1   We thank all reviewers for their helpful feedback. Below we address the questions and comments individually.

2   **R1,R2,R3 - Writing.** We will correct typos in the main text and bibliography, and refer to Figure 1 in the introduction.

3   **R1 - Finite Data Analysis.** We remark that: (i) while our results are asymptotic in nature, our experiments showed
4   that they give accurate description of the prediction risk for dataset with moderate size, as demonstrated in Figure 2
5   ($n = 300$); (ii) we believe that similar to previous works on double descent, our asymptotic characterization can be
6   translated to non-asymptotic guarantees using standard concentration tools.

7   **R1 - Definition of "Same Order".** We apologize for the confusion. As we specify a joint relation between covariances,
8   we use the term "order" to describe the increasing or decreasing trend of a vector, i.e. "same order" of $a$ and $b$ implies
9   $a_i \geq a_j$ iff $b_i \geq b_j$ for all i,j. When $\Sigma_x, \Sigma_\beta$ are codiagonalizable, "same order" of eigenvalues suggests that the features
10  are informative (learning is "easy"); this is analogous to a fast eigenvalue decay of the RKHS "source condition".

11  **R1 - Extension to Neural Network.** We believe that our analysis can be extended to neural nets that are well-described
12  by a kernel or random features model (which is linear regression under different features). This includes a two-layer
13  network with fixed 1st layer (RF model), or with trained 1st layer under overparameterization (the NTK model).

14  **R2 - Novelty & Comparison with [1].** We thank the reviewer for mentioning this reference – we were not aware of
15  this work (which is arxived in May) at the time of submission. However, we believe that the reviewer has misjudged the
16  relation and overlap between [1] and our work. We first note that computing the prediction risk is only a small part of
17  our contribution; we also provided precise characterization of the optimal ridge parameter and weighting matrix. More
18  importantly, the setup and theory in [1] are quite different than and *do not* produce our results. In particular:

19  • **[1] may not cover the case when $\lambda \leq 0$.** When the regularization strength $\lambda$ is negative, it is not guaranteed that
20    limit of VAMP exists due to non-convexity of the objective in the overparameterized regime. Even if we assume
21    VAMP converges for all $\lambda$ and estimation of generalization error is accurate, it is still not clear that VAMP converges
22    to the specific ridge solution studied in our paper when $\lambda$ is negative or 0 (which corresponds to minimum $\ell_2$
23    norm solution), since it is possible that the SE has multiple fixed points when $\lambda \leq 0$. We remark that establishing
24    convergence and uniqueness can also be challenging for analysis based on the Convex Gordon Min-max Theorem.

25  • **The VAMP framework does not capture our "aligned" or "misaligned" cases.** The key assumption in the VAMP
26    analysis in [1] is that the limiting distribution of $\beta_i^\star$, the components of the true signal, is independent to that of the
27    features. As a result, random permutation of $\beta_i^\star$ does not effect the generalization error of the VAMP estimate. This
28    corresponds to the "random order" case in our setup, for which we showed that the corresponding optimal ridge
29    parameter is always non-negative[1]. In contrast, we specified a general joint distribution between the eigenvalues of
30    $\Sigma_x$ and the components of $d_\beta$ (See exact definition in paper). As demonstrated in Figure 2, different joint distribution
31    leads to different generalization error even when the limiting spectral distributions of $\Sigma_x$ and $\Sigma_\beta$ are the same. It is
32    precisely this extension that enables us to explain the "negative ridge" phenomenon.

33  • **The analysis in [1] alone cannot characterize the optimal ridge parameter and weighting matrix.** As mentioned
34    above, the optimal ridge parameter is always non-negative under the assumption in VAMP (it is also not true that SE
35    can always be simplified to exact expressions.). Furthermore, we decided the optimal weighting matrix for optimal $\lambda$
36    or the bias and variance term separately, and also covered the case when it only depends on feature covariance.

37  **R3 - Optimal Choice of Weighting Matrix.** We provide the following clarification and comments on the optimal $\Sigma_w$.

38  • The data covariance $\Sigma_x$ considered in Section 6 is the population covariance matrix, *not* the empirical covariance
39    (which is degenerate when $\gamma > 1$). While it can be difficult to determine $\Sigma_x$ from labeled data alone, the quantity
40    can be estimated from additional unlabeled data (i.e., a semi-supervised setting[2]). We agree with the reviewer that
41    designing a weighting matrix solely based on the empirical covariance is an interesting problem.

42  • We agree that interpolating between weighting matrices is an intuitive strategy: Figure 4 and 5 consider various
43    powers of data covariance ($\Sigma_x^\alpha$ for different $\alpha$), which can be seen as a geometric interpolation between $I_p$ and $\Sigma_x$.

44  • As mentioned in the paragraph starting from line 251, given the data covariance, a reasonable decision rule is to
45    construct the weighting matrix based on certain polynomial transformations to $\Sigma_x$, the parameter of which can be
46    tuned via cross validation. Figure 5 and 8 show that such approach does indeed outperform standard ridge regression.

47  **R3 - On Assumption 2.** We agree with the reviewer that assumption 2 does not cover all possible weighting matrices.
48  We remark that this assumption is *not* needed in our risk calculation, but only in the characterization of optimal
49  weighting matrix (for convenient formulas). In addition, it is worth noting that similar codiagonalizability assumptions
50  is not uncommon in the theoretical study of regression models. For instance, the standard source condition in RKHS
51  regression is analogous to codiagonalizable $\Sigma_x$ and $\Sigma_\beta$ with certain eigenvalue decay in our setup.

---

[1]This partially justifies the restriction to non-negative $\lambda$ in the VAMP framework.

[2]For example see: Tony Cai, T., and Zijian Guo. "Semisupervised inference for explained variance in high dimensional linear regression and its applications." Journal of the Royal Statistical Society: Series B.