

Appendix II: Supplementary Results

YANG LIU, University of Helsinki, Finland

ALAN MEDLAR, University of Helsinki, Finland

DOROTA GŁOWACKA, University of Helsinki, Finland

The document is under construction.

1 BETWEEN-REPLICATE CONSISTENCY

We reported the average results over ten replicates for each experiment in the original paper, in order to improve the reliability of the analysis. In this section, we provide further details related to the standard deviation. To recap, we performed four experiments to understand the impact of different sampling conditions on *model ranking consistency*, *discriminative power*, *sparsity bias* and *popularity bias*, respectively.

Table 1. Mean and Maximum Standard Deviation for All Experiments.

Dataset	Mean Corr.	Std. Max.	Std. Mean
MovieLens-100K	0.829	0.018	0.008
MoveiLens-1M	0.735	0.014	0.005
Amazon-Books	0.74	0.007	0.003
(a) Experiment I: model ranking consistency			
Dataset	Mean Corr.	Std. Max.	Std. Mean
MovieLens-100K	0.929	0.066	0.012
MoveiLens-1M	0.942	0.032	0.009
Amazon-Books	0.986	0.017	0.004
(c) Experiment III: sparsity bias			
Dataset	Mean #Ties	Std. Max.	Std. Mean
MovieLens-100K	330.887	3.885	1.145
MoveiLens-1M	2031.838	12.097	3.633
Amazon-Books	52557.728	30.43	11.002
(b) Experiment II: discriminative power			
Dataset	Mean Corr.	Std. Max.	Std. Mean
MovieLens-100K	0.72	0.029	0.009
MoveiLens-1M	0.699	0.015	0.006
Amazon-Books	0.734	0.012	0.004
(d) Experiment IV: popularity bias			

Table 1 presents the mean and maximum standard deviations observed in these experiments. We have also included an additional column in each subtable, showing the mean correlation (denoted as Mean Corr.) or mean number of ties (denoted as Mean #Ties) calculated across all sampling conditions, facilitating readers to assess the magnitude of observed standard deviations. As evident from each subtable, the highest deviations are noticeably small when compared to the mean metric values, i.e. mean correlation or mean number of ties, across three datasets.