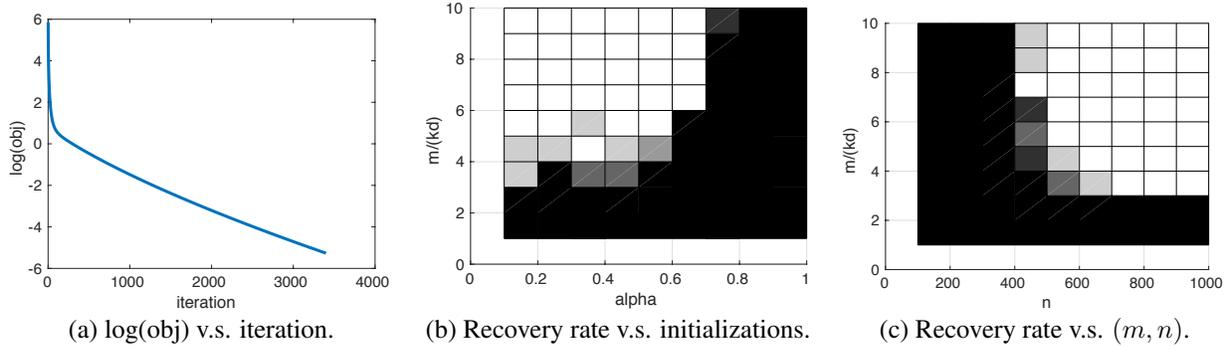


1 We thank all the reviewers for their constructive feedback!

2 **To Reviewer #1.** 1. We agree (and will acknowledge more explicitly) that the overall proof program is similar to
 3 existing results in the area. However, NIMC problem presents two key challenges: a) The Hessian has entangled terms
 4 for items' features and queries' features, which are challenging to handle. 2) In addition to non-convexity arising due to
 5 non-linearity of the activation function which standard 1-2 hidden layer NNs also face, we have to handle additional
 6 noise/uncertainty due to missing ratings, and provide strong sample complexity bounds for the results to be meaningful.

7 2. We'll reduce section 3.4 to a short sentence.

8 3. Here we provide more experimental results in Fig. 1. We use sigmoid as the activation function, and set $k = 10, d =$
 9 100 , which are larger than those in the paper ($k = 5, d = 10$). We set the initialization as $W^{(0)} = (1 - \alpha)W^* + \alpha W^{(r)}$,
 10 where W^* is the ground truth, $W^{(r)}$ is a Gaussian random matrix, and $\alpha \in [0, 1]$. In (a), $\alpha = 0.1, n = 1000$, and
 11 $m = 10000$. In (b), $n = 500$. In (c), $\alpha = 0.1$. The other settings are same as those in the paper. As we can see,
 12 (a) shows how the objective value converges, which is almost linear. (b) shows that when the initialization is purely
 13 random ($\alpha = 1$), gradient descent doesn't converge to the ground truth. In the paper, when $k = 5, d = 10$, pure random
 14 initialization still converges to the ground truth. We believe that it is because when k, d are larger, random initialization
 15 can be further away from the ground truth. Hence, gradient descent can get stuck in local optima more easily. Finally,
 16 comparing (c) with Fig. 1(a) of the paper, we can obtain a similar conclusion, i.e., when n is sufficiently large, the
 17 number of observed ratings required for successful recovery remains the same.



18 4. To remove the redundancy in the ReLU case, we assume that $u_{1,i}^*$ is nonzero for all $i \in [k]$ and know the number of
 19 positives in $\{u_{1,i}^*\}_{i=1,\dots,k}$. Note that if the columns of U and the columns of V do the same permutation, the output
 20 doesn't change. Without loss of generality, we can assume $\{u_{1,i}^*\}_{i=1,\dots,k_+}$ ($0 \leq k_+ \leq k$) are positive and the
 21 remaining $\{u_{1,i}^*\}_{i=k_++1,\dots,k}$ are negative. So if we fix $u_{1,i} = 1$ for all $i \leq k_+$ and $u_{1,i} = -1$ for all $i > k_+$, we can
 22 remove the redundancy and the target solutions for U and V are $u_{:,i} = u_{:,i}^*/|u_{1,i}^*|$ and $v_{:,i} = v_{:,i}^*/|u_{1,i}^*|$ respectively.

23 **To Reviewer #2.** We will like to stress that Non-linear Inductive Matrix Completion is a significantly different
 24 architecture than the 1-layer NNs and hence theoretical analyses for the two are quite different. As mentioned in
 25 response to R1, while the high level approach is same, we have to deal with non-linearity of NNs along with the noise
 26 due to missing ratings and entangled Hessian due to non-linearity in both item's features and query's features. These
 27 challenges require a significantly different analysis than existing results.

28 **To Reviewer #3.** Movielens dataset: our main goal in these experiments is to study the problem in the *inductive*
 29 setting, i.e., to predict ratings for *new* users. R3's observations for collaborative filtering (CF) are valid but they apply
 30 only to *transductive* setting which does not allow for new users.

31 a) SVD based solution: SVD based CF does not predict ratings for new users and hence does not apply in the inductive
 32 setting. Furthermore, as we are predicting ratings for completely new users, for which only weak features are available,
 33 naturally the resulting RMSE is worse than the results for the standard collaborative filtering settings (where several
 34 ratings of a user are available a priori).

35 b) Generalization error: as mentioned above, the only information about new users is their relatively weak features,
 36 hence non-linear methods can extract more information from them compared to the linear ones, and might be the reason
 37 for the superior performance of NIMC over IMC.