1   We thank the reviewers for their comments.

2   **Reviewer 1.** In the final draft, we will edit the exposition to make it friendlier to non-expert readers. Specifically,
3   we will correct the reference to matrix $A$ at the end of the introduction (line 68) and change it to matrix $M$. We will
4   also move the definitions of $D_{W_{i,:}}$ and $D_{W_{:,j}}$ to a point before their first use in line 89. We will provide the standard
5   definition for *non-negative rank*.

6   Like [RSW16], our setting involves a more general family of weight matrices $W$ than just binary matrices. Since
7   our proofs dealt with weight matrices that were not necessarily binary matrices, we wanted our experiments to use
8   non-binary weight matrices to highlight the fact that we weren't just studying matrix completion. We have additional
9   experiments involving the NIPS and synthetic datasets that use binary weight matrices which turn out similarly to our
10  current experiments. For the final draft, we would be happy to add these experiments to the appendix, as well as some
11  additional experiments with varying regularization parameter values.

12  **Reviewer 2.** We will address speed and SVD-related issues in the comments to Reviewer 3.

13  We believe the contribution of this work over [RSW16] is more than incremental because even though the algorithmic
14  steps may be similar, the proof techniques required are quite different. Most provable sketching results for Low Rank
15  Approximation (LRA) problems do not have sketch sizes that can be significantly smaller than the rank. The small
16  sketch size means imitating the analysis of [RSW16] is insufficient because when one solves a regression problem on a
17  matrix with fewer rows than columns one always gets 0. Furthermore, the proof was achieved without the incoherence
18  assumptions on the input matrix that are popular in the matrix completion literature. Thus, we needed a finer analysis
19  based on condition numbers and tail bounds on the singular vectors because directly following the approach of [RSW16]
20  will fail to give sharp enough inequalities. We elaborate on this in lines 101 to 108.

21  We also improve the results of [RSW16] by providing fast $2^{\mathrm{poly}(r \cdot sd)}$ algorithms (as opposed to $n^{\mathrm{poly}(r \cdot sd)}$) in the case
22  when the ratio of the largest to smallest entries of the weight matrix is controlled and the largest singular value of our
23  regression matrices is small relative to lambda. We achieve this by replacing the Cramer's rule-based approach in
24  [RSW16] with different techniques from optimization like Richardson's Iteration.

25  **Reviewer 3.** The primary focus of our work is on the theoretical side rather than the experimental side. We would like
26  to reiterate that our algorithm has a greatly improved running time when compared to other $(1 + \epsilon)$-approximation
27  algorithms for our regularized, weighted setting. The theoretical running time is not being compared to that of singular
28  value decomposition because SVD is not a $(1 + \epsilon)$-approximation algorithm for the regularized, weighted setting.

29  We included SVD in our experiments because it is widely used in practice for LRA type problems and to demonstrate
30  that it results in high objective values for the loss function. Given the significantly higher objective values and the fact
31  that SVD does not provide a $(1 + \epsilon)$-approximation algorithm for our problem (because it is NP-complete) we did
32  not think it was appropriate to compare its speed with our algorithm. However, we did think it was fair to compare
33  the sketched version of our algorithm to an *unsketched* version of our algorithm. As described in line 299, we ran
34  experiments that showed that alternating minimization with sketching was between $1.43$ and $2$ times as fast as alternating
35  minimization without sketching. We can add a table to the final draft.

36  We also wanted to emphasize in line 258 that the purpose of the experiments was to show that even if one sketches
37  down to the statistical dimension, which can potentially be much lower than the rank of the matrix, it is possible to do
38  this without blowing up the objective value in regularized weighted LRA. While the focus on the theoretical side of the
39  paper was on the running time, the focus on the experimental side of the paper was on the dimension reduction.

40  This is because this algorithm and the algorithm in [RSW16] which we improve on both use polynomial system solvers
41  which are costly. In fact, the implementation in [RSW16] could barely handle target ranks and sketching dimensions
42  larger than 2 and matrix dimensions larger than $100$. Our experiments involve target ranks of at least $50$ and matrix
43  dimensions in the 1000s. Thus, we feel that the dataset sizes show a marked improvement over the prior work but we
44  can include an even larger dataset for the final draft.

45  Although they are costly, polynomial system solvers do have provable theoretical guarantees which is why we invoked
46  them in the theoretical part of our paper. In practice, heuristics like alternating minimization are often faster but they
47  lack provable theoretical guarantees without making assumptions on the input, which we do not. Since our experiments
48  were not using the polynomial system solver technique described in the theoretical sections of our papers but they were
49  using the same sketches, we decided to focus the experiments on dimension reduction.

50  We obtained our $\lambda$ value by hand-tuning and felt that it was in a sweet spot that avoided underfitting and overfitting. We
51  can add experiments with varying regularization parameter values, or $\lambda$ values tuned by cross-validation, in the final
52  draft.