1    We thank the reviewers for their feedback and the time spent on our submission.

2    **First,** let us elaborate on the concerns by **Reviewers 1 & 3** regarding restrictions to specific linear classifiers. To quote
3    **Reviewer 4** – who we thank for the encouraging feedback – : "*The theoretical analysis of multiclass classification is an*
4    *open problem at the core of machine learning/statistical modelling. While the specific setting considered seem limited,*
5    *they are insightful and likely an important stepping stone in the full analysis of multi-class classification...Such results*
6    *are a very good lead for future investigations of more general settings*". Below, we remark on the following regarding
7    the motivation behind the studied classifiers and the impact of our results. We will expand upon these in the revision.

8    **(I)** On the averaging estimator: The averaging estimator is Bayes optimal for balanced GMM with Gaussian means
9    (Prop. 3.4). As such, it serves well as a baseline and a "natural algorithm" for this data model. One then might wonder
10   how the performance of this algorithm depends on the data model. Thus, we further analyze its performance for
11   the second basic model considered here: the MLM. **(II)** On the LS classifier: There have been numerous empirical
12   works that investigate the role of the loss function in classification tasks for various data models. Several of these
13   find empirically that simple LS can have comparable performance compared to the (perhaps most commonly used)
14   hinge/logistic losses, e.g. [53,74,75]. Quoting [p.105, 53]: "*Intuitively, it seems that the square loss may be less*
15   *well suited to classification than the hinge loss (...) However, in practice, we have found that the accuracy of RLSC*
16   *(regularized LS classification) is essentially equivalent to that of SVMs*." One of the long-term goals of our project is to
17   provide theoretical evidence against/in-favor of such empirical findings and to characterize what loss is suitable for
18   each setting. As a first step, we naturally ask whether these claims are already justified (or not) in simple linear models,
19   and if so, under what conditions. Along these lines, there are several recent works that theoretically study the role of LS
20   in high-dimensional *binary* linear classification. For example, under the same asymptotic regime as in our paper, [44]
21   proves that LS is optimal for GMM within the family of convex un-regularized empirical-risk minimization, and, [60]
22   proves that LS is approximately optimal (thus, comparable to the ML solution: logistic loss) for logistic data. We take
23   the first steps towards extending these to the more challenging, but more versatile, *multiclass* setting. **(III)** On WLS: (a)
24   We are motivated by recent findings [13] that "weighted" variations of LS can significantly boost the performance over
25   simple LS. (b) Compared to LS, WLS offers the flexibility to adjust the algorithm to balance performance between
26   majority vs minority classes (together with our ability to accurately predict class-wise errors).

27   While there is a lot to do further down the road, our results (model setup, analysis, sharp asymptotics) are the first step
28   towards this direction and facing some of the new challenges in multiclass settings (see lines 45-52). Certain important
29   additions such as *regularized* (W)LS and correlated Gaussian features – while requiring extra work – are almost direct
30   extensions of our current framework. Others, such as the study of cross-entropy minimization or extreme classification
31   will likely require combining elements of our work with new ideas. We believe that our paper sets the fundamentals and
32   will inspire further investigations in this direction. Of course, extensions to non-asymptotic regimes and non-linear
33   models (e.g., RFF, NTK) are highly desirable. Such results, only recently obtained for regression settings, are typically
34   founded on long prior work on simpler regression models – linear, (isotropic) Gaussian, etc.. Our work, together with
35   refs. in (lines 81-82) for binary settings, resemble these essential precursor works for the setting of classification.

36   **Second,** on **Reviewer's 1** question on the relative performance of the averaging and LS estimators for the two data
37   models: This is discussed in Sec. 3.2 and 4.2. In Prop. 3.3, we prove that averaging outperforms LS for balanced GMM
38   with orthogonal means. Intuitively, this is because "compared to the weight vectors $\mathbf{w}_i$ of the class averaging classifier
39   that are also (asymptotically) orthogonal when means are orthogonal, this is not the case for LS" (line 222, pg.6). In
40   fact in Sec. 3.3 we formally study the optimality of the averaging estimator in a Bayesian GMM setting. Similarly, in
41   Prop. 4.3 and Sec. 4.2, we show that LS outperforms the averaging estimator in MLM for large data samples.

42   **Third,** we agree with **Reviewers 1 and 4** that sketching key proof ideas in the main body of the paper will benefit the
43   reader. If accepted, we will use the extra space to move the corresponding discussions from App. F to the main body.

44   **Fourth,** in response to **Reviewer's 4** suggestion. Indeed, results for binary classification can be lifted to characterize
45   the limit of $\mathbf{\Sigma}_{w,\mu}$ and diagonal entries of $\mathbf{\Sigma}_{w,w}$ for one-vs-all classifiers (including LS) (with some additional technical
46   work to capture correlations $\mathbf{\Sigma}_{\mu,\mu}$ for k>2). As mentioned in the paper, this alone does not give any information
47   on the off-diagonals of $\mathbf{\Sigma}_{w,w}$, needed in the exact test-error formulas (2.3)/(2.4). It is possible to derive heuristic
48   approximations and union bound arguments leading to error expressions that depend only on the diagonals of $\mathbf{\Sigma}_{w,w}$.
49   Indeed, Fig. 5 in App. A provides a result of this flavor and gives a sense of how our exact results improve upon such
50   approximations. We will expand upon this comparison in Appendix A in the revision.

51   **Reviewer 2:** With respect to test error, our paper is precisely about characterizing the performance of the studied linear
52   classifiers *in the sense of test error.* The formulas of Thms. 3.1,3.2,4.1,4.2 can be directly plugged in (2.3) and (2.4)
53   to obtain test error. Regarding "calculations in the double asymptotic regime are not quite new": Of course, there are
54   numerous works in this regime under numerous settings over the last decade (lines 77-79). However, ours is *the first*
55   *such work in multiclass classification* This point is well-articulated in the introduction (lines 73-76, 82-91).