

1 We thank all of the reviewers for their thoughtful feedback, and will incorporate their suggestions into the next version
2 of our paper. We detail our responses to their comments below.

3 **R1.** We thank R1 for their comments and will emphasize the broader implications of our work on model explainability.

4 **R2.** R2 asked to contrast using (i) influence functions to measure the importance of training points with (ii) existing
5 techniques for measuring feature importance, namely Datta, Sen, & Zick, 2016; Adler et al., 2016; and Adebayo &
6 Kagal, 2016. These papers address a different problem setting from ours and their methods are correspondingly distinct.

7 The main difference is that the papers above seek to explain a *fixed* model θ , whereas we examine how the *learned*
8 model $\hat{\theta}$ changes as a function of its training data. Our central issue is therefore reasoning about retraining the model,
9 which is not a concern shared by the papers above. Concretely, those papers consider the question: given a fixed model
10 θ , how do its predictions on a test set $\mathcal{D}_{\text{test}}$ depend on the values of some feature x_k in $\mathcal{D}_{\text{test}}$? They investigate this by
11 perturbing the value of x_k in various ways (e.g., by randomizing x_k for each test point x in $\mathcal{D}_{\text{test}}$). In contrast, we are
12 given the training set $\mathcal{D}_{\text{train}}$ and our goal is determining the effect of removing groups of points in $\mathcal{D}_{\text{train}}$ on the learned $\hat{\theta}$.

13 Despite their differences, these methods could be complementary, as R2 suggested. For instance, if we find (with feature
14 importance methods) that a model depends heavily on some feature, we could use influence functions to identify the
15 training data most responsible for that dependence. We will include this discussion and we thank R2 for pointing it out.

16 **R3: Non-convex objectives.** R3 asked if our empirical findings hold for non-convex models. Our initial experiments
17 are promising and suggest that this can be true; we will discuss this question in our next revision and plan to conduct a
18 more extensive study. We are grateful to R3 for highlighting this question. To properly respond, let us first provide
19 context on why influence functions and actual effects are classically only defined for convex models.

20 Recall that the actual effect $\mathcal{I}_f^*(w)$ of a subset w measures the difference between (i) the original model $\hat{\theta}(\mathbf{1})$, which
21 minimizes the loss on the training data $\mathcal{D}_{\text{train}}$, and (ii) the new model $\hat{\theta}(\mathbf{1} - w)$, which minimizes the loss on $\mathcal{D}_{\text{train}}$ with
22 w removed. Thus, for $\mathcal{I}_f^*(w)$ to be well-defined, there must be a unique model $\hat{\theta}(\mathbf{1})$ that globally minimizes the training
23 loss on $\mathcal{D}_{\text{train}}$, and likewise for $\hat{\theta}(\mathbf{1} - w)$. This condition is satisfied when the model is strongly convex. Similarly,
24 the influence $\mathcal{I}_f(w)$ is only well-defined when $\hat{\theta}(\mathbf{1})$ is unique and the model is strongly convex around it. Finally, to
25 measure the actual effect and influence, the models $\hat{\theta}(\mathbf{1})$ and $\hat{\theta}(\mathbf{1} - w)$ must not only be well-defined but computable.

26 Non-convex models unfortunately violate all of these requirements: the global minimizer $\hat{\theta}(\mathbf{1})$ may not be unique for a
27 given $\mathcal{D}_{\text{train}}$, and even if it were, we may not find it. For instance, neural networks are typically trained with SGD-based
28 methods that only guarantee convergence to a local minimum, so in general we cannot compute $\hat{\theta}(\mathbf{1})$ nor $\hat{\theta}(\mathbf{1} - w)$.

29 To address these issues, we propose augmenting the classical definitions as follows. Let the actual effect $\mathcal{I}_f^*(w, \theta_0, r)$ of
30 a subset w given an initial trained model θ_0 and a random seed r to be the change in the model after removing w and
31 retraining the model by starting from θ_0 and running SGD with randomness r . Specifying the initial model θ_0 sidesteps
32 the issue of $\hat{\theta}(\mathbf{1})$ being non-unique or impossible to compute, while specifying r resolves the issue of the retrained
33 model $\hat{\theta}(\mathbf{1} - w)$ being ill-defined. Similarly, we augment the predicted effect $\mathcal{I}_f(w, \theta_0)$ to be the influence of w around
34 the initial model θ_0 . Our goal is then to see if $\mathcal{I}_f(w, \theta_0) \approx \mathcal{I}_f^*(w, \theta_0, r)$ for all r .

35 We tested this with a multi-layer perceptron (MLP) with 2 hidden layers (128 and 32
36 nodes) on 10% of the MNIST 10-class dataset, choosing θ_0 as a model trained from
37 scratch with an arbitrary random seed. As in our submission, we generated 70 coherent
38 groups of training points. For each group, we tried 50 values of r but found negligible
39 variation between the actual effects (max difference of 3×10^{-4} between r 's). Figure
40 S1 shows that influence is highly correlated with actual effects (for an arbitrary r) on
41 the test point with highest loss (Spearman $\rho = 0.96$), even in this non-convex setting.

42 **R3: Case studies.** R3 is right to point out the distinction between removing certain
43 labeling functions (LFs) and crowdworkers from the training data, which we do in
44 our work, and actually changing the LFs or the crowdworker recruiting policy. For
45 example, it is possible that explicitly encouraging LF programmers to create LFs with
46 higher coverage could result in spurious LFs that lower model performance, even
47 though the naturally-obtained high-coverage LFs are the most helpful for the model.
48 Determining the actual effect of manipulating data collection will require systematic
49 user experiments, and we will clarify and emphasize this distinction. We thank R3 for
50 bringing this point up.

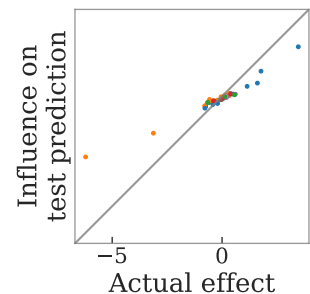


Figure S1: Predicted vs. actual effects on an MLP. Colors are as in our submission.