Summary of the paper: The paper presents hdmax2, an R package designed to facilitate high-dimensional mediation analysis. The package builds upon the HDMAX2 framework, which integrates latent factor mixed models (LFMM) to estimate unobserved confounders and a max-squared test to identify significant mediators. This package represents a significant step forward in making complex mediation analysis more accessible and robust. A key strength of the package lies in its versatility. hdmax2 accommodates a variety of data types, including univariate and multivariate exposures and binary or continuous outcomes. This flexibility makes it a valuable tool for researchers analyzing high-throughput molecular data, such as DNA methylation or gene expression, where the number of mediators often far exceeds the number of samples. The paper showcases the package through two case studies:

1) Breast Cancer and HER2 Status: Explores the mediating role of DNA methylation in the pathway linking HER2-positive breast cancer status to a survival risk score.
2) Gender and Multiple Sclerosis (MS) Subtypes: Investigates gene expression as a potential mediator in the pathway linking gender to MS subtypes.

The package includes visualization tools, helper functions for mediator selection, and options for handling multivariate exposures. However, the paper focuses on univariate exposure models, leaving the multivariate capabilities underexplored.

My comments are primarily minor and aimed at enhancing clarity and providing additional context in a few areas.

Minor comment on Abstract:

The abstract summarizes the purpose and contributions of the study, emphasizing the development of a method that addresses statistical challenges in high-dimensional mediation analysis. However, it would benefit from explicitly connecting the features of the package to the case studies presented, particularly in explaining how the results demonstrate its utility.

Minor comment on Materials and Methods

The methodology is described in sufficient detail, including the use of LFMM for estimating latent confounders, the max-squared test for assessing mediators, and the R package's features. Figure 1 nicely illustrates the workarounds of the package. A significant limitation is the lack of discussion on the selection and interpretation of the number of latent factors K. For example, in the breast cancer case study, K=2 is mentioned without further justification, and for the MS case study, K is not mentioned at all.

Minor comment on Case study 1:

The breast cancer case study presents both the total effect (adjusted for age) and the mediation results. However, the total effect of 0.30 is introduced in a somewhat disconnected manner. If it is part of a preliminary analysis, this should be clarified, and its relevance to the HDMAX2 results (e.g., comparison to indirect effects in Figure 2C) should be explicitly discussed.

Minor comment on Case study 2:

It is worth questioning why the authors chose this case study. While the initial motivation to assess the relationship between gender and MS is clear and compelling, and the discussion of negative results can be valuable, the authors acknowledge that the lack of

significant findings is likely due to the small sample size. They proceed to discuss the top-ranked, albeit not statistically significant, results and provide reasoning for their potential biological relevance. If the authors believe that limited statistical power might be a common challenge when applying the hdmax2 model, it would be beneficial to include a power analysis in the paper to better address such scenarios.

Additional minor comments:

1) Is this sentence accurate: "Upon observing a significant decrease in CIS-MS occurrence among women (see Fig 3A), we sought to investigate this phenomenon further"? (line 187-188) Did the authors intend to refer to RR-MS instead of CIS-MS?

2) The statement, "Remarkably, most of the top 10 identified mediators were associated with genes known to be involved in breast cancer biology, thus supporting the biological relevance of our approach," is likely accurate. However, it appears to rely on "PubMed hits," defined as the number of outputs from the search "(Breast cancer) AND ('Gene Symbol')." I recommend that the authors either explicitly mention in the text the methodology used to link genes and breast cancer or refine their search methodology to provide stronger evidence supporting this claim.