

1 The thorough reviews asked numerous questions and suggested great clarifications/improvements. Due to space  
2 limitations, we are not able to respond to all of them, but we will incorporate all the feedback into the next version.

3 **Reviewer 1.** You make an excellent point that the skin-tone face example is a poor motivating example since one  
4 should collect a diverse set of labeled faces for training. We will give this as an example of how \*not\* to use our  
5 algorithm, and we are so grateful that you identified this issue. We will also note that research on PQ algorithms may  
6 help raise awareness that P and Q may differ. We will switch to the following motivating natural example:

7         One wishes to provide a service classifying medical scans as normal or abnormal. Training data  
8 consists of a set of volunteered examples hand-labeled by multiple radiologists over a period of time.  
9 Test examples are to be classified in large daily batches. A concerning distribution shift may occur  
10 due a new disease, e.g. COVID-19 scans show “the presence of bilateral nodular and peripheral  
11 ground glass opacities” which may not occur in labeled training data [Sawani, 2020]. If there aren’t  
12 enough unlabeled test examples at classification time, periodic rejections may be useful *in hindsight*,  
13 e.g., it could be immensely helpful if a machine can recognize a new disease (or a problematic change  
14 in the scanning pipeline) at the end of a week.

15 About the question of when one cannot simply label a sample of test examples: note that large-scale classification  
16 services may have millions of test examples per hour, but it may be impractical or even unethical for human annotators  
17 to label them. For instance, privacy policies may prohibit random private Twitter or SnapChat messages from being  
18 classified by humans, though one may have curated messages (e.g., public tweets) at train time. One may simply  
19 not know how representative the training examples are of the private messages, and one would prefer to abstain then  
20 misclassify at test (deploy) time.

21 About the concern of whether or not the adversary can choose examples based on  $S$ , in fact, in our model  $S$  is a given  
22 deterministic function of the training data and unlabeled test examples provided by the adversary. Since the adversary is  
23 all-knowing and all-powerful (e.g., knows the training data), they can determine exactly what  $S$  will be based on their  
24 chosen test set. Hence  $S$  is effectively known. Forcing the adversary to provide a batch of test examples avoids the need  
25 for an “arms race” as the classifier always wins the game.

26 In our empirical evaluation, the simple “baseline” worked really well on our toy example, illustrating our fundamental  
27 contribution: *PQ learning is possible but provably requires unlabeled test data*. That is why we did not implement  
28 any of our more advanced algorithms, thus their value remains theoretical at this point. Future work evaluating and  
29 designing new PQ algorithms on real-world datasets is indeed of interest.

30 Finally, thank you for referring us to the work on “unrestricted adversarial examples”, we will add this to the related  
31 work in the updated version.

32 **Reviewer 2.** Thank you for the enthusiasm and presentation suggestions. We will try to clarify the transductive  
33 setting. As you suggest, one may think of a one-to-one map between examples in  $\tilde{x}$  and  $z$ : one way to think about it is  
34 that every example in  $\tilde{x}$  is the result of a manipulation of some example in  $z$ . To see why our transductive guarantees  
35 are about  $z$  and not  $P$ , consider a case where  $\tilde{P} = Q$  is uniform on a huge set  $X$ . It would be unreasonable to reject the  
36 entire test set  $\tilde{x}$  (and thus achieve zero test error), though if this is all you rejected you would have a negligible rejection  
37 rate on  $P$ . Regarding L159, It is true that the notion of arbitrary test examples is more general than the notion of  
38 adversarial examples with some specific perturbation set (e.g.  $\ell_p$  perturbations). However, in this paper we do selective  
39 classification, which is unlike work on adversarial examples where a prediction is usually always required (even on  
40 perturbations). Also, thank you for referring us to the line of work on adversarial examples that uses unlabeled data, we  
41 will cite these papers in the updated version. We will replace the ? symbol with a different symbol, and we will address  
42 your many other questions though space prohibits us from commenting on them here.

43 **Reviewer 3.** Thank you as well for your enthusiasm and suggestions. As you suggest, the rejection rate bounds can  
44 be improved beyond statistical distance. In Appendix B, we give an example where a tighter bound of 0.1 can be shown  
45 while the statistical distance is 0.91, and Equation (2) also gives a refinement. We agree that it would be especially  
46 interesting to consider better bounds specific to the given concept class  $C$ .

47 As the reviewers suggest, this work is clearly not the final word on learning with arbitrary adversarial examples but we  
48 hope the theoretical results on the possibility of PQ learning inspire people to work in this area.

## 49 References

50 Jina Sawani. How does covid 19 appear in the lungs? *U of M Health Blog*, March 2 2020.