

1 We would like to thank reviewers for their time and constructive comments. Following are our responses to the raised
2 questions. We start with common questions followed by individual ones.

3 **Common Q1:** Both **R2** and **R3** question about the claim of the tradeoff between personalization and accuracy. We
4 empirically validated the existence of such tradeoff. For example, by tracking the performance of users in the validation
5 set, we find that the results follow the “80-20 rule”, i.e., 80% users whose RMSE results become worse are within top
6 22.5% over-personalized users. Following reviewers’ suggestion, we will include this result into our revised version.

7 **Common Q2:** Both **R2** and **R3** also question about the applicability of the proposed method. Indeed, we agree with **R2**
8 that the proposed framework is quite generic, and we would like to clarify that the proposed method does not necessarily
9 require a closed form solution such as in matrix factorization. In fact, the proposed framework can be used in many
10 collaborative filtering models as long as the following two theoretical conditions are satisfied. First, the $\nabla_{\Theta} L_{\Gamma}$ term in
11 Eq. (5) can be optimized to (or near) 0 (i.e., the KKT condition). Second, the second partial derivatives of L_{Γ} exist so
12 that we can use the the implicit function existence theorem to obtain Eq. (7). Both conditions can be easily satisfied
13 in many collaborative filtering models including the neural network models mentioned by the reviewers. We plan to
14 extend our idea to such models in our future work.

15 **Q1 from R1:** **R1** has a concern about overfitting the validation set and the inclusion of a test set. We would like to
16 clarify that this is a misunderstanding and we do have a test set which is separated from both the training set and the
17 validation set. To be specific, as stated in Section 4.1, “we randomly select 80% ratings for training and use the rest
18 20% ratings for testing”. For the training set, we then split it into several folds of debug set and validation set. After the
19 CFDebug finishes training on the training set as illustrated in Fig. 1, we evaluate its performance on the test set. Finally,
20 we repeat the above process five times to do cross-validation. We will make it clearer in the revised version.

21 **Q2 from R1:** **R1** also concerns about the technical merit of the proposed method. We would like to point out that the
22 key contributions of this work are two fold. First, the data debugging framework that is potentially applicable for a large
23 set of collaborative filtering models (see the response for **Common Q2**). Second, experimental findings and analysis
24 about the tradeoff between personalization and accuracy (e.g., modifying the over-personalized ratings would help
25 improve the overall accuracy).

26 **Other questions from R1:** Thanks for the advice, we will add more explanations about the C matrix. In Section 4 (D),
27 we did display the movie titles in the case study. We will make it clearer in the revised version.

28 **Q1 from R2:** **R2** encourages us to experiment on larger datasets. Thank you for the suggestion. We have done
29 experiments on ML-10M. The results also show the effectiveness of the proposed method. For example, with 10%
30 modified ratings, our method gives 2.1% RMSE improvement compared to the original performance (from 0.8336 to
31 0.8162). We will include these results into our future version.

32 **Other questions from R2:** As suggested by the reviewer, we will include the top-N studies in the main paper instead
33 of the appendix, and will consider to include the comparison with [7] in our journal version.

34 **Q1 from R3:** **R3** also asks about the optimal ratio of the results in Table 1 and 2. Empirically, the optimal ratio is
35 around 5%-10% for the studied datasets. We will make this clearer in our revised version.

36 **Q1 from R4.** **R4** encourages us to analyze the limitation of the proposed method. Thanks for pointing this out. For
37 example, when the ratio of “over-personalized” ratings increases to a certain extent, the overall performance might start
38 to suffer. Much more theoretic work is needed to understand or identify such a theoretic ‘transition’ point.

39 **Q2 from R4:** **R4** also suggests using matrix completion view to provide theoretical guarantees. We agree that matrix
40 completion has the potential to provide additional theoretical guarantee.