

1 We thank the reviewers for their thorough feedback. We address each of their concerns as follows.

2 **@R1, R3, R4 Numerical experiments:** The main focus of this paper is to extend the "theory" of escaping from saddle  
 3 points to the decentralized setting. However, following the reviewers' suggestion, we'll provide some experiments in  
 4 the revised paper. Here, we briefly mention the setup and results and will include more details in the revised paper. Our  
 5 goal is to show that PDGT is able to escape from saddle points quickly. We compare PDGT with D-GET which is a  
 6 decentralized gradient tracking method that "doesn't use the perturbation idea" [37]. We focus on a matrix factorization  
 7 problem for the MovieLens dataset, where the goal is to find a rank  $r$  approximation of a matrix  $\mathbf{M} \in \mathcal{M}^{l \times n}$ ,  
 8 representing the ratings from 943 users to 1682 movies. Each user has rated at least 20 movies for a total of 9990 known  
 9 ratings. This problem can be written as:  $(\mathbf{U}^*, \mathbf{V}^*) := \underset{\mathbf{U} \in \mathcal{M}^{l \times r}, \mathbf{V} \in \mathcal{M}^{n \times r}}{\operatorname{argmin}} f(\mathbf{U}, \mathbf{V}) = \underset{\mathbf{U} \in \mathcal{M}^{l \times r}, \mathbf{V} \in \mathcal{M}^{n \times r}}{\operatorname{argmin}} \|\mathbf{M} - \mathbf{U}\mathbf{V}^\top\|_F^2$ .

10 We consider different values of target rank and number of nodes. Both methods are given the same randomly generated  
 11 connected graph, mixing matrix, and step size. Further, they are initialized at the same point which lies in the  
 12 neighborhood of a saddle point. Note that in this problem all saddles are escapable and each local min is a global min.

13 In Fig. 1 the experiment is run for 10 nodes, and the target  
 14 rank is 20. Initially both algorithms are stuck close to a  
 15 saddle point and make very little progress. However, since  
 16 the theoretical criterion for PDGT is satisfied in the very  
 17 first rounds (small average gradient and consensus error) we  
 18 have injection of noise. This nudge is sufficient to accelerate  
 19 substantially the escape of PDGT. As we see in the plot, D-  
 20 GET remains close to the saddle point at least until iteration  
 21 1400 where we can see the gradient increasing somewhat  
 22 faster. At the same time PDGT escapes the saddle point,  
 23 decreases the loss and approaches a local minimum.

24 In Fig. 2, the experiment is run for 30 nodes and the target  
 25 rank is 30. Similarly, PDGT escapes from the saddle point  
 26 much faster and decreases the loss substantially before it  
 27 reaches the local minimum. We observe that D-GET also es-  
 28 capes the saddle point eventually following a similar trace to  
 29 PDGT after spending a lot longer at the saddle. Interestingly,  
 30 for this experiment, we observed that some parameters such  
 31 as the stepsize of the first and the second phase, the injected  
 32 noise and the threshold before we inject noise can afford to  
 33 be substantially greater than the theoretical propositions casting PDGT useful for a series of practical applications.

34 **@R1 Importance of the results considering [15]:** On a high level the PGD method in [15] (designed for centralized  
 35 optimization) is not applicable to decentralized settings as it requires coordination between all nodes at each iteration  
 36 which has a prohibitive communication cost. Moreover, notice that PGD is a descent algorithm whose behavior is rather  
 37 well understood, whereas our proposed PDGT method is a non-monotonically decreasing algorithm which requires the  
 38 *careful construction and analysis of a potential function* in order to argue about its convergence. Specifically, when  
 39 the PDG algorithm adds noise to the iterate, there exist a specific number of steps and function decrease threshold  
 40 that characterize the space topology at the point that the noise was injected. When the same injection happens in a  
 41 distributed algorithm as PDGT, the consensus error (which potentially increases exponentially fast) and more generally  
 42 the fact that the algorithm is not strictly descent deems the escape from the saddle point non trivial and the existence  
 43 of appropriate thresholds unclear. Diving into more technical details, one can observe that in comparison to Lemma  
 44 17 in [15], our corresponding Lemma 15 utilizes a new potential function and controls the consensus error among the  
 45 nodes deriving a novel and more complex analysis. Specifically, the last terms in (148) and (149) correspond to the  
 46 consensus error of the escaping sequences. To control these term we need to prove an induction with two parts as shown  
 47 in (151), (152) and invoke novel Lemma 14. It follows that equations (159)-(163) should hold, and an interesting  
 48 connection is derived between the potential function parameter  $\alpha$ , the second phase step size  $\eta_2$  and the dimension  $d$ .  
 49 The aforementioned relation presents a trade-off between the first and second phase stepsizes,  $\eta_1$  and  $\eta_2$ , through  $\alpha$   
 50 and thus achieving the optimal overall convergence rate requires fine-tuning the stepsizes as well as coming up with  
 51 the tightest bounds for quantities such as the target decrease of the potential function  $\mathcal{F}$ , the bound on the norm of the  
 52 iterates  $P$  and the radius of the noise ball  $\mathcal{R}$ . All these values in our analysis are different from the ones in [15] and  
 53 our main proofs are more complex and technical confronting new challenges presented in the distributed framework.  
 54 Finally, the polynomial dependence on  $d$  appears to be crucial in controlling the consensus error after the injection of  
 55 noise. To be more precise, notice that to prove the base of the induction in (151) we need to lower bound  $\zeta\psi_0 = \zeta\mathcal{R}\mu$   
 56 where  $\mu \in \left[\frac{\delta_3}{2\sqrt{d}}, 1\right]$ , deriving the dependence on dimension. The importance of  $\mu$  belonging in this interval becomes  
 57 apparent in Lemma 16 and the corresponding Lemma 15 in [15], where the volume of the stuck region is bounded.

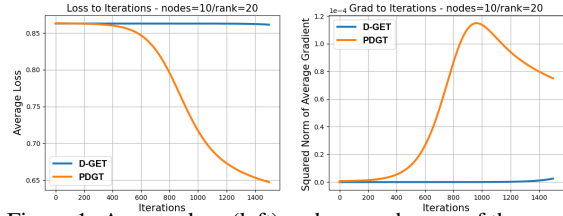


Figure 1: Average loss (left) and squared norm of the average gradient (right) vs. iteration (10 nodes and target rank 20).

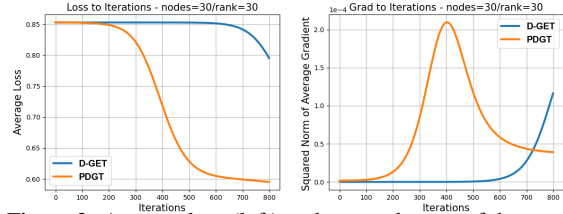


Figure 2: Average loss (left) and squared norm of the average gradient (right) vs. iteration (30 nodes and target rank 30).