

1 We thank the reviewers for their positive and constructive feedback. We will address all minor suggestions in the final
2 version of the paper. We now reply to the main points raised by the individual reviewers.

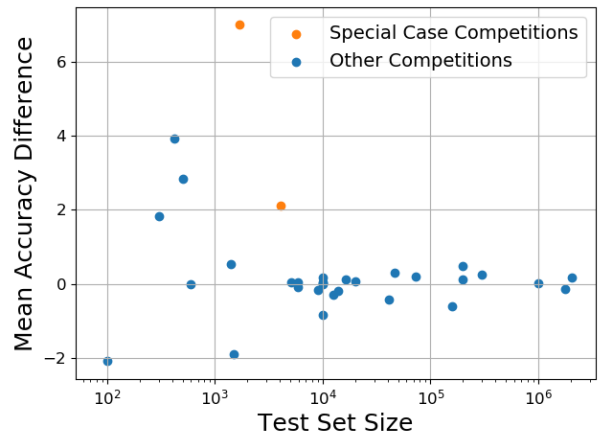
3 Reviewer 1

4 **Point 1:** We agree with the overall conclusion of the reviewer, with the exception that the non-uniform p-values in
5 Section 3.3 are only limited evidence for adaptive overfitting due to the strong null hypothesis. Nevertheless, some of
6 the plots in Section 3.2 also indicate small differences between public and private accuracy scores that are indicative of
7 mild overfitting. We will update our discussion in Section 3.5 to more clearly reflect this.

8 **Point 2:** We will update our figures with alpha values for the scatter plots. Figure 1 does show confidence intervals, but
9 they are small due to the scale of the axes and the test set sizes. We will clarify this in the figure captions.

10 **Point 3:** Figure 4 of our submitted paper already points to-
11 wards one covariate that is correlated with adaptive overfit-
12 ting: in-class vs. other competitions (mainly featured competi-
13 tions). This indicates that quality control performed
14 by Kaggle for the featured competitions may avoid certain
15 causes of public vs. private accuracy differences.

16 Moreover, we will include the figure on the right in the fi-
17 nal version of our paper. The figure shows the relationship
18 between overfitting and test set size for accuracy competi-
19 tions. It also points towards 10,000 examples as a possible
20 recommendation for minimum test set size. While a reli-
21 able recommendation for applied machine learning will
22 require a broader investigation that goes beyond the scope
23 of our current paper, we believe that our candidate hypoth-
24 esis can still inform future work on adaptive overfitting.



25 **Point 4:** While we agree that ranking submissions is an important aspect of the Kaggle platform, we believe that score
26 differences are more fundamental from the perspective of adaptive overfitting. If many models with essentially the
27 same performance are submitted to a competition, small variations can lead to large rank changes. However, these rank
28 changes are not problematic if the performance of each method stays the same up to statistical fluctuations. Hence we
29 view ranking mainly as an artifact of the competition format and not inherently important for safely deploying machine
30 learning in real applications. Having said that, we have compared our results to those of the “Meta Kaggle: Competition
31 Shake-up” notebook on Kaggle and found similar conclusions when considering changes in ranking. We will comment
32 on this in the final version of our paper.

33 Reviewer 2

34 Regarding the relationship between competition type or test set size and overfitting, please see Point 3 in our response
35 to Reviewer 1.

36 We certainly plan to investigate further competition platforms and machine learning tasks (including regression). In
37 fact, we have already contacted 20 individual ML competitions hosted outside Kaggle in order to gather data.

38 Regarding the ICML’19 workshop submission mentioned by Reviewer 2: please note that the corresponding workshop
39 submission deadline was after the NeurIPS submission deadline. Due to the NeurIPS guidelines for anonymous
40 submissions, we cannot comment on this point beyond stating that all applicable rules are being followed to the best of
41 our knowledge.

42 Reviewer 3

43 We agree with the reviewer that the non-i.i.d. case (i.e., distribution shifts) is crucial for deploying machine learning and
44 plan to investigate this point in future work. But as stated by the reviewer, the purpose of the current paper is the i.i.d.
45 case since it allows us to clearly separate adaptive overfitting from distribution shift and show that adaptive overfitting
46 is currently not substantively affecting machine learning practice.

47 Regarding lines 275–278: with “not enough data”, we mean that the MetaKaggle dataset only provides the aggregate
48 score for each submission and not the loss on each example in the test set. This makes it impossible to use techniques
49 such as bootstrapping to compute error bounds since we cannot re-sample the test set. In contrast, the accuracy metric
50 is simple enough so that the aggregate score is a sufficient statistic to compute tight (even exact) confidence intervals
51 (Clopper-Pearson confidence intervals).