

## 1 **Author Response: Stochastic Optimization for Performative Prediction – Paper #28**

2 We thank the reviewers for their feedback and look forward to incorporating these comments into our revised manuscript.  
3 We address below the remaining questions and comments.

### 4 **Reviewer 1**

5 **Further intuition & examples.** We appreciate the suggestion of including an additional example in the main body of  
6 the paper illustrating performative stability as well as the behavior of lazy/greedy deploy on this example. We agree that  
7 this would aid the reader who is unfamiliar with the framework, and will think about a good example to incorporate.

8 We see two main reasons why stable solutions are desirable. First, in most real-world systems frequent retraining comes  
9 at a significant cost; stability, on the other hand, removes the need for retraining. Second, once the distribution shifts  
10 as a response to model deployment, we in general have no guarantees as to the performance of the deployed model.  
11 Performative stability ensures the model will have nearly optimal predictive power, as alluded to in L:130-136.

12 **Connection to reinforcement learning (RL).** The discussion section in [15] analyzes connections between perform-  
13 mative prediction and RL in detail. One way of understanding performative prediction is as a particular case of a  
14 reinforcement learning problem with special structure ( $\epsilon$ -sensitivity, restricted reward functions) that makes it tractable.  
15 To provide more context we will elaborate on these connections in our revised version, and additionally discuss how our  
16 ideas connect with the stochastic optimization literature in RL.

17 **Scale-invariance of the sensitivity parameter.** Thank you for the careful observation, we will clarify this. Namely,  
18 it is not possible to reduce sensitivity by scaling the parameter  $\theta$ . The reason is that the notion of *joint* smoothness  
19 we consider does not scale like strong convexity with the rescaling of  $\theta$ . For example, rescaling  $\theta \mapsto 2\theta$  (thus making  
20  $\epsilon \mapsto \epsilon/2$ ) would downscale the strong convexity parameter and the parameter corresponding to the usual notion of  
21 smoothness in optimization by 4, however the smoothness in  $z$  (second inequality in L:143) would downscale by 2.  
22 Therefore, the critical ratio necessary for convergence,  $\epsilon \frac{\beta}{\gamma} < 1$ , is unaltered by scaling.

### 23 **Reviewer 2**

24 **Contribution over prior work.** As outlined in the introduction and later emphasized in L:88-91, our work does build on  
25 the framework of performative prediction introduced in [15]. In [15], the authors largely focus on proving convergence  
26 of repeated risk minimization/gradient descent in settings where the learner has access to the *full* distribution. While  
27 they provide an extension of their results to the finite-sample regime, their results in this regime are quite weak since  
28 their analysis relies on concentration of the empirical distribution to the true distribution in the Wasserstein metric. As a  
29 result, they require the learner to collect *exponentially* many samples in the dimension at every step (and, in fact, their  
30 rate is only asymptotic).

31 In contrast, our analysis ensures convergence even if the learner collects a *single* sample at every step, something that is  
32 not at all guaranteed in [15]. To achieve this result, we rely on a different proof technique and a fundamentally new  
33 analysis of the stochastic gradient method in performative contexts. We will extend this comparison in the related work  
34 section to make the technical novelty of our proofs more clear.

35 **Experimental comparison to prior work.** We hope that the above discussion clarifies why a comparison to [15] is not  
36 meaningful. All experiments in [15] evaluate convergence when the learner has access to the full distribution at every  
37 step, while we carry out experiments in a finite-sample setting. Indeed, 4/5 of our plots have the number of samples  
38 collected on the x-axis, a quantity that is not well-defined within the experimental setup of [15].

39 That said, we will happily carry out a new experiment at the population level in Figure 3b to illustrate to the reader the  
40 slowdown in convergence rate caused by the stochastic (versus exact) nature of the gradient updates.

### 41 **Reviewer 3**

42 **Extension to non-convex setting.** We believe that extending our results to non-convex settings is an important and  
43 exciting direction for future work. In Proposition 2.4 we show that repeated gradient descent (essentially the population-  
44 level analog of greedy deploy) need not converge even for weakly convex losses — thus, studying non-convex settings  
45 through the lens of performative prediction would likely require a completely new set of algorithmic tools.

### 46 **Reviewer 4**

47 We thank the reviewer for the positive assessment of our work.