, b', c', \beta, \kappa \in \mathbb{R}$ be constants. Also, let $\mathcal{S} \subset \{1, \dots, d\}$ denote a subset of indices, with $|\mathcal{S}| = p$. Next, we define matrices

$$\begin{aligned} \mathbf{A} &= c\mathbb{I}_p + b\mathbf{W}_p, & \mathbf{B} &= a\mathbb{I}_p + \mathbf{A}^{-1} = a\mathbb{I}_p + (c\mathbb{I}_p + b\mathbf{W}_p)^{-1}, \\ \boldsymbol{\Sigma} &= \mathbf{A} \otimes \mathbb{I}_p = c\mathbb{I}_{p^2} + b\mathbf{W}_p \otimes \mathbb{I}_p, \\ \mathbf{M} &= \mathbf{B} \otimes \mathbb{I}_p = a\mathbb{I}_{p^2} + \boldsymbol{\Sigma}^{-1}, \\ \boldsymbol{\Sigma}(\mathbf{w}) &= \left(c' \frac{\|\mathbf{w}\|^2}{d} + b'\right) \mathbb{I}_{p^2} + (\mathbf{w}[\mathcal{S}]\mathbf{w}[\mathcal{S}]^\top) \otimes \mathbb{I}_p, \\ \boldsymbol{\Sigma}(\mathbf{w}_j) &= \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b'\right) \mathbb{I}_{p^2} + (\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p, \end{aligned}$$

where $\mathbf{w}_j = \mathbf{g}_j$, the j^{th} row of \mathbf{G} . Next, define vectors

$$\begin{aligned} \mathbf{v} &= \text{vec}(\mathbf{W}_p), & \boldsymbol{\theta} &= \boldsymbol{\Sigma}^{-1} \left(\frac{1}{M} \sum_{i=1}^M \mathbf{g}_i \otimes \mathbf{g}_i \right) = \boldsymbol{\Sigma}^{-1} \mathbf{v}, \\ \mathbf{u} &= \beta \mathbf{g} - d^{1/2} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta}, \end{aligned}$$

and scalars

$$z_a = -\frac{(1+ac)}{ab}, \quad z_c = -\frac{c}{b}.$$

We consider the asymptotic regime where $d, p, M \rightarrow \infty$ with ratios

$$\frac{p}{d} =: \rho, \quad \frac{M}{d} =: \mu,$$

where $\rho, \mu \in (0, \infty)$. The following lemmas are in the asymptotic regime $d, p, M \rightarrow \infty$ with fixed ratios, we will use shorthand $\lim_{d \rightarrow \infty}$ to denote $\lim_{d, p, M \rightarrow \infty}$.

E.3 ASYMPTOTIC LIMIT OF $S(\tau, \beta)$

Lemma 2. *The normalized trace of the $p^2 \times p^2$ random matrix \mathbf{M}^{-1} converges almost surely to a deterministic limit:*

$$\lim_{d \rightarrow \infty} \frac{1}{d^2} \text{tr}(\mathbf{M}^{-1}) = \frac{\rho^2}{a} \left(1 - \frac{1}{ab} m_\gamma(z_a) \right).$$

Proof. Using properties of the Kronecker product (Eqs. (29) and (30)), we have:

$$\text{tr}(\mathbf{M}^{-1}) = \text{tr}(\mathbf{B}^{-1} \otimes \mathbb{I}_p) = p \text{tr}(\mathbf{B}^{-1}) = p \text{tr} \left((a\mathbb{I}_p + (c\mathbb{I}_p + b\mathbf{W}_p)^{-1})^{-1} \right).$$

The normalized trace is:

$$\frac{1}{d^2} \text{tr}(\mathbf{M}^{-1}) = \frac{p\rho^2}{p^2} \sum_{i=1}^p \left(a + \frac{1}{c + b\lambda_i} \right)^{-1} = \frac{\rho^2}{p} \sum_{i=1}^p \frac{b\lambda_i + c}{ab\lambda_i + ac + 1}.$$

As $p, M \rightarrow \infty$, and the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ . The sum converges to the integral:

$$\mathcal{I} = \frac{\rho^2}{a} \int \frac{\lambda + c/b}{\lambda + (ac + 1)/ab} d\mu_\gamma(\lambda).$$

Using Lemma 10 (\mathcal{I}_1 with $z_1 = z_c, z_2 = z_a$) gives the final expression. \square

E.4 ASYMPTOTIC LIMIT OF THE QUADRATIC FORMS

Lemma 3. Let $Q_p := \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta}$. Almost surely,

$$\lim_{d \rightarrow \infty} \frac{1}{d} Q_p = \rho \left(\frac{z_c^2}{b^2} m'_\gamma(z_c) + \frac{a}{b} \left(z_c(z_c - 2z_a) m_\gamma(z_c) + z_a^2 m_\gamma(z_a) \right) \right),$$

$$\lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} = \frac{\rho}{b} \left(1 + z_c + z_c^2 m_\gamma(z_c) \right).$$

Proof. First, we simplify Q_d :

$$\begin{aligned} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta} &= \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}^{-1} \mathbf{v} \\ &= \mathbf{v}^\top \boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-3/2} \mathbf{v}. \end{aligned}$$

Since \mathbf{M} shares eigenvectors with $\boldsymbol{\Sigma}$, it commutes with $\boldsymbol{\Sigma}$, and we have that

$$\boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-3/2} = \boldsymbol{\Sigma}^{-3} \mathbf{M}^{-1} = (a\boldsymbol{\Sigma}^3 + \boldsymbol{\Sigma}^2)^{-1}.$$

Hence, $Q_p = \mathbf{v}^\top (a\boldsymbol{\Sigma}^3 + \boldsymbol{\Sigma}^2)^{-1} \mathbf{v}$. Using Eqs. (29) and (31), we have

$$(a\boldsymbol{\Sigma}^3 + \boldsymbol{\Sigma}^2)^{-1} \mathbf{v} = \text{vec} \left((a\mathbf{A}^3 + \mathbf{A}^2)^{-1} \mathbf{W}_d \right).$$

Using Eq. (34), we get

$$Q_p = \text{tr} \left(\mathbf{W}_p (a\mathbf{A}^3 + \mathbf{A}^2)^{-1} \mathbf{W}_p \right) = \text{tr} \left(\mathbf{W}_p^2 (a\mathbf{A}^3 + \mathbf{A}^2)^{-1} \right).$$

Since \mathbf{W}_p and \mathbf{A} share eigenvectors, with eigenvalues λ_i and $c + b\lambda_i$, respectively,

$$Q_p = \sum_{i=1}^p \frac{\lambda_i^2}{(c + b\lambda_i)^2 (a(c + b\lambda_i) + 1)}.$$

As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

$$\lim_{d \rightarrow \infty} \frac{1}{d} Q_p = \frac{\rho}{ab^3} \int \frac{\lambda^2}{(\lambda + c/b)^2 (\lambda + (ac + 1)/ab)} d\mu_\gamma(x).$$

Using Lemma 10 (\mathcal{I}_3 with $z_1 = z_c, z_2 = z_a$) gives the final expression for the first result.

For the second part, we have that

$$\boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} = \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} \mathbf{v} = \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v}.$$

Using Eqs. (29) and (31),

$$\boldsymbol{\Sigma}^{-1} \mathbf{v} = \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbf{W}_d) = \text{vec}(\mathbf{A}^{-1} \mathbf{W}_d).$$

Using Eq. (34),

$$\boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} = \text{vec}(\mathbf{W}_p)^\top \text{vec}(\mathbf{A}^{-1} \mathbf{W}_p) = \text{tr}(\mathbf{W}_d \mathbf{A}^{-1} \mathbf{W}_d) = \text{tr}(\mathbf{W}_d^2 \mathbf{A}^{-1}).$$

Since \mathbf{W}_p and \mathbf{A} share eigenvectors, with eigenvalues λ_i and $c + b\lambda_i$, respectively,

$$\boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} = \sum_{i=1}^p \frac{\lambda_i^2}{c + b\lambda_i}.$$

As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

$$\lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} = \frac{\rho}{b} \int \frac{\lambda^2}{\lambda + c/b} d\mu_\gamma(x).$$

Using Lemma 10 (\mathcal{I}_1 with $z_1 = 0, z_2 = z_c$), we get

$$\int \frac{\lambda}{\lambda + c/b} d\mu_\gamma(x) = 1 + z_c m_\gamma(z_c).$$

Then using Lemma 10 ((v), with $\alpha_0 = 1, \alpha_1 = z_1 = z_c, \alpha_2 = 0$), we get the final expression. \square

E.5 FIRST ORDER OPTIMALITY CONDITIONS OF $D(\bar{u}, \bar{\beta})$

Theorem 5. Consider $D(\bar{u}, \bar{\beta})$ as in Eq. (25). Let $a := \frac{\bar{\beta}\nu}{\bar{u}}$. Then, the solution of $\min_{\bar{u}} \max_{\bar{\beta}} D(\bar{u}, \bar{\beta})$ is given by

$$\begin{aligned} \bar{u} &= \frac{\bar{\beta}\nu}{a}, \quad a = \frac{\gamma\nu/\rho^2 - (c+1) + \sqrt{(\gamma\nu/\rho^2 - (c+1))^2 + 4\gamma c\nu/\rho^2}}{2\gamma c(1 - \nu/\rho^2)}, \quad \gamma = \mu^{-1}, \\ \bar{\beta}^2 &= \frac{a^2(\sigma_n^2 + 1 - \rho + \rho c) + \rho(1 - a^2 c^2)m_\gamma(z_a) - \rho a z_a^2 m'_\gamma(z_a)}{\rho^2 \left(1 - \frac{\nu}{\rho^2} - \frac{m'_\gamma(z_a)}{a^2}\right)}. \end{aligned} \quad (35)$$

Proof. We derive the first order optimality conditions of $D(\bar{u}, \bar{\beta})$. First, we have

$$\begin{aligned} \frac{\partial D(\bar{u}, \bar{\beta})}{\partial \bar{u}} &= 0 \\ \implies \frac{\bar{\beta}}{2} - \frac{\bar{\beta}c_\infty}{2\bar{u}^2} - \frac{\bar{\beta}^2}{2} \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} - \frac{1}{2} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} &= 0 \end{aligned} \quad (36)$$

Similarly, we have

$$\begin{aligned} \frac{\partial D(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} &= 0 \\ \implies \frac{\bar{u}}{2} + \frac{c_\infty}{2\bar{u}} - \bar{\beta} s_\infty(\bar{u}, \bar{\beta}) - \frac{\bar{\beta}^2}{2} \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} - \frac{1}{2} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} &= 0 \end{aligned} \quad (37)$$

We first simplify the limits v_∞, s_∞ before we get the partial derivatives. First recall that

$$\begin{aligned} \bar{\Sigma}_{\mathbf{H}_{S_p}} &:= d^3 \Sigma_{\mathbf{H}_{S_p}} = \frac{1}{\tau} \left(\sigma_n^2 + \frac{1}{M} \sum_{i=1}^M \frac{\|\mathbf{w}_i\|^2}{d} \right) \mathbb{I}_{p^2} + \frac{1}{M} \sum_{i=1}^M (\mathbf{w}_i[\mathbf{S}] \mathbf{w}_i[\mathbf{S}]^\top) \otimes \mathbb{I}_p, \\ \bar{\mathbf{M}} &:= \frac{\mathbf{M}}{d^3} = \frac{\bar{\beta}\nu}{\bar{u}} \mathbb{I} + \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1} \\ \bar{\boldsymbol{\theta}}_\star &:= \frac{\boldsymbol{\theta}_\star}{d} = \frac{1}{dn} \Sigma_{\mathbf{H}_{S_p}}^{-1} \boldsymbol{\sigma}_{\mathbf{H}_{S_p} \mathbf{y}} = \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1} \left(\frac{1}{M} \sum_{i=1}^M \mathbf{w}_i[\mathbf{S}] \otimes \mathbf{w}_i[\mathbf{S}] \right). \end{aligned}$$

Then, since $\lim_{d \rightarrow \infty} \frac{\|\mathbf{w}_i\|^2}{d} = 1$, using Lemma 2 with $\mathbf{M} = \bar{\mathbf{M}}$, $a = \frac{\bar{\beta}\nu}{\bar{u}}$, $b = 1$, $c = \frac{\sigma_n^2 + 1}{\tau}$, $\gamma = \mu^{-1}$, $z_a = -\frac{1}{a} - c$, we get

$$s_\infty(\bar{u}, \bar{\beta}) = \lim_{d \rightarrow \infty} d \operatorname{tr}(\bar{\mathbf{M}}^{-1}) = \lim_{d \rightarrow \infty} \frac{1}{d^2} \operatorname{tr}(\bar{\mathbf{M}}^{-1}) = \rho^2 \left(\frac{1}{a} - \frac{1}{a^2} m_\gamma(z_a) \right). \quad (38)$$

Using Lemma 3 with $\Sigma = \bar{\Sigma}_{\mathbf{H}_{S_p}}$, $\boldsymbol{\theta} = \bar{\boldsymbol{\theta}}_\star$, we get

$$\begin{aligned} v_\infty(\bar{u}, \bar{\beta}) &= \lim_{d \rightarrow \infty} \frac{1}{d} (\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_\star)^\top \bar{\mathbf{M}}^{-1} (\Sigma_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_\star) = \lim_{d \rightarrow \infty} \frac{1}{d} (\bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_\star)^\top \bar{\mathbf{M}}^{-1} (\bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_\star) \\ &= \rho(c^2 m'_\gamma(z_c) + ac(c + 2z_a) m_\gamma(z_c) + a z_a^2 m_\gamma(z_a)). \end{aligned}$$

As both $s_\infty(\bar{u}, \bar{\beta})$ and $v_\infty(\bar{u}, \bar{\beta})$ depend on $\bar{u}, \bar{\beta}$ through a , we use the derivatives with respect to a to find the partial derivatives wrt $\bar{u}, \bar{\beta}$.

$$\frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} = \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \frac{\partial a}{\partial \bar{u}} = \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \left(-\frac{\bar{\beta}\nu}{\bar{u}^2} \right).$$

$$\frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} = \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \frac{\partial a}{\partial \bar{\beta}} = \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \left(\frac{\nu}{\bar{u}} \right).$$

$$\begin{aligned} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} &= \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \frac{\partial a}{\partial \bar{u}} = \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \left(-\frac{\bar{\beta}\nu}{\bar{u}^2} \right). \\ \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} &= \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \frac{\partial a}{\partial \bar{\beta}} = \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \left(\frac{\nu}{\bar{u}} \right). \end{aligned}$$

Substituting back in Eq. (36)

$$\begin{aligned} \frac{\partial D(\bar{u}, \bar{\beta})}{\partial \bar{u}} &= \frac{\bar{\beta}}{2} - \frac{\bar{\beta}c_\infty}{2\bar{u}^2} - \frac{\bar{\beta}^2}{2} \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} - \frac{1}{2} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{u}} \\ &= \frac{\bar{\beta}}{2} - \frac{\bar{\beta}c_\infty}{2\bar{u}^2} - \frac{\bar{\beta}^2}{2} \left(\frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \right) \left(-\frac{\bar{\beta}\nu}{\bar{u}^2} \right) - \frac{1}{2} \left(\frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \right) \left(-\frac{\bar{\beta}\nu}{\bar{u}^2} \right) \\ &= \frac{\bar{\beta}}{2} - \frac{\bar{\beta}c_\infty}{2\bar{u}^2} + \frac{\bar{\beta}^3\nu}{2\bar{u}^2} \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} + \frac{\bar{\beta}\nu}{2\bar{u}^2} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} = 0. \end{aligned}$$

Simplifying it further we get this is

$$\begin{aligned} \bar{u}^2 - c_\infty + \bar{\beta}^2\nu \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} + \nu \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} &= 0 \\ \iff \bar{u}^2 = c_\infty - \nu \left(\bar{\beta}^2 \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} + \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \right). \end{aligned} \quad (39)$$

Similarly, substituting in Eq. (37)

$$\begin{aligned} \frac{\partial D(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} &= \frac{\bar{u}}{2} + \frac{c_\infty}{2\bar{u}} - \bar{\beta} s_\infty(\bar{u}, \bar{\beta}) - \frac{\bar{\beta}^2}{2} \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} - \frac{1}{2} \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial \bar{\beta}} \\ &= \frac{\bar{u}}{2} + \frac{c_\infty}{2\bar{u}} - \bar{\beta} s_\infty(\bar{u}, \bar{\beta}) - \frac{\bar{\beta}^2}{2} \left(\frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} \right) \left(\frac{\nu}{\bar{u}} \right) - \frac{1}{2} \left(\frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \right) \left(\frac{\nu}{\bar{u}} \right) = 0. \end{aligned}$$

Simplifying it further we have

$$\bar{u} + \frac{c_\infty}{\bar{u}} - 2\bar{\beta} s_\infty(\bar{u}, \bar{\beta}) - \frac{\nu}{\bar{u}} \left(\bar{\beta}^2 \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} + \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \right) = 0.$$

Using Eq. (39), we get

$$\begin{aligned} 0 &= \bar{u} + \frac{c_\infty}{\bar{u}} - 2\bar{\beta} s_\infty(\bar{u}, \bar{\beta}) - \frac{1}{\bar{u}} (c_\infty - \bar{u}^2) \\ &= 2\bar{u} - 2\bar{\beta} s_\infty(\bar{u}, \bar{\beta}) \\ \iff \bar{u} &= \bar{\beta} s_\infty(\bar{u}, \bar{\beta}). \end{aligned} \quad (40)$$

Using Eq. (38) and simplifying Eq. (40), we get

$$\begin{aligned} \bar{u} &= \frac{\bar{\beta}\rho^2}{a} \left(1 - \frac{1}{a} m_\gamma(z_a) \right) \\ \iff \frac{\rho^2}{\nu} \left(1 - \frac{1}{a} m_\gamma(z_a) \right) &= 1, \\ \iff \left(1 - \frac{\nu}{\rho^2} \right) a &= m_\gamma(z_a). \end{aligned} \quad (41)$$

Eq. (41) can be simplified as follows. MP quadratic for $m_\gamma(z)$:

$$z\gamma m_\gamma(z)^2 + (z + \gamma - 1)m_\gamma(z) + 1 = 0.$$

Substituting $z = z_a = -1/a - c$, $m_\gamma(z_a) = a(1 - \nu)$, we have

$$\begin{aligned} \left(-\frac{1}{a} - c \right) \gamma a^2 (1 - \nu/\rho^2)^2 + \left(-\frac{1}{a} - c + \gamma - 1 \right) a(1 - \nu/\rho^2) + 1 &= 0 \\ \iff - (1 + ac)\gamma a(1 - \nu/\rho^2)^2 - (1 + ac)(1 - \nu/\rho^2) + (\gamma - 1)a(1 - \nu/\rho^2) + 1 &= 0 \\ \iff -\gamma c(1 - \nu/\rho^2)^2 a^2 + (1 - \nu/\rho^2)(\gamma\nu/\rho^2 - (c + 1))a + \nu/\rho^2 &= 0 \\ \iff c\gamma(1 - \nu/\rho^2)a^2 + ((c + 1) - \gamma\nu/\rho^2)a - \frac{\nu/\rho^2}{1 - \nu/\rho^2} &= 0. \end{aligned}$$

Solving for a using the quadratic formula finishes the proof for the first part.

Define

$$q_\infty(\bar{u}, \bar{\beta}) := - \left(\bar{\beta}^2 \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} + \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} \right).$$

Next to get a simplified form of this, we use $\frac{dz}{da} = \frac{1}{a^2}$ and get the partials

$$\begin{aligned} s_\infty(\bar{u}, \bar{\beta}) &= \rho^2 \left(\frac{1}{a} - \frac{1}{a^2} m_\gamma(z_a) \right), \\ \frac{\partial s_\infty(\bar{u}, \bar{\beta})}{\partial a} &= \rho^2 \left(-\frac{1}{a^2} + \frac{2}{a^3} m_\gamma(z_a) - \frac{1}{a^4} m'_\gamma(z_a) \right). \end{aligned}$$

$$\begin{aligned} v_\infty(\bar{u}, \bar{\beta}) &= \rho(c^2 m'_\gamma(-c) + ac(c + 2z_a)m_\gamma(-c) + a(z_a)^2 m_\gamma(z_a)) \\ \frac{\partial v_\infty(\bar{u}, \bar{\beta})}{\partial a} &= \rho(cm_\gamma(-c) \left(c + 2z_a + \frac{2}{a} \right) + z_a m_\gamma(z_a) \left(z_a + \frac{2}{a} \right) + \frac{z_a^2}{a} m'_\gamma(z_a)) \\ &= \rho(-c^2 m_\gamma(-c) + z_a m_\gamma(z_a) \left(\frac{1}{a} - c \right) + \frac{z_a^2}{a} m'_\gamma(z_a)). \end{aligned}$$

$$\begin{aligned} q_\infty(\bar{u}, \bar{\beta}) &= - \left[\bar{\beta}^2 \rho^2 \left(-\frac{1}{a^2} + \frac{2}{a^3} m_\gamma(z_a) - \frac{1}{a^4} m'_\gamma(z_a) \right) \right. \\ &\quad \left. + \rho \left(-c^2 m_\gamma(-c) + z_a m_\gamma(z_a) \left(\frac{1}{a} - c \right) + \frac{z_a^2}{a} m'_\gamma(z_a) \right) \right] \\ &= \frac{\bar{\beta}^2 \rho^2}{a^2} \left(1 - \frac{2}{a} m_\gamma(z_a) + \frac{1}{a^2} m'_\gamma(z_a) \right) + \rho \left(c^2 m_\gamma(z_c) + \left(\frac{1}{a^2} - c^2 \right) m_\gamma(z_a) - \frac{z_a^2}{a} m'_\gamma(z_a) \right). \end{aligned} \quad (42)$$

Next, we simplify c_∞ . Using Eq. (11) and Lemma 3 with $c = \frac{\sigma_n^2 + 1}{\tau}$, $b = 1$, we have

$$\begin{aligned} c_\infty &= \lim_{d \rightarrow \infty} \frac{1}{d^2} \left(\sum_{i=1}^n d_i^2 \right) = \lim_{d \rightarrow \infty} \frac{1}{d^2} \left(\sum_{i=1}^n \sigma_x^2 \|\mathbf{w}_i\|^2 + n\sigma_n^2 - n\boldsymbol{\theta}_*^\top \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_* \right) \\ &= \lim_{d \rightarrow \infty} \frac{1}{d^2} \left(\sum_{i=1}^M \frac{n}{M} \sigma_x^2 \|\mathbf{w}_i\|^2 \right) - \nu \lim_{d \rightarrow \infty} \boldsymbol{\theta}_*^\top \boldsymbol{\Sigma}_{\mathbf{H}_{S_p}} \boldsymbol{\theta}_* + \nu \sigma_n^2 \\ &= \nu \lim_{d \rightarrow \infty} \frac{1}{Md} \left(\sum_{i=1}^M \|\mathbf{w}_i\|^2 \right) - \nu \lim_{d \rightarrow \infty} \frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{H}_{S_p}} \bar{\boldsymbol{\theta}}_* + \nu \sigma_n^2 \\ &= \nu(1 + \sigma_n^2) - \nu\rho(1 - c + c^2 m_\gamma(-c)) = \nu(\sigma_n^2 + 1 - \rho + \rho c - \rho c^2 m_\gamma(-c)). \end{aligned} \quad (43)$$

We can write Eq. (39) as

$$\bar{u}^2 = c_\infty + \nu q_\infty(\bar{u}, \bar{\beta}),$$

Substituting c_∞ and q_∞ using Eqs. (42) and (43), and using $a = \frac{\bar{\beta}\nu}{\bar{u}}$, we get

$$\bar{\beta}_*^2 \nu = a^2(\sigma_n^2 + 1 - \rho + \rho c) + \bar{\beta}^2 \rho^2 \left(1 - \frac{2}{a} m_\gamma(z_a) + \frac{1}{a^2} m'_\gamma(z_a) \right) + \rho(1 - a^2 c^2) m_\gamma(z_a) - \rho a z_a^2 m'_\gamma(z_a).$$

Solving for $\bar{\beta}^2$ gives Eq. (35).

□

1242 F PROOF FOR LS AO

1243
1244 The PO is

$$1245 \min_{\mathbf{a}} \frac{1}{2} \|\mathbf{G}\mathbf{a} - \mathbf{D}_w \mathbf{g}'\|^2 = \min_{\mathbf{a}} \max_{\mathbf{u}} \mathbf{u}^\top \mathbf{G}\mathbf{a} - \mathbf{u}^\top \mathbf{D}_w \mathbf{g}' - \frac{1}{2} \|\mathbf{u}\|^2.$$

1246
1247 By direct application of the CGMT, the AO is

$$1248 \min_{\mathbf{a}} \max_{\mathbf{u}} \|\mathbf{u}\| \mathbf{g}^\top \mathbf{a} + \|\mathbf{a}\| \mathbf{h}^\top \mathbf{u} - \mathbf{u}^\top \mathbf{D}_w \mathbf{g}' - \frac{1}{2} \|\mathbf{u}\|^2.$$

1249
1250 Let $\alpha = \|\mathbf{a}\|$, $\beta = \|\mathbf{u}\|$, and $\hat{\mathbf{a}}$ and $\hat{\mathbf{u}}$ be the corresponding unit norm directions. Then the above can
1251 be written as

$$1252 \min_{\alpha \geq 0, \hat{\mathbf{a}}} \max_{\beta \geq 0, \hat{\mathbf{u}}} \beta \alpha \mathbf{g}^\top \hat{\mathbf{a}} + \beta \hat{\mathbf{u}}^\top (\alpha \mathbf{h} - \mathbf{D}_w \mathbf{g}') - \frac{1}{2} \beta^2.$$

1253
1254 Solving for $\hat{\mathbf{a}}$ and $\hat{\mathbf{u}}$, the optimal is at

$$1255 \hat{\mathbf{a}} = -\frac{\mathbf{g}}{\|\mathbf{g}\|}, \quad \hat{\mathbf{u}} = \frac{\alpha \mathbf{h} - \mathbf{D}_w \mathbf{g}'}{\|\alpha \mathbf{h} - \mathbf{D}_w \mathbf{g}'\|}.$$

1256
1257 Plugging these back, we have

$$1258 \min_{\alpha \geq 0} \max_{\beta \geq 0} \beta \|\alpha \mathbf{h} - \mathbf{D}_w \mathbf{g}'\| - \beta \alpha \|\mathbf{g}\| - \frac{1}{2} \beta^2$$

$$1259 = \min_{\alpha \geq 0} \max_{\beta \geq 0} \beta \left(\sqrt{n\alpha^2 + \sum_{i \in [n]} d_i^2} - \alpha \|\mathbf{g}\| \right) - \frac{1}{2} \beta^2$$

$$1260 = \frac{1}{2} \left(\min_{\alpha \geq 0} \sqrt{n\alpha^2 + \sum_{i=1}^n d_i^2} - \alpha \|\mathbf{g}\| \right)_+^2$$

1261
1262 where $(\cdot)_+ = \max\{x, 0\}$. Thus, using Eq. (1),

$$1263 \alpha_*^2 = \frac{\sum_{i=1}^n d_i^2}{n \left(\frac{n}{\|\mathbf{g}\|^2} - 1 \right)} = \frac{\sum_{i=1}^n d_i^2}{\nu d^2 \left(\frac{\nu p^2}{\rho^2 \|\mathbf{g}\|^2} - 1 \right)}.$$

1264
1265 Then, we get $\mathbf{a}_* = -\alpha_* \hat{\mathbf{g}}$. Using Eq. (24),

$$1266 \kappa_\infty := \lim_{d \rightarrow \infty} \alpha_* = \sqrt{\frac{c_\infty}{\nu \left(\frac{\nu}{\rho^2} - 1 \right)}}.$$

1267
1268 In the limit, $\mathbf{a}_* = \kappa_\infty \frac{\mathbf{g}}{\|\mathbf{g}\|}$.

1269 G ASYMPTOTIC RISK

1270 G.1 ASYMPTOTIC RISK OF MINIMUM-NORM LGP

1271
1272 **Lemma 4** (Term-1, OOD Loss). *Using \mathbf{M} , Σ , $\Sigma(\mathbf{w})$, \mathbf{u} , $\boldsymbol{\theta}$ of the same form as Appendix E.2, let*

$$1273 T_1(\mathbf{w}) := \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \Sigma(\mathbf{w}) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u},$$

$$1274 T_2(\mathbf{w}) := \frac{1}{d} \boldsymbol{\theta}^\top \Sigma(\mathbf{w}) \boldsymbol{\theta}, \quad T_3(\mathbf{w}) := \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \Sigma(\mathbf{w}) \boldsymbol{\theta}.$$

1275
1276 Then,

$$1277 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}} [T_1(\mathbf{w})] = (1 + b' + c') J_1,$$

$$1278 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}} [T_2(\mathbf{w})] = \frac{(1 + b' + c') \rho}{b^2} (1 + 2z_c m_\gamma(z_c) + z_c^2 m'_\gamma(z_c)), \quad (44)$$

$$1279 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}} [T_3(\mathbf{w})] = -(1 + b' + c') \rho \left(-\frac{c(ac + 2)}{b^3} m_\gamma(z_c) + \frac{a}{b} z_a^2 m_\gamma(z_a) + \frac{1}{b^2} z_c^2 m'_\gamma(z_c) \right), \quad (45)$$

1296

where

1297

1298

1299

1300

1301

$$J_1 = \frac{\rho}{b^2} \left[\frac{\beta^2 \rho}{a^2} \left(b m_\gamma(z_a) - \frac{1}{a} m'_\gamma(z_a) \right) - 2 a b z_c z_a (m_\gamma(z_c) - m_\gamma(z_a)) + z_c^2 m'_\gamma(z_c) + z_a^2 m'_\gamma(z_a) \right]. \quad (46)$$

1302

Proof. We first work with $T_1(\mathbf{w})$. Using Lemma 11, since \mathbf{w} is independent of \mathbf{g} , \mathbf{W}_p , we get

1303

1304

1305

$$\mathbb{E}_{\mathbf{w}}[T_1(\mathbf{w})] = \frac{c' + b' + 1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-1} \mathbf{u}.$$

1306

Since \mathbf{M} and $\boldsymbol{\Sigma}$ share eigenvectors, they commute, and we have that $\mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-1} = \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}$.

1307

Further, we have that

1308

1309

1310

1311

1312

1313

$$\begin{aligned} \mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2} \mathbf{u} &= \beta^2 \mathbf{g}^\top (\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}) \mathbf{g} + d \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta} \\ &\quad - \beta d^{1/2} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2} \mathbf{g} - \beta d^{1/2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta} \\ &= \beta^2 \mathbf{g}^\top (\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}) \mathbf{g} + d \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-2} \mathbf{M}^{-2} \boldsymbol{\theta} - 2 \beta d^{1/2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-2} \boldsymbol{\theta}. \end{aligned}$$

1314

1315

1316

1317

1318

1319

$$\begin{aligned} J_1 &= \lim_{d \rightarrow \infty} \frac{1}{d^2} \mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2} \mathbf{u} \\ &= \lim_{p \rightarrow \infty} \frac{\rho^2}{p^2} \beta^2 \mathbf{g}^\top (\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}) \mathbf{g} + \lim_{p \rightarrow \infty} \frac{\rho}{p} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-2} \mathbf{M}^{-2} \boldsymbol{\theta} - \lim_{p \rightarrow \infty} \frac{\rho^{3/2}}{p^{3/2}} 2 \beta \mathbf{g}^\top \boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-2} \boldsymbol{\theta} \\ &= \lim_{p \rightarrow \infty} \beta^2 \mathcal{I}'_1(p) + \lim_{p \rightarrow \infty} \mathcal{I}'_2(p), \end{aligned} \quad (47)$$

1320

where

1321

1322

$$\mathcal{I}'_1(p) := \frac{1}{d^2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}), \quad \mathcal{I}'_2(p) := \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-2} \mathbf{M}^{-2} \boldsymbol{\theta},$$

1323

For the first term as $\mathbf{g} \in \mathbb{R}^{p^2}$ we apply Lemma 14 with $d_1 = p^2$ and get

1324

1325

1326

$$\lim_{p \rightarrow \infty} \frac{1}{p^2} \mathbf{g}^\top (\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}) \mathbf{g} = \lim_{p \rightarrow \infty} \frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{M}^{-2}).$$

1327

Also, for the third term, from Lemma 14 it follows that

1328

1329

1330

$$\lim_{p \rightarrow \infty} \frac{1}{p^{3/2}} \mathbf{g}^\top \boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-2} \boldsymbol{\theta} = 0, \quad (48)$$

1331

as $\|\boldsymbol{\Sigma}^{-3/2} \mathbf{M}^{-2} \boldsymbol{\theta}\| = O(\|\boldsymbol{\theta}\|) = O(\text{vec}(\mathbf{W}_p)) = O(\sqrt{p})$.

1332

1333

For $\mathcal{I}'_1(p)$, using Eqs. (29), (30) and (32) we have that

1334

1335

1336

1337

1338

1339

1340

1341

$$\begin{aligned} \mathcal{I}'_1(p) &= \frac{\rho^2}{p^2} \text{tr}((\mathbf{A}^{-1} \mathbf{B}^{-2}) \otimes \mathbb{I}_p) \\ &= \frac{\rho^2}{p} \text{tr}((c \mathbb{I}_p + b \mathbf{W}_p)^{-1} (a \mathbb{I}_p + (c \mathbb{I}_p + b \mathbf{W}_p)^{-1})^{-2}) \\ &= \frac{\rho^2}{p} \sum_{i=1}^p \frac{1}{c + b \lambda_i} \frac{1}{(a + \frac{1}{c + b \lambda_i})^2} = \frac{\rho^2}{p} \sum_{i=1}^p \frac{c + b \lambda_i}{(1 + a(c + b \lambda_i))^2}. \end{aligned}$$

1342

As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

1344

1345

1346

$$\lim_{p \rightarrow \infty} \mathcal{I}'_1(p) = \frac{\rho^2}{a^2 b} \int \frac{\lambda + c/b}{(\lambda + (1 + ac)/ab)^2} d\mu_{\text{MP}, \gamma}(\lambda).$$

1347

Using Lemma 10 (\mathcal{I}_2 with $z_1 = z_c, z_2 = z_a$) gives

1348

1349

$$\lim_{p \rightarrow \infty} \mathcal{I}'_1(p) = \frac{\rho^2}{a^2 b} (m_\gamma(z_a) - \frac{1}{ab} m'_\gamma(z_a)). \quad (49)$$

For $\mathcal{I}'_2(p)$, since \mathbf{M} and Σ share eigenvectors, and using Eqs. (29), (31), (32) and (34), we have

$$\begin{aligned} \boldsymbol{\theta}^\top \Sigma^{-2} \mathbf{M}^{-2} \boldsymbol{\theta} &= \mathbf{v}^\top \Sigma^{-1} \Sigma^{-2} \mathbf{M}^{-2} \Sigma^{-1} \mathbf{v} = \mathbf{v}^\top \Sigma^{-4} \mathbf{M}^{-2} \mathbf{v} \\ &= \mathbf{v}^\top (\mathbf{A}^{-4} \mathbf{B}^{-2} \otimes \mathbb{I}) \text{vec}(\mathbf{W}_p) = \text{vec}(\mathbf{W}_p)^\top \text{vec}(\mathbf{A}^{-4} \mathbf{B}^{-2} \mathbf{W}_p) \\ &= \text{tr}(\mathbf{W}_p^2 \mathbf{A}^{-4} \mathbf{B}^{-2}). \end{aligned}$$

Using this, we get

$$\mathcal{I}'_2(p) = \frac{\rho}{p} \sum_{i=1}^p \frac{\lambda_i^2}{(c + b\lambda_i)^4} \frac{1}{(a + (c + b\lambda_i)^{-1})^2} = \frac{\rho}{p} \sum_{i=1}^p \frac{\lambda_i^2}{(c + b\lambda_i)^2 (a(c + b\lambda_i) + 1)^2}.$$

As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

$$\lim_{p \rightarrow \infty} \mathcal{I}'_2(p) = \frac{\rho}{a^2 b^4} \int \frac{\lambda^2}{(\lambda + c/b)^2 (\lambda + (1 + ac)/ab)^2} d\mu_{\text{MP}, \gamma}(\lambda)$$

Using Lemma 10 (\mathcal{I}_7 with $z_1 = z_a, z_2 = z_c$), we get

$$\lim_{d \rightarrow \infty} \mathcal{I}'_2(p) = \frac{\rho}{a^2 b^4} \left[\frac{2z_a z_c}{(z_a - z_c)^3} (m_\gamma(z_c) - m_\gamma(z_a)) + \frac{z_c^2}{(z_a - z_c)^2} m'_\gamma(z_c) + \frac{z_a^2}{(z_a - z_c)^2} m'_\gamma(z_a) \right]. \quad (50)$$

Using $z_a - z_c = -1/(ab)$ to simplify Eq. (50), and substituting it with Eq. (49) in Eq. (47) gives the final expression in Eq. (46).

Next, working with $T_2(\mathbf{w})$, using Lemma 11, we have

$$\mathbb{E}_{\mathbf{w}} \left[\frac{1}{d} \boldsymbol{\theta}^\top \Sigma(\mathbf{w}) \boldsymbol{\theta} \right] = \frac{(1 + b' + c')}{d} \boldsymbol{\theta}^\top \boldsymbol{\theta} = \frac{(1 + b' + c')}{d} \mathbf{v}^\top \Sigma^{-2} \mathbf{v}.$$

Using Eqs. (29), (31) and (34), we have

$$\mathbf{v}^\top \Sigma^{-2} \mathbf{v} = \text{vec}(\mathbf{W}_p)^\top \text{vec}(\mathbf{A}^{-2} \mathbf{W}_p) = \text{tr}(\mathbf{W}_p^2 (c\mathbb{I}_p + b\mathbf{W}_p)^{-2}). \quad (51)$$

Therefore,

$$\mathbb{E}_{\mathbf{w}} \left[\frac{1}{d} \boldsymbol{\theta}^\top \Sigma(\mathbf{w}) \boldsymbol{\theta} \right] = \frac{(1 + b' + c') \rho}{p} \sum_{i=1}^p \frac{\lambda_i^2}{(c + b\lambda_i)^2}.$$

As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

$$\lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}} \left[\frac{1}{d} \boldsymbol{\theta}^\top \Sigma(\mathbf{w}) \boldsymbol{\theta} \right] = \frac{(1 + b' + c') \rho}{b^2} \int \frac{\lambda^2}{(\lambda + c/b)^2} d\mu_\gamma(\lambda). \quad (52)$$

Using Lemma 10 (\mathcal{I}_6 with $z_1 = z_c$) then gives the expression in Eq. (44).

Next, working with $T_3(\mathbf{w})$, using Lemma 11, we have that

$$\begin{aligned} \mathbb{E}_{\mathbf{w}}[T_3(\mathbf{w})] &= \frac{1 + b' + c'}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} \\ &= \frac{1 + b' + c'}{d^{3/2}} (\beta \mathbf{g} - d^{1/2} \Sigma^{-1/2} \boldsymbol{\theta})^\top \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} \\ &= \frac{1 + b' + c'}{d} \left(\beta d^{-1/2} \mathbf{g}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} - \boldsymbol{\theta}^\top \Sigma^{-1/2} \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} \right). \end{aligned}$$

Using Lemma 14, we have that

$$\lim_{d \rightarrow \infty} \frac{1}{d^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} = \lim_{p \rightarrow \infty} \frac{\rho^{3/2}}{p^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \boldsymbol{\theta} = 0,$$

1404 as $\|\mathbf{M}^{-1}\boldsymbol{\Sigma}^{-1/2}\boldsymbol{\theta}\| = O(\|\boldsymbol{\theta}\|) = O(\|\boldsymbol{\Sigma}^{-1}\text{vec}(\mathbf{W}_p)\|) = O(\sqrt{p})$.

1405
1406 Using this we have

$$1407 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}}[T_3(\mathbf{w})] = -(1 + b' + c') \lim_{d \rightarrow \infty} \mathcal{I}'_3(p), \quad \mathcal{I}'_3(p) = \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta}. \quad (53)$$

1408 Since $\boldsymbol{\Sigma}$ and \mathbf{M} share eigenvectors, and using Eqs. (29), (31), (32) and (34), we have

$$1409 \begin{aligned} 1410 \mathcal{I}'_3(p) &= \frac{1}{d} \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}^{-1} \mathbf{v} = \frac{1}{d} \mathbf{v}^\top \boldsymbol{\Sigma}^{-3} \mathbf{M}^{-1} \mathbf{v} \\ 1411 &= \frac{1}{d} \text{vec}(\mathbf{W}_p)^\top \text{vec}(\mathbf{A}^{-3} \mathbf{B}^{-1} \mathbf{W}_p) = \frac{1}{d} \text{tr}(\mathbf{W}_p^2 \mathbf{A}^{-3} \mathbf{B}^{-1}) \\ 1412 &= \frac{1}{d} \text{tr}(\mathbf{W}_p^2 (c\mathbb{I} + b\mathbf{W}_p)^{-3} (a\mathbb{I} + (c\mathbb{I} + b\mathbf{W}_p)^{-1})^{-1}) \\ 1413 &= \frac{\rho}{p} \sum_{i=1}^p \frac{\lambda_i^2}{(c + b\lambda_i)^2 (a(c + b\lambda_i) + 1)}. \end{aligned} \quad (54)$$

1414
1415 As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding,

$$1416 \lim_{d \rightarrow \infty} \mathcal{I}'_3(p) = \frac{\rho}{ab^2} \int \frac{\lambda^2}{(\lambda + c/b)^2 (\lambda + (1 + ac)/ab)} d\mu_\gamma(\lambda). \quad (55)$$

1417 Using Lemma 10 (\mathcal{I}_3 with $z_1 = z_c, z_2 = z_a$) and simplifying gives the final expression in Eq. (45). \square

1418
1419 **Lemma 5** (Term-2, OOD Loss). *Using $\mathbf{M}, \boldsymbol{\Sigma}, \mathbf{u}, \boldsymbol{\theta}$ of the same form as Appendix E.2, consider the scalar random variable*

$$1420 U_p(\mathbf{w}) := \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \boldsymbol{\theta}^\top (\mathbf{w} \otimes \mathbf{w}). \quad (56)$$

1421 Then, almost surely,

$$1422 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}}[U_p(\mathbf{w})] = \frac{\rho}{b} \left(1 + z_a m_\gamma(z_a)\right).$$

1423 *Proof.* Using $\mathbb{E}_{\mathbf{w}}[\mathbf{w} \otimes \mathbf{w}] = \text{vec}(\mathbb{I}_p)$ and since \mathbf{w} is independent of \mathbf{g} and \mathbf{W}_p , we have that

$$1424 \mathbb{E}_{\mathbf{w}}[U_p(\mathbf{w})] = \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) + \frac{1}{d} \boldsymbol{\theta}^\top \text{vec}(\mathbb{I}_p).$$

1425 Considering the first term, we have

$$1426 \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) = \frac{\beta}{d^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) - \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p).$$

1427 Using Lemma 14, we have that

$$1428 \lim_{d \rightarrow \infty} \frac{\beta}{d^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) = \lim_{p \rightarrow \infty} \frac{\beta \rho^{3/2}}{p^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) = 0,$$

1429 as $\|\mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbb{I}_p)\| = O(\sqrt{p})$. Since \mathbf{M} and $\boldsymbol{\Sigma}$ share eigenvectors, they commute and we have that $\boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} = \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-1}$. Therefore

$$1430 \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}}[U_p(\mathbf{w})] = - \lim_{d \rightarrow \infty} \mathcal{I}'_1(p) + \lim_{d \rightarrow \infty} \mathcal{I}'_2(p) \quad (57)$$

$$1431 \mathcal{I}'_1(p) = \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}^{-1} \text{vec}(\mathbb{I}_p), \quad \mathcal{I}'_2(p) = \frac{1}{d} \boldsymbol{\theta}^\top \text{vec}(\mathbb{I}_p). \quad (58)$$

1432 For $\mathcal{I}'_1(p)$, using Eqs. (29), (31), (32) and (34) we have that

$$1433 \begin{aligned} 1434 \mathcal{I}'_1(p) &= \frac{1}{d} \mathbf{v}^\top \boldsymbol{\Sigma}^{-2} \mathbf{M}^{-1} \text{vec}(\mathbb{I}_p) = \frac{1}{d} \mathbf{v}^\top ((\mathbf{A}^{-2} \mathbf{B}^{-1}) \otimes \mathbb{I}_p) \text{vec}(\mathbb{I}_p) = \frac{1}{d} \text{vec}(\mathbf{W}_p) \text{vec}(\mathbf{A}^{-2} \mathbf{B}^{-1}) \\ 1435 &= \frac{1}{d} \text{tr}(\mathbf{W}_p \mathbf{A}^{-2} \mathbf{B}^{-1}) = \frac{\rho}{p} \sum_{i=1}^p \frac{\lambda_i}{(c + b\lambda_i)(a(c + b\lambda_i) + 1)} \end{aligned}$$

1458 As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ ,
 1459 yielding

$$1461 \lim_{d \rightarrow \infty} \mathcal{I}'_1(p) = \frac{\rho}{ab^2} \int \frac{\lambda}{(\lambda + c/b)(\lambda + (1 + ac)/ab)} d\mu_\gamma(\lambda). \quad 1462$$

1463 Using Lemma 10 (\mathcal{I}_4 with $z_1 = z_c, z_2 = z_a$), and $z_c - z_a = 1/ab$ we get

$$1465 \lim_{d \rightarrow \infty} \mathcal{I}'_1(p) = \frac{\rho}{b} (z_c m_\gamma(z_c) - z_a m_\gamma(z_a)). \quad (59) \quad 1466$$

1467 Next, for $\mathcal{I}'_2(p)$, using Eqs. (29), (31), (32) and (34), we have that

$$1469 \mathcal{I}'_2(p) = \frac{1}{d} \mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \text{vec}(\mathbb{I}_p) = \frac{1}{d} \text{tr}(\mathbf{W}_p \mathbf{A}^{-1}) = \frac{\rho}{p} \sum_{i=1}^d \frac{\lambda_i}{c + b\lambda_i}. \quad 1470$$

1472 As $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ ,
 1473 yielding

$$1474 \lim_{d \rightarrow \infty} \mathcal{I}'_2(p) = \frac{\rho}{b} \int \frac{\lambda}{(\lambda + c/b)} d\mu_\gamma(\lambda). \quad 1475$$

1477 Using Lemma 10 (\mathcal{I}_1 with $z_1 = z_c$), we get

$$1478 \lim_{d \rightarrow \infty} \mathcal{I}'_2(p) = \frac{\rho}{b} (1 + z_c m_\gamma(z_c)). \quad (60) \quad 1479$$

1480 Substituting Eq. (59) and Eq. (60) in Eq. (57) finishes the proof. \square

1481 **Lemma 6** (Term-1, ID Loss). *Using $\mathbf{M}, \boldsymbol{\Sigma}, \boldsymbol{\Sigma}(\mathbf{w}_j), \mathbf{u}, \boldsymbol{\theta}$ of the same form as Appendix E.2, let*

$$1483 T_1(\mathbf{w}_j) := \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\Sigma}^{-1/2} \mathbf{M}^{-1} \mathbf{u}, \quad 1484$$

$$1485 T_2(\mathbf{w}_j) := \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\theta}, \quad T_3(\mathbf{w}_j) := \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\theta}. \quad 1486$$

1487 Then, almost surely

$$1489 \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_1(\mathbf{w}_j) = (b' + c') T_{1,a} + T_{1,b}, \quad 1490$$

$$1492 \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_2(\mathbf{w}_j) = (b' + c') T_{2,a} + T_{2,b}, \quad 1493$$

$$1495 \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_3(\mathbf{w}_j) = (b' + c') T_{3,a} + T_{3,b} \quad 1497$$

1498 where $T_{1,a} = J_1$ from Eq. (46),

$$1500 T_{1,b} = \frac{\rho}{b^2} \left[-2abz_a z_c (z_c m_\gamma(z_c) - z_a m_\gamma(z_a)) + z_c^2 (m_\gamma(z_c) + z_c m'_\gamma(z_c)) + z_a^2 (m_\gamma(z_a) + z_a m'_\gamma(z_a)) \right] \quad 1501$$

$$1502 + \frac{\rho^2 \beta^2}{a^2 b} \left(1 + \left(z_a - \frac{1}{ab} \right) m_\gamma(z_a) - \frac{1}{ab} z_a m'_\gamma(z_a) \right), \quad 1503$$

$$1504 T_{2,a} = \frac{\rho}{b^2} (1 + 2z_c m_\gamma(z_c) + z_c^2 m'_\gamma(z_c)), \quad 1505$$

$$1506 T_{2,b} = \frac{\rho}{b^2} (1 + 2z_c + 3z_c^2 m_\gamma(z_c) + z_c^3 m'_\gamma(z_c)), \quad 1507$$

$$1508 T_{3,a} = \rho \left(-\frac{c(ac+2)}{b^3} m_\gamma(z_c) + \frac{a}{b} z_a^2 m_\gamma(z_a) + \frac{1}{b^2} z_c^2 m'_\gamma(z_c) \right), \quad 1509$$

$$1510 T_{3,b} = \rho \left(-\frac{c(ac+2)}{b^3} (1 + z_c m_\gamma(z_c)) + \frac{a}{b} z_a^2 (1 + z_a m_\gamma(z_a)) + \frac{1}{b^2} z_c^2 (m_\gamma(z_c) + z_c m'_\gamma(z_c)) \right). \quad 1511$$

1512 *Proof.* Recall

$$1513 \quad \Sigma(\mathbf{w}_j) = \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right) \mathbb{I}_{p^2} + (\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p.$$

1514 We have

$$1515 \quad T_1(\mathbf{w}_j) = \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \Sigma(\mathbf{w}_j) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}$$

$$1516 \quad = \underbrace{\frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}}_{T_{1,a}(\mathbf{w}_j)} + \underbrace{\frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} ((\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}}_{T_{1,b}(\mathbf{w}_j)}.$$

1517 Let us first consider the limit of $T_{1,a}(\mathbf{w}_j)$. We have

$$1518 \quad \lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_{1,a}(\mathbf{w}_j) = \lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{i=1}^M \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}$$

$$1519 \quad = (b' + c') \lim_{d \rightarrow \infty} \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1} \mathbf{M}^{-1} \mathbf{u}, \quad (61)$$

1520 where we use the fact that $\lim_{d \rightarrow \infty} \|\mathbf{w}\|^2/d = 1$ for $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbb{I}_d)$. The above limit is same as Eq. (46) derived in OOD loss.

1521 Next, we solve the limit of $T_{1,b}(\mathbf{w}_j)$. We have

$$1522 \quad \lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_{1,b}(\mathbf{w}_j) = \lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{i=1}^M \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} ((\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}$$

$$1523 \quad = \lim_{d \rightarrow \infty} \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \left(\frac{1}{M} \sum_{j=1}^M (\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p \right) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}$$

$$1524 \quad = \lim_{d \rightarrow \infty} \frac{1}{d^2} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} (\mathbf{W}_p \otimes \mathbb{I}_p) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u}.$$

1525 Using the definitions of Σ, \mathbf{M} we have

$$1526 \quad \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} (\mathbf{W}_p \otimes \mathbb{I}_p) \Sigma^{-1/2} \mathbf{M}^{-1} \mathbf{u} = \mathbf{u}^\top \Sigma^{-1} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \mathbf{u},$$

1527 where we use the fact that $\Sigma, \mathbf{M}, (\mathbf{W}_p \otimes \mathbb{I}_p)$ share Eigenvectors and hence commute. Plugging the definition of \mathbf{u} , we have

$$1528 \quad \mathbf{u}^\top \mathbf{M}^{-2} \Sigma^{-1} \mathbf{W}_p \otimes \mathbb{I}_p \mathbf{u} = \beta^2 \mathbf{g}^\top \Sigma^{-1} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \mathbf{g} + d \boldsymbol{\theta}^\top \Sigma^{-2} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \boldsymbol{\theta}$$

$$1529 \quad - 2\beta d^{1/2} \mathbf{g}^\top \Sigma^{-3/2} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \boldsymbol{\theta}.$$

1530 Following similar steps to Eq. (48),

$$1531 \quad \lim_{d \rightarrow \infty} \frac{1}{d^{3/2}} \mathbf{g}^\top \Sigma^{-3/2} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \boldsymbol{\theta} = \lim_{p \rightarrow \infty} \frac{\rho^{3/2}}{p^{3/2}} \mathbf{g}^\top \Sigma^{-3/2} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \boldsymbol{\theta} = 0,$$

1532 Using Lemma 14 again, we get

$$1533 \quad \frac{1}{d^2} \lim_{d \rightarrow \infty} \mathbf{g}^\top (\Sigma^{-1} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p) \mathbf{g} = \lim_{p \rightarrow \infty} \frac{\rho^2}{p^2} \text{tr}(\Sigma^{-1} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p).$$

1534 Following similar steps to the calculation of Eq. (49) in OOD loss, we get the following integral (with an additional λ) factor in the limit $p, M \rightarrow \infty$

$$1535 \quad \lim_{p \rightarrow \infty} \frac{\rho^2}{p^2} \text{tr}(\Sigma^{-1} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p) = \frac{\rho^2}{a^2 b} \int \lambda \frac{\lambda + c/b}{(\lambda + (1 + ac)/ab)^2} d\mu_{\text{MP}, \gamma}(\lambda).$$

1536 Simplifying this using result in Eq. (49) and Lemma 10 (v) we get

$$1537 \quad \frac{\rho^2}{a^2 b} \int \lambda \frac{\lambda + c/b}{(\lambda + (1 + ac)/ab)^2} d\mu_{\text{MP}, \gamma}(\lambda) = \frac{\rho^2}{a^2 b} \left(1 + \left(z_a - \frac{1}{ab} \right) m_\gamma(z_a) - \frac{1}{ab} z_a m'_\gamma(z_a) \right).$$

1566 Again, using similar steps to result derived in Eq. (50), we get
 1567

$$1568 \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}^{-2} \mathbf{M}^{-2} \mathbf{W}_p \otimes \mathbb{I}_p \boldsymbol{\theta} = \frac{\rho}{a^2 b^4} \int \lambda \frac{\lambda^2}{(\lambda + c/b)^2 (\lambda + (1 + ac)/ab)^2} d\mu_{\text{MP}, \gamma}(\lambda)$$

1571 Using Lemma 10 (v) with Eq. (50) gives the following expression for the above integral
 1572

$$1573 \frac{\rho}{a^2 b^4} \left[\frac{2z_a z_c}{(z_a - z_c)^3} (z_c m_\gamma(z_c) - z_a m_\gamma(z_a)) + \frac{z_c^2}{(z_a - z_c)^2} (m_\gamma(z_c) + z_c m'_\gamma(z_c)) \right. \\ 1574 \left. + \frac{z_a^2}{(z_a - z_c)^2} (m_\gamma(z_a) + z_a m'_\gamma(z_a)) \right] \quad (62)$$

1579 Combining Eq. (61) and Eq. (62), gives the result. Next, we consider $T_2(\mathbf{w}_j)$. It is easy to see that it
 1580 decomposes to the following similar to $T_1(\mathbf{w}_j)$
 1581

$$1582 T_2(\mathbf{w}_j) = \frac{1}{d} \boldsymbol{\theta}^\top \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\theta} \\ 1583 = \underbrace{\frac{1}{d} \boldsymbol{\theta} \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right)}_{T_{2,a}(\mathbf{w}_j)} \boldsymbol{\theta} + \underbrace{\frac{1}{d} \boldsymbol{\theta}^\top ((\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p)}_{T_{2,b}(\mathbf{w}_j)} \boldsymbol{\theta}.$$

1589 The limiting average $\lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_2(\mathbf{w}_j)$, similar to the average for $T_1(\mathbf{w}_j)$, it is easy to see
 1590 that $\lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_{2,a}(\mathbf{w}_j)$ uses the same integral as derived in Eq. (44) for the OOD loss.
 1591

1592 The term $T_{2,b}(\mathbf{w}_j)$ average in the limit
 1593

$$1594 \lim_{d \rightarrow \infty} \frac{1}{dM} \sum_{j=1}^M T_{2,b}(\mathbf{w}_j) = \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \left(\frac{1}{M} \sum_{j=1}^M (\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p \right) \boldsymbol{\theta} \\ 1595 = \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top (\mathbf{W}_p \otimes \mathbb{I}_p) \boldsymbol{\theta}.$$

1600 Following similar steps to the calculation in Eq. (51), we get the following trace expression (with an
 1601 additional \mathbf{W}_p factor)
 1602

$$1603 \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top (\mathbf{W}_p \otimes \mathbb{I}_p) \boldsymbol{\theta} = \lim_{d \rightarrow \infty} \frac{1}{d} \text{tr}(\mathbf{W}_p^3 (c\mathbb{I}_p + b\mathbf{W}_p)^{-2}) \\ 1604 = \rho \int \lambda \frac{\lambda^2}{(\lambda + c/b)^2} d\mu_\gamma(\lambda), \\ 1605 = \frac{\rho}{b^2} (1 + 2z_c m_\gamma(z_c) + 3z_c^2 m_\gamma(z_c) + z_c^3 m'_\gamma(z_c)).$$

1610 where the integral follows from calculation in Eq. (52). The final equality follows by combining
 1611 Lemma 10 (v) with integral solved in Eq. (52).

1612 Next, for $T_3(\mathbf{w}_j)$, we get a similar decomposition to $T_1(\mathbf{w}_j)$ and $T_2(\mathbf{w}_j)$
 1613

$$1614 T_3(\mathbf{w}_j) = \underbrace{\mathbf{u}^\top \mathbf{M}^{-1} \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right)}_{T_{3,a}(\mathbf{w}_j)} \boldsymbol{\theta} + \underbrace{\mathbf{u}^\top \mathbf{M}^{-1} ((\mathbf{w}_j[\mathcal{S}]\mathbf{w}_j[\mathcal{S}]^\top) \otimes \mathbb{I}_p)}_{T_{3,b}(\mathbf{w}_j)} \boldsymbol{\theta}.$$

1618 The limiting average $\lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_3(\mathbf{w}_j)$, similar to the limiting averages computed for
 1619 $T_1(\mathbf{w}_j)$ and $T_2(\mathbf{w}_j)$, it is easy to see that $\lim_{d \rightarrow \infty} \frac{1}{d^2 M} \sum_{j=1}^M T_{3,a}(\mathbf{w}_j)$ uses the integral Eq. (55).

For $\lim_{d \rightarrow \infty} \frac{1}{d^{2M}} \sum_{j=1}^M T_{3,b}(\mathbf{w}_j)$, following similar calculations to Eq. (54), we have the expression below (with an additional \mathbf{W}_p factor)

$$\begin{aligned} \lim_{d \rightarrow \infty} \frac{1}{d^{2M}} \sum_{j=1}^M T_{3,b}(\mathbf{w}_j) &= \frac{1}{d} \text{tr}(\mathbf{W}_p^3 \mathbf{A}^{-3} \mathbf{B}^{-1}) \\ &= \frac{\rho}{ab^2} \int \lambda \frac{\lambda^2}{(\lambda + c/b)^2 (\lambda + (1 + ac)/ab)} d\mu_\gamma(\lambda) \\ &= \rho \left(-\frac{c(ac + 2)}{b^3} (1 + z_c m_\gamma(z_c)) + \frac{a}{b} z_a^2 (1 + z_a m_\gamma(z_a)) + \frac{1}{b^2} z_c^2 (m_\gamma(z_c) + z_c m'_\gamma(z_c)) \right), \end{aligned}$$

where the second equality follows similar steps from Eq. (54) to Eq. (55), and the last inequality follows by combining the result in Eq. (55) with Lemma 10 (v). \square

Lemma 7 (Term-2, ID Loss). *Using $M, \Sigma, \Sigma(\mathbf{w}), \mathbf{u}, \boldsymbol{\theta}$ of the same form as Appendix E.2, let*

$$T_2(\mathbf{w}_j) := \frac{1}{d^{3/2}} \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} (\mathbf{w}_j[\mathcal{S}] \otimes \mathbf{w}_j[\mathcal{S}]) + \frac{1}{d} \boldsymbol{\theta}^\top (\mathbf{w}_j[\mathcal{S}] \otimes \mathbf{w}_j[\mathcal{S}]).$$

Then, almost surely

$$\lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_2(\mathbf{w}_j) = -\frac{\rho}{b} (z_c (1 + z_c m_\gamma(z_c)) - z_a (1 + z_a m_\gamma(z_a))) + \frac{\rho}{b} (1 + z_c (1 + z_c m_\gamma(z_c)))$$

Proof. Consider the average of the first term. Using $\frac{1}{M} \sum_{j=1}^M \mathbf{w}_j[\mathcal{S}] \otimes \mathbf{w}_j[\mathcal{S}] = \text{vec}(\mathbf{W}_p)$, we have

$$\frac{1}{Md^{3/2}} \sum_{j=1}^M \mathbf{u}^\top \mathbf{M}^{-1} \Sigma^{-1/2} (\mathbf{w}_j[\mathcal{S}] \otimes \mathbf{w}_j[\mathcal{S}]) = \frac{\beta}{d^{3/2}} \mathbf{g}^\top \mathbf{M}^{-1} \Sigma^{-1/2} \text{vec}(\mathbf{W}_p) - \frac{1}{d} \boldsymbol{\theta}^\top \Sigma^{-1/2} \mathbf{M}^{-1} \Sigma^{-1/2} \text{vec}(\mathbf{W}_p).$$

Term 1 above is 0 in the limit $d \rightarrow \infty$, using Lemma 14. Term 2 simplifies to

$$\begin{aligned} \frac{1}{d} \boldsymbol{\theta}^\top \Sigma^{-1/2} \mathbf{M}^{-1} \Sigma^{-1/2} \text{vec}(\mathbf{W}_p) &= \frac{1}{d} \text{vec}(\mathbf{W}_p)^\top \Sigma^{-2} \mathbf{M}^{-1} \text{vec}(\mathbf{W}_p) = \frac{1}{d} \text{vec}(\mathbf{W}_p)^\top \mathbf{A}^{-2} \mathbf{B}^{-1} \otimes \mathbb{I}_p \text{vec}(\mathbf{W}_p) \\ &= \frac{1}{d} \text{vec}(\mathbf{W}_p)^\top \text{vec}(\mathbf{A}^{-2} \mathbf{B}^{-1} \mathbf{W}_p) \\ &= \frac{1}{d} \text{tr}(\mathbf{W}_p^2 \mathbf{A}^{-2} \mathbf{B}^{-1}). \end{aligned}$$

Here, second equality follows from the definition of $\boldsymbol{\theta}$ and the fact that \mathbf{M} and Σ commute. Third and fourth equality use the definition of Σ, \mathbf{M} and Eq. (31) and Eq. (32). In the limit

$$\begin{aligned} \lim_{d \rightarrow \infty} \frac{1}{d} \text{tr}(\mathbf{W}_p^2 \mathbf{A}^{-2} \mathbf{B}^{-1}) &= \frac{\rho}{ab^2} \int \frac{\lambda}{(\lambda + c/b)(\lambda + (1 + ac)/ab)} d\mu_\gamma(\lambda) \\ &= \frac{\rho}{b} (z_c (1 + z_c m_\gamma(z_c)) - z_a (1 + z_a m_\gamma(z_a))), \end{aligned}$$

where second equality follows by Lemma 10 with Eq. (59). Next, for term 2

$$\begin{aligned} \lim_{d \rightarrow \infty} \frac{1}{d} \boldsymbol{\theta}^\top \text{vec}(\mathbf{W}_p) &= \lim_{d \rightarrow \infty} \frac{1}{d} \text{tr}(\mathbf{W}_p^2 \mathbf{A}^{-1}) \\ &= \frac{\rho}{b} \int \lambda \frac{\lambda}{(\lambda + c/b)} d\mu_\gamma(\lambda) \\ &= \frac{\rho}{b} (1 + z_c (1 + z_c m_\gamma(z_c))). \end{aligned}$$

The first equality follows by using definition of $\boldsymbol{\theta}, \Sigma$ and Eq. (31) and Eq. (32). Second equality follows by the fact the empirical distribution of \mathbf{W}_p converges to Marchenko-Pastur law μ_γ in the limit, and the final equality uses Lemma 10 with Eq. (60). \square

We restate Theorem 3 below followed by the proof.

Theorem 6 (Asymptotic Risk of Minimum-Norm LGP). *Consider the minimum norm estimator $\hat{\theta}_{S_p}$ for the linear Gaussian equivalent problem in Definition 3 under the overparameterized regime ($\nu < \rho^2$). Under the asymptotic scaling defined in Eq. (1), as $d \rightarrow \infty$, the ID and OOD risks converge in probability to the following deterministic limits:*

$$\mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) \xrightarrow{P} \left(1 - \frac{\nu}{\rho^2}\right) \rho(1+ac)(1-c(1+a+ac)) + (1+a+ac)(\sigma_n^2 + \rho c + 1 - \rho), \quad (63)$$

$$\mathcal{L}_{\text{ID}}(\hat{\theta}_{S_p}) \xrightarrow{P} \frac{\nu}{a^2} \bar{\beta}^2. \quad (64)$$

Proof. From Eq. (18) and using Eqs. (19) to (22), we have

$$\mathbf{a} = \mathbf{M}^{-1}(\beta \mathbf{g} - \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \boldsymbol{\theta}_*) = \frac{1}{d} \bar{\mathbf{M}}^{-1}(\bar{\beta} \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*).$$

Substituting this in Lemma 1, we get Term-1(a) as

$$\frac{1}{d^2} (\bar{\beta}_* \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*)^\top \bar{\mathbf{M}}^{-1} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\mathbf{M}}^{-1} (\bar{\beta}_* \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*).$$

Similarly, Term-1(c) is

$$\mathbf{a}^\top \Sigma_{\mathbf{H}_{S_p}}^{-1/2} \Sigma(\mathbf{w}) \boldsymbol{\theta}_* = \frac{1}{d^{3/2}} (\bar{\beta}_* \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*)^\top \bar{\mathbf{M}}^{-1} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\boldsymbol{\theta}}_*$$

Similarly, Term-2 is

$$\hat{\boldsymbol{\theta}}_{S_p}^\top \mathbb{E}_{X,y}[\mathbf{h}_{S_p}(X) y_q] = \frac{1}{d^{3/2}} (\bar{\beta}_* \mathbf{g} - d^{1/2} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\boldsymbol{\theta}}_*)^\top \bar{\mathbf{M}}^{-1} \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \bar{\boldsymbol{\theta}}_*^\top \mathbf{w} \otimes \mathbf{w}.$$

OOD Loss. Using these, and Lemma 4 and Lemma 5 with $\Sigma = \bar{\Sigma}_{\mathbf{H}_{S_p}}$, $\boldsymbol{\theta} = \bar{\boldsymbol{\theta}}_*$, $\mathbf{M} = \bar{\mathbf{M}}$, $b = 1$, $b' + c' = c$, $\gamma = \mu^{-1}$, $z_a = -\frac{1}{a} - c$, $z_c = -c$, and defining $m_a := m_\gamma(z_a)$, $m'_a := m'_\gamma(z_a)$, $m_c := m_\gamma(-c)$, $m'_c := m'_\gamma(-c)$, we have

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) &= (1+c) \underbrace{\left(\bar{\beta}^2 \rho^2 \left(\frac{1}{a^2} m_a - \frac{1}{a^3} m'_a \right) + \rho(c^2 m'_c + 2z_a a c(m_c - m_a) + z_a^2 m'_a) \right)}_{\text{Term-1(a)}} \\ &\quad + \underbrace{(1+c)\rho(1-2cm_c + c^2 m'_c)}_{\text{Term-1(b)}} \underbrace{-2(1+c)\rho[-c(ac+2)m_c + az_a^2 m_a + c^2 m'_c]}_{\text{Term-1(c)}} \\ &\quad \underbrace{-2\rho(1+z_a m_a)}_{\text{Term-2}} + 1 + \sigma_n^2. \end{aligned}$$

Using the expression for $\bar{\beta}$ from Eq. (35), Term-1(a) can be simplified as:

$$(1+c) \left(a(\sigma_n^2 + 1 - \rho + \rho c) + \rho z_a^2 a m_a + \rho c^2 m'_c - 2\rho(1+ac)cm_c \right)$$

Combining all terms, we get

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) &= m_a(-z_a)\rho(-az_a(1+c) + 2az_a(1+c) + 2) + m_c(1+c)\rho(2z_a ac - 2c + 2c(ac+2)) \\ &\quad + m'_c(1+c)\rho c^2(1+1-2) + (1+c)a(\sigma_n^2 + 1 - \rho + \rho c) + \rho(1+c-2) + 1 + \sigma_n^2. \end{aligned}$$

Simplifying this and using $m_a = a(1 - \nu/\rho^2)$, we get Eq. (63).

ID Loss. Using the expressions for Term-1(a), Term-1(c) and Term-2 computed above, Lemmas 1, 6 and 7, we have

$$\begin{aligned}
\mathcal{L}_{\text{ID}}(\hat{\boldsymbol{\theta}}_{S_p}) &= c \underbrace{\left(\bar{\beta}^2 \rho^2 \left(\frac{1}{a^2} m_a - \frac{1}{a^3} m'_a \right) + \rho(c^2 m'_c + 2z_a a c(m_c - m_a) + z_a^2 m'_a) \right)}_{\text{Term-1(a)-old part}} \\
&+ \underbrace{\bar{\beta}^2 \rho^2 \left(\frac{1}{a^2} (1 + z_a m_a) - \frac{1}{a^3} (m_a + z_a m'_a) \right) + \rho(c^2(m_c + z_c m'_c) + 2z_a a c(z_c m_c - z_a m_a) + z_a^2(m_a + z_a m'_a))}_{\text{Term-1(a)-new part}} \\
&+ \underbrace{c\rho(1 - 2c m_c + c^2 m'_c)}_{\text{Term-1(b)-old}} + \underbrace{\rho(1 - 2c + 3c^2 m_c - c^3 m'_c)}_{\text{Term-1(b)-new}} \\
&- \underbrace{2c\rho[-c(ac + 2)m_c + az_a^2 m_a + c^2 m'_c]}_{\text{Term-1(c)-old}} - \underbrace{2\rho[a^{-1} + az_a^3 m_a + (3c^2 + ac^3)m_c - c^3 m'_c]}_{\text{Term-1(c)-new}} \\
&- \underbrace{2\rho(1 - c + c^2 m_c - a^{-1} - z_c^2 m_c + z_a^2 m_a)}_{\text{Term-2-new}} + 1 + \sigma_n^2.
\end{aligned}$$

Combining terms, we have

$$\begin{aligned}
\mathcal{L}_{\text{ID}}(\hat{\boldsymbol{\theta}}_{S_p}) &= \bar{\beta}^2 \rho^2 \frac{1}{a^2} \left((z_a + c) \left(m_a - \frac{1}{a} m'_a \right) + 1 - \frac{1}{a} m_a \right) + m_a \rho z_a (-2ac^2 - 2z_a a c + z_a - 2ac z_a - 2az_a^2 - 2z_a) \\
&+ \rho m'_a (cz_a^2 + z_a^3) + \rho m_c (2ac^2 z_a + c^2 - 2z_a a c^2 - 2c^2 + 3c^2 + 2c^2(ac + 2) - 2(3c^2 + ac^3)) \\
&+ \rho m'_c (c^3 - c^3 + c^3 - c^3 - 2c^3 + 2c^3) + 1 + \sigma_n^2 + \rho(-2(1 - c) + 1 - 2c + c) \\
&= \bar{\beta}^2 \rho^2 \frac{1}{a^2} \left(-\frac{1}{a} \left(m_a - \frac{1}{a} m'_a \right) + \frac{\nu}{\rho^2} \right) + \frac{\rho(1 - a^2 c^2)}{a^2} m_a - \frac{\rho z_a^2}{a} m'_a + \sigma_n^2 + 1 - \rho + \rho c.
\end{aligned}$$

Using the expression for $\bar{\beta}$ from Eq. (35), we simplify and get Eq. (64). □

G.2 ASYMPTOTIC RISK OF LS LGP

Lemma 8 (OOD Loss). *Using $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}(\boldsymbol{w})$, $\boldsymbol{\theta}$ of the same form as Appendix E.2, let*

$$\begin{aligned}
T_1(\boldsymbol{w}) &:= \frac{1}{\|\boldsymbol{g}\|^2} \boldsymbol{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\boldsymbol{w}) \boldsymbol{\Sigma}^{-1/2} \boldsymbol{g}, \\
T_2(\boldsymbol{w}) &:= \frac{1}{\sqrt{d}\|\boldsymbol{g}\|} \boldsymbol{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\boldsymbol{w}) \boldsymbol{\theta}, \quad T_3(\boldsymbol{w}) := \frac{1}{\sqrt{d}\|\boldsymbol{g}\|} \boldsymbol{g}^\top \boldsymbol{\Sigma}^{-1/2} (\boldsymbol{w} \otimes \boldsymbol{w}).
\end{aligned}$$

then

$$\lim_{d \rightarrow \infty} \mathbb{E}_{\boldsymbol{w}}[T_1(\boldsymbol{w})] = \frac{1 + b' + c'}{b} m_\gamma(z_c), \quad \lim_{d \rightarrow \infty} \mathbb{E}_{\boldsymbol{w}}[T_2(\boldsymbol{w})] = 0, \quad \lim_{d \rightarrow \infty} \mathbb{E}_{\boldsymbol{w}}[T_3(\boldsymbol{w})] = 0.$$

Proof. We first work with $T_1(\boldsymbol{w})$. Using Lemma 11, since \boldsymbol{w} is independent of \boldsymbol{g} , \boldsymbol{W}_p , we get

$$\mathbb{E}_{\boldsymbol{w}}[T_1(\boldsymbol{w})] = \frac{(c' + b' + 1)}{\|\boldsymbol{g}\|^2} \boldsymbol{g}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{g}.$$

Using Lemma 14, and since $\lim_{d \rightarrow \infty} \frac{p}{\|\boldsymbol{g}\|} = 1$, we have that

$$\lim_{d \rightarrow \infty} \frac{1}{\|\boldsymbol{g}\|^2} \boldsymbol{g}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{g} = \lim_{d \rightarrow \infty} \frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^{-1}).$$

Using Eqs. (29) and (30), we have

$$\frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^{-1}) = \frac{1}{p} \text{tr}(\boldsymbol{A}^{-1}) = \frac{1}{p} \sum_{i=1}^p \frac{1}{c + b\lambda_i}. \tag{65}$$

As $p \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ , yielding

$$\lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}}[T_1(\mathbf{w})] = \frac{1}{b} \int \frac{1}{\lambda + c/b} d\mu_\gamma(\lambda) = \frac{1}{b} m_\gamma(z_c).$$

Next, for $T_2(\mathbf{w})$, using Lemma 11, since \mathbf{w} is independent of \mathbf{g} , \mathbf{W}_p , we have that

$$\mathbb{E}_{\mathbf{w}}[T_2(\mathbf{w})] = \frac{(c' + b' + 1)}{\sqrt{d}\|\mathbf{g}\|} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta}.$$

Using Lemma 14, it follows that

$$\lim_{d \rightarrow \infty} \frac{1}{p^{3/2}} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\theta} = 0,$$

which concludes the proof for the second part.

Finally, working with $T_3(\mathbf{w})$, Using $\mathbb{E}_{\mathbf{w}}[\mathbf{w} \otimes \mathbf{w}] = \text{vec}(\mathbb{I}_p)$ and since \mathbf{w} is independent of \mathbf{g} and \mathbf{W}_p , we have that

$$\mathbb{E}_{\mathbf{w}}[T_3(\mathbf{w})] = \frac{1}{\sqrt{d}\|\mathbf{g}\|} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p).$$

Using Lemma 14, it follows that

$$\lim_{d \rightarrow \infty} \frac{1}{p^{3/2}} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \text{vec}(\mathbb{I}_p) = 0,$$

which concludes the proof for this part. \square

Lemma 9 (ID Loss). *Using $\boldsymbol{\Sigma}$, $\boldsymbol{\Sigma}(\mathbf{w}_j)$, \mathbf{u} , $\boldsymbol{\theta}$ of the same form as Appendix E.2, let*

$$T_1(\mathbf{w}_j) := \frac{1}{\|\mathbf{g}\|^2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\Sigma}^{-1/2} \mathbf{g},$$

$$T_2(\mathbf{w}_j) := \frac{1}{\sqrt{d}\|\mathbf{g}\|} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \boldsymbol{\Sigma}(\mathbf{w}_j) \boldsymbol{\theta}, \quad T_3(\mathbf{w}_j) := \frac{1}{\sqrt{d}\|\mathbf{g}\|} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} (\mathbf{w}_j \otimes \mathbf{w}_j).$$

Then, almost surely

$$\lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_1(\mathbf{w}_j) = \frac{1}{b} (1 + z_c m_\gamma(z_c)), \quad \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_2(\mathbf{w}_j) = 0, \quad \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{j=1}^M T_3(\mathbf{w}_j) = 0.$$

Proof. We first work with $T_1(\mathbf{w}_j)$. The average is

$$\frac{1}{M} \sum_{j=1}^M T_1(\mathbf{w}_j) = \frac{1}{M} \frac{1}{\|\mathbf{g}\|^2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} \left(c' \frac{\|\mathbf{w}_j\|^2}{d} + b' \right) \boldsymbol{\Sigma}^{-1/2} \mathbf{g} + \frac{1}{\|\mathbf{g}\|^2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} (\mathbf{W}_p \otimes \mathbb{I}_p) \boldsymbol{\Sigma}^{-1/2} \mathbf{g}.$$

This just follows using definition of $\boldsymbol{\Sigma}(\mathbf{w}_j)$. Next, in the limit $d \rightarrow \infty$, term 1 follows using Eq. (65). The limit of second quantity above

$$\begin{aligned} \lim_{d \rightarrow \infty} \frac{1}{\|\mathbf{g}\|^2} \mathbf{g}^\top \boldsymbol{\Sigma}^{-1/2} (\mathbf{W}_p \otimes \mathbb{I}_p) \boldsymbol{\Sigma}^{-1/2} \mathbf{g} &= \lim_{d \rightarrow \infty} \frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^{-1/2} (\mathbf{W}_p \otimes \mathbb{I}_p) \boldsymbol{\Sigma}^{-1/2}) \\ &= \lim_{p \rightarrow \infty} \frac{1}{p^2} \text{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{W}_p \otimes \mathbb{I}_p) \\ &= \lim_{p \rightarrow \infty} \frac{1}{p} \text{tr}(\mathbf{A}^{-1} \mathbf{W}_p) \\ &= \frac{1}{b} \int \frac{\lambda}{\lambda + c/b} d\mu_\gamma(\lambda) = \frac{1}{b} (1 + z_c m_\gamma(z_c)). \end{aligned}$$

Here, we use commutation of $\boldsymbol{\Sigma}$ and $\mathbf{W}_p \otimes \mathbb{I}_p$ and properties Eq. (30), Eq. (32). The final equality follows as when $p, M \rightarrow \infty$, the empirical spectral distribution of \mathbf{W}_p converges to the Marchenko-Pastur law μ_γ .

Next, the limiting averages of $T_2(\mathbf{w}_j)$ and $T_3(\mathbf{w}_j)$, using Lemma 14, and similar steps used in Lemma 8 both evaluate to 0. \square

We restate Theorem 4 below followed by its proof.

Theorem 7 (Asymptotic Risk of LS LGP). *Consider the least-squares estimator $\hat{\theta}_{S_p}$ for the linear Gaussian equivalent problem in the underparameterized regime ($\nu > \rho^2$). Under the asymptotic scaling in Eq. (1), as $d \rightarrow \infty$, the ID and OOD risks converge in probability to:*

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) &\xrightarrow{P} m_\gamma(z_c)(\kappa_\infty^2(1+c) - 2c^2\rho) + (1+c)\rho c^2 m'_\gamma(z_c) + \sigma_n^2 + 1 - \rho + \rho c, \\ \mathcal{L}_{\text{ID}}(\hat{\theta}_{S_p}) &\xrightarrow{P} \frac{\nu}{\rho^2} \kappa_\infty^2, \end{aligned}$$

where κ_∞ is defined in Eq. (9).

Proof. Substituting $\mathbf{a} = \frac{\kappa_\infty}{\|\mathbf{g}\|} \mathbf{g}$ and using Eq. (27) and Eq. (28) in Eq. (26), we have that

$$\begin{aligned} \mathcal{L}(\mathbf{a}; \mathbf{w}) &= \frac{\kappa_\infty^2}{\|\mathbf{g}\|^2} \mathbf{g}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \mathbf{g} + \frac{1}{d} \bar{\theta}_*^\top \bar{\Sigma}(\mathbf{w}) \bar{\theta}_* + 2 \frac{\kappa_\infty}{\sqrt{d} \|\mathbf{g}\|} \mathbf{g}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} \bar{\Sigma}(\mathbf{w}) \bar{\theta}_* \\ &\quad - 2 \left(\frac{\kappa_\infty}{d^{1/2} \|\mathbf{g}\|} \mathbf{g}^\top \bar{\Sigma}_{\mathbf{H}_{S_p}}^{-1/2} (\mathbf{w} \otimes \mathbf{w}) + \frac{1}{d} \bar{\theta}_*^\top (\mathbf{w} \otimes \mathbf{w}) \right) + 1 + \sigma_n^2. \end{aligned}$$

Using Lemma 8, Eq. (44), Eq. (60), and $b = 1, c' + b' = c$, we have that

$$\begin{aligned} \mathcal{L}_{\text{OOD}}(\hat{\theta}_{S_p}) &:= \lim_{d \rightarrow \infty} \mathbb{E}_{\mathbf{w}}[\mathcal{L}(\mathbf{a}; \mathbf{w})] \\ &= \kappa_\infty^2(1+c)m_\gamma(z_c) + (1+c)\rho(1+2z_c m_\gamma(z_c) + z_c^2 m'_\gamma(z_c)) - 2\rho(1+z_c m_\gamma(z_c)) + 1 + \sigma_n^2 \\ &= m_\gamma(z_c)(\kappa_\infty^2(1+c) - 2c^2\rho) + (1+c)\rho c^2 m'_\gamma(z_c) + \sigma_n^2 + 1 - \rho + \rho c, \end{aligned}$$

where using Eq. (43),

$$\kappa_\infty^2 = \frac{c_\infty}{\nu \left(\frac{\nu}{\rho^2} - 1 \right)} = \frac{\sigma_n^2 + 1 - \rho + \rho c - \rho c^2 m_\gamma(z_c)}{\frac{\nu}{\rho^2} - 1}.$$

Similarly, using Lemma 9, we have

$$\begin{aligned} \mathcal{L}_{\text{ID}}(\hat{\theta}_{S_p}) &:= \lim_{d \rightarrow \infty} \frac{1}{M} \sum_{i=1}^M \mathcal{L}(\mathbf{a}; \mathbf{w}_i) \\ &= \kappa_\infty^2 (c m_\gamma(z_c) + 1 - c m_\gamma(z_c)) + c\rho(1 - 2c m_\gamma(z_c) + c^2 m'_\gamma(z_c)) \\ &\quad + \rho(1 - 2c + 3c^2 m_\gamma(z_c) - c^3 m'_\gamma(z_c)) - 2\rho(1 - c + c^2 m_\gamma(z_c)) + 1 + \sigma_n^2 \\ &= \kappa_\infty^2 - \rho c^2 m_\gamma(z_c) + \sigma_n^2 + 1 - \rho + \rho c = \frac{\nu}{\rho^2} \kappa_\infty^2. \end{aligned}$$

□

H HELPER LEMMAS

Lemma 10. *Let $z_1, z_2, z_3, \alpha_0, \alpha_1, \alpha_2 \in \mathbb{R}$ be constants and $f(\cdot)$ be a function. Using*

$$m_\gamma(z) = \int \frac{1}{\lambda - z} d\mu_\gamma(\lambda), \quad m'_\gamma(z) = \int \frac{1}{(\lambda - z)^2} d\mu_\gamma(\lambda),$$

$$(i) \mathcal{I}_1 := \int \frac{\lambda - z_1}{\lambda - z_2} d\mu_\gamma(\lambda) = 1 + (z_2 - z_1)m_\gamma(z_2),$$

$$(ii) \mathcal{I}_2 := \int \frac{\lambda - z_1}{(\lambda - z_2)^2} d\mu_\gamma(\lambda) = m_\gamma(z_2) + (z_2 - z_1)m'_\gamma(z_2),$$

$$(iii) \mathcal{I}_3 := \int \frac{\lambda^2}{(\lambda - z_1)^2 (\lambda - z_2)} d\mu_\gamma(\lambda) = \frac{1}{(z_1 - z_2)^2} (z_1(z_1 - 2z_2)m_\gamma(z_1) + z_2^2 m_\gamma(z_2) + z_1^2 (z_1 - z_2)m'_\gamma(z_1)),$$

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

$$(iv) \mathcal{I}_4 := \int \frac{\lambda}{(\lambda-z_1)(\lambda-z_2)} d\mu_\gamma(\lambda) = \frac{1}{z_1-z_2} (z_1 m_\gamma(z_1) - z_2 m_\gamma(z_2)),$$

(v) if $\int f(\lambda) d\mu_\gamma(\lambda) = \alpha_0 + \alpha_1 m_\gamma(z_1) + \alpha_2 m'_\gamma(z_1)$, then

$$\int \lambda f(\lambda) d\mu_\gamma(\lambda) = \alpha_0 + \alpha_1 + (\alpha_1 z_1 + \alpha_2) m_\gamma(z_1) + \alpha_2 z_1 m'_\gamma(z_1),$$

$$(vi) \mathcal{I}_6 := \int \frac{\lambda^2}{(\lambda-z_1)^2} d\mu_\gamma(\lambda) = 1 + 2z_1 m_\gamma(z_1) + z_1^2 m'_\gamma(z_1),$$

$$(vii) \mathcal{I}_7 := \int \frac{\lambda^2}{(\lambda-z_1)^2(\lambda-z_2)^2} d\mu_\gamma(\lambda) = \frac{1}{(z_1-z_2)^2} \left(z_1^2 m'_\gamma(z_1) + z_2^2 m'_\gamma(z_2) - \frac{2z_1 z_2}{z_1-z_2} (m_\gamma(z_1) - m_\gamma(z_2)) \right).$$

Proof. The proof relies on obtaining partial fraction decompositions and using the definitions of $m_\gamma(z), m'_\gamma(z)$.

First, for \mathcal{I}_1 , note that

$$\frac{\lambda - z_1}{\lambda - z_2} = 1 + \frac{z_2 - z_1}{\lambda - z_2}, \quad \mathcal{I}_1 = 1 + (z_2 - z_1) m_\gamma(z_2).$$

Next, for \mathcal{I}_2 , we have

$$\frac{\lambda - z_1}{(\lambda - z_2)^2} = \frac{1}{\lambda - z_2} + \frac{z_2 - z_1}{(\lambda - z_2)^2} \implies \mathcal{I}_2 = m_\gamma(z_2) + (z_2 - z_1) m'_\gamma(z_2).$$

Next, for \mathcal{I}_3 , we have

$$\frac{\lambda^2}{(\lambda - z_1)^2(\lambda - z_2)} = \frac{z_1^2}{z_1 - z_2} \frac{1}{(\lambda - z_1)^2} + \frac{z_1(z_1 - 2z_2)}{(z_1 - z_2)^2} \frac{1}{\lambda - z_1} + \frac{z_2^2}{(z_1 - z_2)^2} \frac{1}{\lambda - z_2},$$

which gives $\mathcal{I}_3 = \frac{1}{(z_1-z_2)^2} (z_1(z_1 - 2z_2) m_\gamma(z_1) + z_2^2 m_\gamma(z_2) + z_1^2 (z_1 - z_2) m'_\gamma(z_1))$.

Next, for \mathcal{I}_4 , we have

$$\frac{\lambda}{(\lambda - z_1)(\lambda - z_2)} = \frac{z_1}{z_1 - z_2} \frac{1}{\lambda - z_1} + \frac{z_2}{z_2 - z_1} \frac{1}{\lambda - z_2},$$

which gives $\mathcal{I}_4 = \frac{z_1}{z_1-z_2} m_\gamma(z_1) + \frac{z_2}{z_2-z_1} m_\gamma(z_2)$.

For the next part, we have

$$\alpha_1 m_\gamma(z_1) + \alpha_2 m'_\gamma(z_1) = \alpha_0 + \alpha_1 \int \frac{1}{\lambda - z_1} d\mu_\gamma(\lambda) + \alpha_2 \int \frac{1}{(\lambda - z_1)^2} d\mu_\gamma(\lambda)$$

This gives

$$\begin{aligned} \int \lambda f(\lambda) d\mu_\gamma(\lambda) &= \alpha_0 \int \lambda d\mu_\gamma(\lambda) + \alpha_1 \int \frac{\lambda}{\lambda - z_1} d\mu_\gamma(\lambda) + \alpha_2 \int \frac{\lambda}{(\lambda - z_1)^2} d\mu_\gamma(\lambda) \\ &= \alpha_0 + \alpha_1 (1 + z_1 m_\gamma(z_1)) + \alpha_2 (m_\gamma(z_1) + z_1 m'_\gamma(z_1)), \end{aligned}$$

where we used \mathcal{I}_1 for the first term and \mathcal{I}_2 for second term. Simplifying this finishes the proof for this part.

Next, for \mathcal{I}_6 , we have

$$\frac{\lambda^2}{(\lambda - z_1)^2} = \left(1 + \frac{z_1}{\lambda - z_1} \right)^2 = 1 + \frac{2z_1}{\lambda - z_1} + \frac{z_1^2}{(\lambda - z_1)^2},$$

which gives $\mathcal{I}_6 = 1 + 2z_1 m_\gamma(z_1) + z_1^2 m'_\gamma(z_1)$.

Next, for \mathcal{I}_7 , we have

$$\begin{aligned} \frac{\lambda^2}{(\lambda - z_1)^2(\lambda - z_2)^2} &= \frac{1}{(z_1 - z_2)^2} \left(\frac{z_1}{\lambda - z_1} - \frac{z_2}{\lambda - z_2} \right)^2 \\ &= \frac{1}{(z_1 - z_2)^2} \left(\frac{z_1^2}{(\lambda - z_1)^2} + \frac{z_2^2}{(\lambda - z_2)^2} - \frac{2z_1 z_2}{z_1 - z_2} \left(\frac{1}{\lambda - z_1} - \frac{1}{\lambda - z_2} \right) \right), \end{aligned}$$

which gives $\mathcal{I}_7 = \frac{1}{(z_1-z_2)^2} \left(z_1^2 m'_\gamma(z_1) + z_2^2 m'_\gamma(z_2) - \frac{2z_1 z_2}{z_1-z_2} (m_\gamma(z_1) - m_\gamma(z_2)) \right)$. \square

1944 **Lemma 11.** $\mathbb{E}_{\mathbf{w}}[\Sigma(\mathbf{w})] = (c' + b' + 1)\mathbb{I}_{d^2}$.

1945

1946

1947 *Proof.* Since $\mathbf{w} \sim \mathcal{N}(0, \mathbb{I}_d)$, we have $\mathbb{E}[\mathbf{w}\mathbf{w}^\top] = \mathbb{I}_d$ and $\mathbb{E}[\|\mathbf{w}\|^2/d] = 1$, hence

1948

$$1949 \quad \mathbb{E}_{\mathbf{w}}[\Sigma(\mathbf{w})] = \left(c'\mathbb{E}\left[\frac{\|\mathbf{w}\|^2}{d}\right] + b'\right)\mathbb{I}_{d^2} + \mathbb{E}[(\mathbf{w}\mathbf{w}^\top) \otimes \mathbb{I}_d] = (c' + b' + 1)\mathbb{I}_{d^2}.$$

1950

1951

1952 **Lemma 12** (Gaussian tail bound). *If $X \sim \mathcal{N}(0, \sigma^2)$, then for all $t > 0$,*

1953

1954

$$1955 \quad \Pr(|X| \geq t) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

1956

1957 **Lemma 13** (Sub-exponential tail bound). *If X_i are i.i.d., zero mean, sub-exponential random variables, then there exist universal constants $K, c_0 > 0$ such that $\|X_i\|_{\psi_1} \leq K$ and*

1958

1959

$$1960 \quad \Pr\left(\left|\sum_{i=1}^d a_i X_i\right| \geq t\right) \leq 2 \exp\left(-c_0 \min\left\{\frac{t^2}{K^2\|a\|_2^2}, \frac{t}{K\|a\|_\infty}\right\}\right) \quad \text{for all } t > 0.$$

1961

1962

1963

1964

Lemma 14. *Let $\mathbf{a}_1 \in \mathbb{R}^{d_1}$ and $\mathbf{A}_1 \in \mathbb{R}^{d_1 \times d_1}$, s.t. $\|\mathbf{a}_1\| = O(1/d_1^q)$, $q > 0$ and eigenvalues of \mathbf{A}_1 are $\Theta(1)$. Then,*

1965

1966

1967

1968

1969

$$\lim_{d_1 \rightarrow \infty} \mathbf{g}^\top \mathbf{a}_1 = 0$$

$$\lim_{d_1 \rightarrow \infty} \frac{1}{d_1} \mathbf{g}^\top \mathbf{A}_1 \mathbf{g} = \lim_{d_1 \rightarrow \infty} \frac{1}{d_1} \text{tr}(\mathbf{A}_1).$$

1970

1971

1972

Proof. For the first, we know that if $\mathbf{g} \sim \mathcal{N}(0, \mathbb{I}_{d_1})$, then $\mathbf{g}^\top \mathbf{a}_1 \sim \mathcal{N}(0, \|\mathbf{a}_1\|^2)$. Using standard Gaussian tail bounds from Lemma 12 with $\|\mathbf{a}_1\| = O(1/d_1^q)$, we have $\lim_{d_1 \rightarrow \infty} \mathbf{g}^\top \mathbf{a}_1 = 0$.

1973

Next, let $\mathbf{A}_1 = \mathbf{V}\Lambda\mathbf{V}^\top$ be the Eigenvalue decomposition of \mathbf{A}_1 . We have

1974

1975

1976

1977

$$\mathbf{g}^\top \mathbf{A}_1 \mathbf{g} = (\mathbf{V}^\top \mathbf{g})^\top \Lambda (\mathbf{V}^\top \mathbf{g}) = \frac{1}{d_1} \sum_{i=1}^d \lambda_i \tilde{g}_i^2,$$

1978

1979

1980

where $\mathbf{V}^\top \mathbf{g} =: \tilde{\mathbf{g}} \sim \mathcal{N}(0, \mathbb{I}_{d_1})$. We know that $X_i := \tilde{g}_i^2 - 1$ are 0 mean i.i.d sub-exponential random variables for all $i \in [d_1]$. Applying sub-exponential tail bound from Lemma 13 to $\sum_{i=1}^d \lambda_i X_i$, we have

1981

1982

1983

1984

1985

1986

$$\Pr\left(\left|\frac{1}{d_1} \sum_{i=1}^{d_1} \lambda_i X_i\right| \geq t\right) = 2 \exp\left(-c_0 d_1 \min\left\{\frac{d_1 t^2}{\|\Lambda\|_F^2}, \frac{t}{\lambda_{\max}}\right\}\right) \quad \text{for all } t > 0.$$

$$\leq 2 \exp\left(-c_0 d_1 \min\left\{\frac{t^2}{C^2}, \frac{t}{C}\right\}\right),$$

1987

where inequality uses $\lambda_i = \Theta(1)$. Therefore

1988

1989

1990

1991

$$\lim_{d_1 \rightarrow \infty} \frac{1}{d_1} \sum_{i=1}^{d_1} \lambda_i X_i = \lim_{d_1 \rightarrow \infty} \frac{1}{d_1} (\mathbf{g}^\top \mathbf{A}_1 \mathbf{g} - \text{tr}(\mathbf{A}_1)) = 0,$$

1992

where we use that $\sum_{i=1}^{d_1} \lambda_i = \text{tr}(\Lambda) = \text{tr}(\mathbf{A}_1)$. □

1993

1994

1995

1996

1997

Lemma 15. *Let $\mathbf{x}_t \in \mathbb{R}^d$ have i.i.d. entries $\mathbf{x}_t[i] \sim \mathcal{N}(0, \sigma_x^2)$, and let $\mathbf{w} \in \mathbb{R}^d$ be fixed. For a subset $\mathcal{S} \subseteq [d]$ with $|\mathcal{S}| = p$, denote by $\mathbf{x}_t[\mathcal{S}] \in \mathbb{R}^p$ and $\mathbf{w}[\mathcal{S}] \in \mathbb{R}^p$ the corresponding coordinate sub-vectors. Then*

$$\mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] = \sigma_x^4 \|\mathbf{w}\|^2 \mathbb{I}_p + 2\sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top.$$

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Proof. Write $\mathbf{x}_t = \sigma_x \mathbf{g}$ where $\mathbf{g} \sim \mathcal{N}(0, I_d)$ has independent standard normal entries. Then

$$(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[i] \mathbf{x}_t[j] = \sigma_x^4 (\mathbf{w}^\top \mathbf{g})^2 g_i g_j = \sigma_x^4 \left(\sum_{k=1}^d w_k g_k \right)^2 g_i g_j.$$

Expanding the square and taking expectations,

$$\mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[i] \mathbf{x}_t[j]] = \sigma_x^4 \sum_{k, \ell=1}^d w_k w_\ell \mathbb{E}[g_k g_\ell g_i g_j].$$

By Isserlis' (Wick's) theorem for a zero-mean Gaussian vector,

$$\mathbb{E}[g_k g_\ell g_i g_j] = \delta_{k\ell} \delta_{ij} + \delta_{ki} \delta_{\ell j} + \delta_{kj} \delta_{\ell i},$$

where δ_{ab} is the Kronecker delta. Hence

$$\begin{aligned} \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[i] \mathbf{x}_t[j]] &= \sigma_x^4 \sum_{k, \ell=1}^d w_k w_\ell (\delta_{k\ell} \delta_{ij} + \delta_{ki} \delta_{\ell j} + \delta_{kj} \delta_{\ell i}) \\ &= \sigma_x^4 \left(\delta_{ij} \sum_{k=1}^d w_k^2 + w_i w_j + w_j w_i \right) \\ &= \sigma_x^4 \left(\delta_{ij} \|\mathbf{w}\|^2 + 2w_i w_j \right). \end{aligned}$$

Restricting to indices $i, j \in \mathcal{S}$ and collecting these entries into a $p \times p$ matrix yields

$$\mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] = \sigma_x^4 \|\mathbf{w}\|^2 \mathbb{I}_p + 2\sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top,$$

as claimed. \square

Lemma 16. *Under the setup described in Section 2, where $\mathcal{S} \subseteq [d]$ be any subset of indices of size p , $\mathbf{w} \in \mathbb{R}^d$ be a task vector, $\mathbf{x}_1, \dots, \mathbf{x}_T, \mathbf{x}_q \in \mathbb{R}^d$ be i.i.d. samples from $\mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbb{I}_d)$, labels $y_t = \mathbf{w}^\top \mathbf{x}_t + \epsilon_t$, where $\epsilon_t \sim \mathcal{N}(0, \sigma_n^2)$, and*

$$\hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}] := \frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t[\mathcal{S}] \in \mathbb{R}^p, \quad \mathbf{h}_{\mathcal{S}_p}(X) := \text{vec}(\mathbf{x}_q[\mathcal{S}] \hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}]^\top).$$

Then the following hold:

$$\begin{aligned} \underbrace{\mathbb{E}_{X, \epsilon}[\mathbf{h}_{\mathcal{S}_p}(X) \mathbf{h}_{\mathcal{S}_p}(X)^\top]}_{\text{Term-I}} &= \frac{\sigma_x^4}{T} [(\sigma_x^2 \|\mathbf{w}\|^2 + \sigma_n^2) \mathbb{I}_{p^2} + (T+1) \sigma_x^2 (\mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top) \otimes \mathbb{I}_p], \\ \underbrace{\mathbb{E}_{X, y}[\mathbf{h}_{\mathcal{S}_p}(X) y_q]}_{\text{Term-II}} &= \sigma_x^4 \mathbf{w}[\mathcal{S}] \otimes \mathbf{w}[\mathcal{S}]. \end{aligned}$$

Proof. We first simplify Term-I. We have that

$$\mathbb{E}_{X, \epsilon}[\mathbf{h}_{\mathcal{S}_p}(X) \mathbf{h}_{\mathcal{S}_p}(X)^\top] = \mathbb{E}_{X, \epsilon} [(\text{vec}(\mathbf{x}_q[\mathcal{S}] \hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}]^\top) (\text{vec}(\mathbf{x}_q[\mathcal{S}] \hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}]^\top)^\top)^\top].$$

Using Eq. (33), and since \mathbf{x}_q is independent of $\hat{\mathbf{w}}_{\text{avg}} = \frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t$, separating the expectation, we get

$$\mathbb{E}_{X, \epsilon}[\mathbf{h}_{\mathcal{S}_p}(X) \mathbf{h}_{\mathcal{S}_p}(X)^\top] = \frac{1}{T^2} \mathbb{E} \left[\underbrace{\left(\sum_{t=1}^T y_t \mathbf{x}_t[\mathcal{S}] \right) \left(\sum_{t'=1}^T y_{t'} \mathbf{x}_{t'}[\mathcal{S}] \right)^\top}_{\text{Term-I(i)}} \otimes \underbrace{\mathbb{E}[\mathbf{x}_q[\mathcal{S}] \mathbf{x}_q[\mathcal{S}]^\top]}_{\text{Term-I(ii)}} \right]$$

For Term-I(ii), since $\mathbf{x}_q[\mathcal{S}] \sim \mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbb{I}_p)$, we have $\mathbb{E}[\mathbf{x}_q[\mathcal{S}] \mathbf{x}_q[\mathcal{S}]^\top] = \sigma_x^2 \mathbb{I}_p$. Next, we simplify Term-I(i).

$$\begin{aligned} \mathbb{E} \left[\left(\sum_{t=1}^T y_t \mathbf{x}_t[\mathcal{S}] \right) \left(\sum_{t'=1}^T y_{t'} \mathbf{x}_{t'}[\mathcal{S}] \right)^\top \right] &= \sum_{t=1}^T \sum_{t'=1}^T \mathbb{E}[y_t y_{t'} \mathbf{x}_t[\mathcal{S}] \mathbf{x}_{t'}[\mathcal{S}]^\top] \\ &= \sum_{t=1}^T \mathbb{E}[y_t^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] + \sum_{t \neq t'} \mathbb{E}[y_t y_{t'} \mathbf{x}_t[\mathcal{S}] \mathbf{x}_{t'}[\mathcal{S}]^\top] \end{aligned}$$

2052 First, looking at $t \neq t'$ terms, we have:

$$\begin{aligned}
2053 & \mathbb{E}[y_t y_{t'} \mathbf{x}_t[\mathcal{S}] \mathbf{x}_{t'}[\mathcal{S}]^\top] = \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t + \epsilon_t) (\mathbf{w}^\top \mathbf{x}_{t'} + \epsilon_{t'}) \mathbf{x}_t[\mathcal{S}] \mathbf{x}_{t'}[\mathcal{S}]^\top] \\
2054 & = \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t) \mathbf{x}_t[\mathcal{S}]] \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_{t'}) \mathbf{x}_{t'}[\mathcal{S}]^\top] \\
2055 & = \sigma_x^4 \Pi_{p \times d} \mathbf{w} \mathbf{w}^\top \Pi_{p \times d}^\top \\
2056 & = \sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top,
\end{aligned}$$

2059 where the second equality follows by independence and $\mathbb{E}[\epsilon] = 0$, and for the third equality, we use the definition of $\Pi_{p \times d}$, and that $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbb{I})$.

2062 Next, consider the terms with $t = t'$:

$$\begin{aligned}
2063 & \mathbb{E}[y_t^2 \mathbf{x}_t \mathbf{x}_t^\top] = \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t + \epsilon_t)^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] \\
2064 & = \mathbb{E}[(\mathbf{w}^\top \mathbf{x}_t)^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] + \mathbb{E}[\epsilon_t^2 \mathbf{x}_t[\mathcal{S}] \mathbf{x}_t[\mathcal{S}]^\top] \\
2065 & = \sigma_x^4 \|\mathbf{w}\|^2 \mathbb{I}_p + 2\sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top + \sigma_n^2 \sigma_x^2 \mathbb{I}_p \\
2066 & = (\sigma_x^4 \|\mathbf{w}\|^2 + \sigma_x^2 \sigma_n^2) \mathbb{I}_p + 2\sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top,
\end{aligned}$$

2069 where the third equality follows by using Lemma 15 and that $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \sigma_x^2 \mathbb{I}_d)$.

2071 Combining the T diagonal terms and $T(T-1)$ off-diagonal terms, we have

$$\begin{aligned}
2072 & \mathbb{E} \left[\left(\sum_{t=1}^T y_t \mathbf{x}_t[\mathcal{S}] \right) \left(\sum_{t'=1}^T y_{t'} \mathbf{x}_{t'}[\mathcal{S}] \right)^\top \right] = T [(\sigma_x^4 \|\mathbf{w}\|^2 + \sigma_x^2 \sigma_n^2) \mathbb{I}_p + 2\sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top] + T(T-1) \sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top \\
2073 & = T(\sigma_x^4 \|\mathbf{w}\|^2 + \sigma_x^2 \sigma_n^2) \mathbb{I}_p + (2T + T^2 - T) \sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top \\
2074 & = T(\sigma_x^4 \|\mathbf{w}\|^2 + \sigma_x^2 \sigma_n^2) \mathbb{I}_p + (T^2 + T) \sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top.
\end{aligned}$$

2079 Combining this resultant Term-1(i) with Term-1(ii), we get Term-1:

$$\begin{aligned}
2081 & \mathbb{E}_{X, \epsilon} [\mathbf{h}_{\mathcal{S}_p}(X) \mathbf{h}_{\mathcal{S}_p}(X)^\top] = \frac{1}{T^2} [T(\sigma_x^4 \|\mathbf{w}\|^2 + \sigma_x^2 \sigma_n^2) \mathbb{I}_p + (T^2 + T) \sigma_x^4 \mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top] \otimes \sigma_x^2 \mathbb{I}_p \\
2082 & = \frac{\sigma_x^4}{T} [(\sigma_x^2 \|\mathbf{w}\|^2 + \sigma_n^2) \mathbb{I}_{p^2} + (T+1) \sigma_x^2 (\mathbf{w}[\mathcal{S}] \mathbf{w}[\mathcal{S}]^\top) \otimes \mathbb{I}_p].
\end{aligned}$$

2086 Next, let us simplify Term-II:

$$\begin{aligned}
2087 & \mathbb{E}_{X, y} [\mathbf{h}_{\mathcal{S}_p}(X) y_q] = \mathbb{E}_{X, \epsilon} [\text{vec}(\mathbf{x}_q[\mathcal{S}] \hat{\mathbf{w}}_{\text{avg}}[\mathcal{S}]^\top) (\mathbf{w}^\top \mathbf{x}_q + \epsilon)] \\
2088 & = \frac{1}{T} \mathbb{E} \left[\underbrace{\left[\sum_t y_t \mathbf{x}_t[\mathcal{S}] \right]}_{\text{Term-II(i)}} \otimes \underbrace{\mathbb{E} [\mathbf{x}_q[\mathcal{S}] \mathbf{x}_q^\top] \mathbf{w}}_{\text{Term-II(ii)}} \right].
\end{aligned}$$

2094 First, let us simplify Term-II(i). We have

$$\begin{aligned}
2095 & \mathbb{E} \left[\sum_t y_t \mathbf{x}_t[\mathcal{S}] \right] = \mathbb{E} \left[\sum_t (\mathbf{w}^\top \mathbf{x}_t + \epsilon) \mathbf{x}_t[\mathcal{S}] \right] \\
2096 & = \sum_t \mathbb{E} [\mathbf{x}_t[\mathcal{S}] \mathbf{x}_t^\top] \mathbf{w} = T \sigma_x^2 \Pi_{p \times d} \mathbf{w} = T \sigma_x^2 \mathbf{w}[\mathcal{S}].
\end{aligned}$$

2101 Similarly, Term-II(ii) is $\sigma_x^2 \mathbf{w}[\mathcal{S}]$. Combining these two, we get Term-II:

$$\mathbb{E}_{X, y} [\mathbf{h}_{\mathcal{S}_p}(X) y_q] = \sigma_x^4 \mathbf{w}[\mathcal{S}] \otimes \mathbf{w}[\mathcal{S}].$$

2105 \square

I ADDITIONAL RESULTS AND DETAILS OF EXPERIMENTAL SETTINGS

Comparison with Lu et al. (2025). Next, in Fig. 3, we isolate the setting studied by Lu et al. (2025) by fixing $\rho = 1$ (i.e., using all features) and vary the overparameterization ratio via $1/\nu$. We visualize the theoretical and empirical ID (left) and OOD (right) loss values. Crucially, in this setting, we do not observe in-context benign overfitting. While the OOD loss decreases as expected, the ID loss monotonically increases with overparameterization. This divergence stems from the nature of the scaling: varying $1/\nu$ requires altering the ambient dimension d relative to N , which fundamentally changes the underlying data distribution. In contrast, our feature-selection model (varying $\rho < 1$) allows us to scale the model capacity p while keeping the underlying data distribution fixed—a proxy that we argue more faithfully captures the effect of scaling model size. Consistent with our perspective, this parallels observations in supervised learning: standard linear regression without feature selection does not exhibit benign overfitting (Hastie et al., 2019).

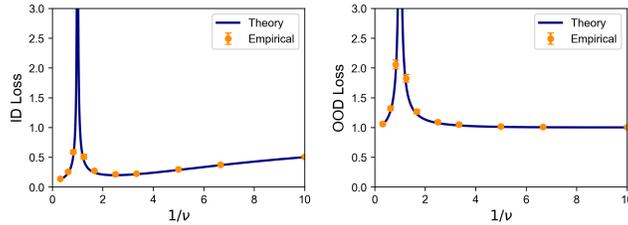


Figure 3: **Absence of In-Context Benign Overfitting.** Theoretical (lines) vs. empirical (points) ID (left) and OOD (right) risk as a function of overparameterization $1/\nu$ with fixed $\rho = 1$ (recovering the setting of Lu et al. (2025)). In contrast to Fig. 2, this setting exhibits a retrieval-learning tradeoff: as the model becomes more overparameterized, OOD generalization improves, but ID performance degrades (higher loss). This underscores that varying ambient dimension d is distinct from scaling model capacity p against a fixed data distribution.

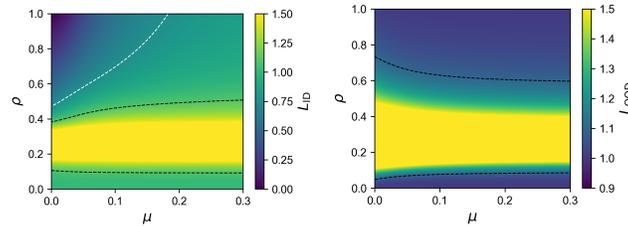


Figure 4: **Task Diversity vs. Model Scale.** Theoretical heatmaps for ID Loss (left) and OOD Loss (right) as a function of task diversity μ and model scale ρ . Lighter colors indicate higher loss. Vertical cuts (fixed μ) demonstrate in-context benign overfitting: larger models ($\rho \approx 1$) minimize ID loss without sacrificing OOD performance. Horizontal cuts (fixed ρ) reveal the task diversity threshold: increasing diversity forces a transition from memorization (low ID, high OOD) loss to learning-based ICL (high ID, low OOD risk).

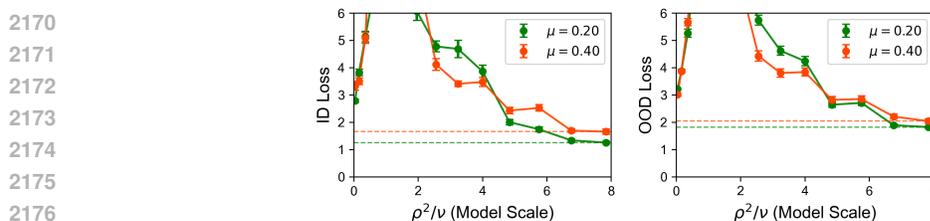
Impact of Task Diversity and Scale. Next, we investigate the role of varying the number of training tasks M , which we call task diversity μ and its interplay with model scale ρ . Figure 4 visualizes our theoretically derived ID (left) and OOD (right) risk landscapes as a joint function of model scale (ρ) and task diversity (μ). Our analysis recovers and clarifies the task-diversity thresholds identified in prior work (Raventos et al., 2023; Park et al., 2025):

- **Fixed Task Diversity (μ):** Consistent with Fig. 2, for any fixed number of pre-training tasks, increasing model scale (ρ) induces double descent in both ID and OOD risks. Consequently, sufficiently large models ($\rho \approx 1$) achieve lower ID loss than their underparameterized counterparts, and similar OOD loss confirming that in-context benign overfitting persists across different diversity regimes.
- **Fixed scale (ρ):** In contrast, for a fixed overparameterized model ($\rho > \sqrt{\nu}$), increasing task diversity (μ) presents a trade-off. Higher diversity improves OOD generalization but degrades ID performance, as memorizing more tasks becomes increasingly difficult for a fixed-capacity model.

2160 The critical point where ID performance begins to degrade while OOD performance improves
 2161 corresponds to the task diversity threshold established in previous literature (Raventos et al., 2023;
 2162 Park et al., 2025).

2163 **Power Law Covariance.** Thus far, we have assumed isotropic covariance matrices for Σ_x and Σ_w ,
 2164 which suffices to establish in-context benign overfitting. For completeness, in Fig. 5, we replace the
 2165 isotropic covariance matrices Σ_x and Σ_w with power-law decay (diagonal entries decaying with
 2166 exponent $\alpha = 1$; see Appendix I for details). In this setting, we observe a stronger form of in-context
 2167 benign overfitting: as scale increases, both ID and OOD losses converge to values strictly lower than
 2168 their respective minima in the underparameterized regime.

2169



2170

2171 Figure 5: **Power Law Covariance (Random Selection).** Experimental (Left) ID and (Right)
 2172 OOD loss curves where covariance eigenvalues (Σ_w, Σ_x) follow a power-law decay ($\alpha = 1$).
 2173 Using random feature selection, the overparameterized regime achieves strictly lower risk than the
 2174 underparameterized minimum for both ID and OOD tasks.

2175

2176

2177

2178

2179

2180

2181

2182

2183

2184

2185

2186

2187

2188

2189

2190

2191

2192

2193

2194

2195

2196

2197

2198

2199

2200

2201

2202

2203

2204

2205

2206

2207

2208

2209

2210

2211

2212

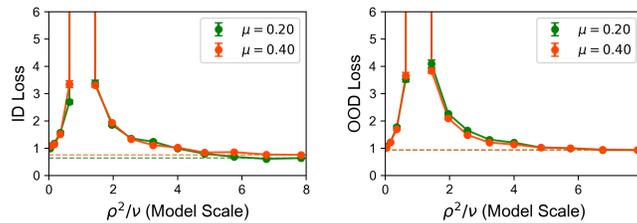
2213

2191 Figure 6: **Power Law Covariance (Top- p Features).** Experimental (Left) ID and (Right) OOD loss
 2192 curves using top- p feature selection. In contrast to random features, we do not observe in-context
 2193 benign overfitting here; the overparameterized performance is strictly worse than the underparameter-
 2194 ized minimum for both ID and OOD risks.

2196 In Fig. 6, instead of random p features, we consider top p features. In this case, consistent with
 2197 observations in the supervised learning setting (Belkin et al., 2020), we do not observe in-context
 2198 benign overfitting: As the model scale is increased, both ID and OOD loss values converge to a higher
 2199 value than in the underparameterized regime.

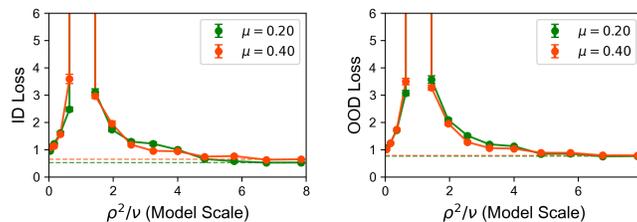
2200 **Settings.** For Fig. 2, we set $d = 80, T = 120, M = 10, n = 400$, and sweep over $p = 2, 6, \dots, 62$.
 2201 For Fig. 3, we set $d = p = 80, T = 120, M = 10$, and sweep over $n = 640, 960, 1280, \dots, 20480$.
 2202 For the experiments in Figs. 5 and 6 and those in this section, we set $d = 30, T = 40, M = 6, 12, n =$
 2203 100 , and sweep over $p = 2, 4, \dots, 28$.

2214
2215
2216
2217
2218
2219
2220
2221
2222



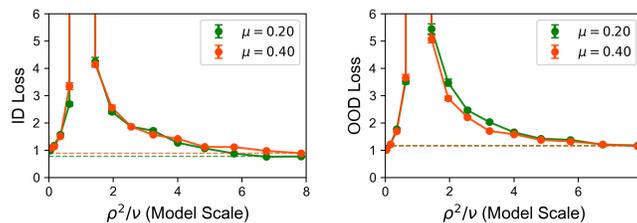
2223 **Figure 7: Power Law Covariance (Only Σ_w).** Experimental (Left) ID and (Right) OOD loss curves
2224 where task vector covariance eigenvalues follow a power-law decay ($\alpha = 1$). Using random feature
2225 selection, the overparameterized regime achieves strictly lower risk than the underparameterized
2226 minimum for ID tasks but not for OOD tasks.

2228
2229
2230
2231
2232
2233
2234
2235
2236



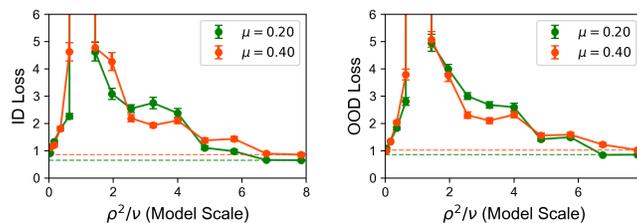
2237 **Figure 8: Power Law Covariance (Only Σ_w).** Experimental (Left) ID and (Right) OOD loss curves
2238 where task vector covariance eigenvalues follow a power-law decay ($\alpha = 1.5$). Using random feature
2239 selection, the overparameterized regime achieves strictly lower risk than the underparameterized
2240 minimum for both ID and OOD tasks.

2242
2243
2244
2245
2246
2247
2248
2249



2250 **Figure 9: Power Law Covariance (Only Σ_x).** Experimental (Left) ID and (Right) OOD loss curves
2251 where feature covariance eigenvalues follow a power-law decay ($\alpha = 1$). Using random feature
2252 selection, the overparameterized regime achieves strictly lower risk than the underparameterized
2253 minimum for ID tasks but not for OOD tasks.

2255
2256
2257
2258
2259
2260
2261
2262
2263



2264 **Figure 10: Power Law Covariance (Only Σ_x).** Experimental (Left) ID and (Right) OOD loss curves
2265 where feature covariance eigenvalues follow a power-law decay ($\alpha = 2.5$). Using random feature
2266 selection, the overparameterized regime achieves strictly lower risk than the underparameterized
2267 minimum for ID tasks and for lower μ for OOD tasks.