

1 We want to thank the reviewers for their careful reading of the paper and their many constructive comments. Due to
 2 space constraints, we are not able to reply to every point made by the reviewers and focus our response on the most
 3 important points. Naturally, all of our responses here will be incorporated into the final version of the paper.

4 **Additional baselines and the two-coin model:** We would like to thank Reviewer 3 for suggesting several highly
 5 relevant references with additional algorithms. We therefore compared the M-MSR method with 8 new baseline
 6 methods, where 6 of them are two-coin methods (Minmax¹, Entropy(O)², EM-twocoin³, BP-twocoin³, MFA-twocoin³,
 7 and M-MSR-twocoin) and 2 of them are tensor methods (EM-MV⁴, EM-SM⁴).

8 Additionally, all reviewers suggest we should test our method on the two-coin model: here the labels are binary and the
 9 probability of giving the correct answer depends on the label, so that each worker is parametrized by two parameters.
 10 Besides the above methods and M-MSR, we will also simulate a natural extension of M-MSR to the rank-2 case:

$$x_i^{t+1} = \arg \min_x \sum_j |_{(i,j) \in \Omega} (x^T y_j - C_{ij})^2; \quad y_j^{t+1} = \arg \min_y \sum_i |_{(i,j) \in \Omega} (x_i^T y - C_{ij})^2,$$

11 where $x_i \in \mathbb{R}^2$ denotes the i th row of X and $y_j \in \mathbb{R}^2$ denotes the j th row of Y ; to deal with the corrupted entries, we
 12 throw out the largest and smallest F values among the quantities $|C_{ij}| / \|y_j\|_2$ inside the minimization problem in each
 13 step. We call this algorithm the M-MSR-twocoin algorithm. It is not hard to see that, if the true proportion of “0” and
 14 “1” labels in the data-set is known (more on this later), skill estimation in the two-coin can be reduced to a rank-2 matrix
 15 recovery problem to which this algorithm can be applied.

16 We had time to implement the same experiments as in the paper on 5 real datasets and one synthetic dataset. The
 17 synthetic dataset is created following the two-coin model. The results are shown in Figure 1. *We can see in the figure
 18 that M-MSR outperforms all the baseline methods on all 6 datasets. Surprisingly, M-MSR even beats its own natural
 19 generalization to rank-2, even though this generalization has the implicit advantage of knowing the true proportion of
 20 “0” and “1” answers; and it can also beat all the methods tailored to the two-coin model on the synthetic data set once
 21 adversaries are introduced.*

22 We have spent substantial efforts during the rebuttal period implementing these new baselines and we believe the
 23 results considerably strengthen the case for the effectiveness of M-MSR. While it may not be surprising that M-MSR
 24 outperforms various two-coin methods (since M-MSR is designed with adversaries in mind) it is quite surprising that
 25 M-MSR outperforms its rank-2 generalization.

26 **Reviewer 2: More discussion for the consensus literature.** Indeed, the reviewer is correct. We will move some of
 27 Section 5.1 to the supplementary information which will give us the space to discuss this connection.

28 **Reviewer 3: Reproducibility.** We will publish all the code with the final version paper on github.

29 **Differences between our work and literature⁵, literature⁶.** One difference from literature⁵ is that we aim to label
 30 the class of the tasks, and they aim to rate the tasks (not necessarily integer) and find the ones with high quality. The
 31 main difference, however, is that they require the ratings of a small portion of tasks are known. Literature⁶ and our work
 32 both study the issue of worker ability estimation, however, we mainly focus on the adversarial setting while they do not.

33 **Use of “gold standard” questions.** One advantage of M-MSR is that it can tolerate arbitrary adversaries, including
 34 adversaries that return large fraction of correct answers (this can be seen in Figure 1(g) in the paper), which means
 35 M-MSR can deal even with smart adversaries that manage to pass the gold question test.

36 **Experimental results using raw datasets.** Among the 17 real datasets we used, 11 of them had each worker with 10
 37 tasks or more, so no removal of workers was necessary. For the other 6 datasets, we implemented experiments on 4 of
 38 them (Dog, Adult, Fashion1, and Fashion2) using raw data without removals; this did not change the results. The results
 39 in Figure 1 below also use raw data without removals. Experiments on remaining 2 raw datasets (Web and TREC) will
 40 be added in final version of the paper (Web is a multiclass data set and some of the two-coin methods cannot deal with
 41 that; TREC is very large and some of the methods take a very long time to run on it).

42 **Reviewer 4: Notation.** Note that we consider two problems: recovery of an $m \times n$ rank-1 matrix and crowdsourcing
 43 with W workers and M classes. We reduce the crowdsourcing problem to rank-1 recovery of a $W \times W$ matrix.

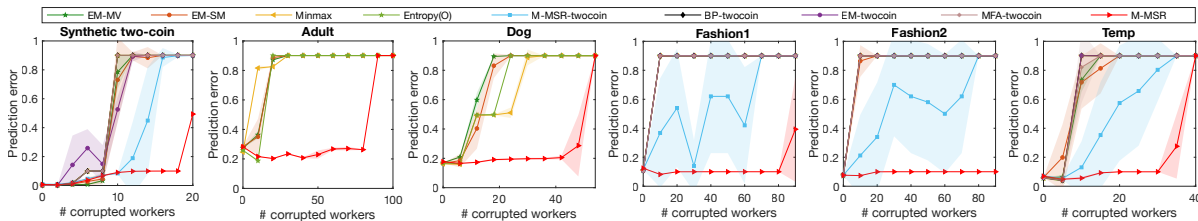


Figure 1: Experimental results. Solid lines are the average of 20 runs, and the shaded region show the standard deviation. The total number of workers of the datasets: Adult 269, Dog 109, Fashion1 196, Fashion2 198, Temp 76, Synthetic two-coin 40. Adult and Dog are multiclass datasets, hence the methods which are designed only for binary datasets are not shown. The synthetic dataset is created following two-coin model rule, where the worker probabilities of giving correct answers are uniform in $[0.5, 1]$ for both classes, and the true answers have equal probabilities of coming from each class.

¹ Zhou et al. Learning from the Wisdom of Crowds... ² Zhou et al. Aggregating Ordinal Labels... ³ Liu et al. Variational Inference for Crowdsourcing ⁴ Zhang et al. Spectral Methods meet EM... ⁵ Steinhardt et al. Avoiding Imposters... ⁶ Zhou et al. MultiC²: An Optimization framework...