

1 We are grateful to the reviewers for their time and comments. For the reviewers’ convenience, we briefly state below
2 the novel contributions of our work, as summarized by R4 (whom we thank for the expert summary):

3 “*This paper shows that under mild conditions, SGD converges to a critical point of general non-
4 convex functions and avoids all strict saddle points, with probability 1. It also presents a convergence
5 rate analysis of SGD once it enters the neighborhood of a local minimum.*”

6 In what follows, we address the reviewers’ comments in order, tagging the reviewers concerned in each as **#RX**.

7 **#R1: Relation to Lee et al [17,18]**. As we explain in Lines 66-68 of the introduction, **[17,18] do not study stochastic
8 gradient descent**, but gradient descent with full, *perfect* gradients – i.e., a **deterministic** algorithm. Specifically, [17,18]
9 show that **deterministic** gradient descent avoids strict saddles from almost every initial condition. The reviewer is
10 therefore not correct in interpreting this statement as an “in probability” result for SGD: *the results of [17,18] provide
11 no guarantees for SGD, from any initial condition.*

12 *Additional comments:* The stochasticity in SGD makes for a drastically different, much more difficult setting. In the full
13 gradient case, there is a well-defined drift that drives GD away from saddle points. This persistent push is no longer
14 present in SGD: this is a crucial difference which we feel may be at the source of this misunderstanding.

15 **#R1#R2: Relation to Pemantle [27]**. Pemantle’s work only applies to **isolated, linearly unstable saddle points**, it
16 does not cover saddle points with a **non-trivial center manifold**. In the deterministic case, strict saddles can indeed
17 be excluded thanks to the existence of local diffeomorphism results based on the center manifold theorem. However,
18 in the stochastic case, the presence of a non-trivial center manifold requires fundamentally different techniques from
19 differential geometry, as we explain in detail in Appendices D.2 and D.3. The reason for this is that there is no longer
20 a persistent drift away from the center stable manifold (in technical terms, there is no “shadowing”). **This major
21 difficulty is not present in Pemantle’s work** (which, again, cannot deal with non-trivial center manifolds); the only
22 relation with [27] is two technical lemmas on random numerical sequences (Lemmas D.1 and D.2).

23 *Additional comments:* The reviewers may have thought that we are making a significantly more restrictive Morse-Smale
24 assumption for the problem’s objective – we emphasize here that *this is not the case*.

25 **#R1#R2: On the rates of Jin et al [14]**. First, as can be seen from (E.3) and (E.42), Thm. 4 gives the precise bound

$$\mathbb{E}[f(X_n) - f(x^*) | X_1 \in \mathcal{U}_1] \leq \frac{2}{\beta} \frac{2\gamma^2}{1-\delta} \frac{G^2 + \sigma^2}{2\alpha\gamma - 1} \frac{1}{n} + o\left(\frac{1}{n}\right). \quad (1)$$

26 We will put this expression for $p = 1$ in the main text. Beyond this, there are two key factual misunderstandings:

- 27 1. **The statements for SGD in [14] and related papers are also asymptotic** because they involve an unknown,
28 probabilistic constant hidden in the $\mathcal{O}(\cdot)$ notation; see Theorem 3, Corollary 4 and Theorem 5 in [14], as well as
29 the corresponding statements in the papers mentioned by R1.
- 30 2. The asymptotic value convergence rate of [14] and related papers is $\mathcal{O}(1/\sqrt{n})$; by contrast, **the value convergence
31 guarantee that we provide is $\mathcal{O}(1/n)$** . The reviewers are therefore incorrect in stating that our rates are similar
32 to those of [14] and related works.

33 **#R1: From high probability to probability 1 via Borel-Cantelli**. This is not possible for (at least) two reasons:

- 34 1. The target probability threshold ζ of Ge et al. is hard-coded in the algorithm’s step-size. Therefore, getting
35 results for different probability thresholds (in order to apply Borel-Cantelli) would necessitate running different
36 algorithms, destroying in this way the validity of the results of Ge et al.
- 37 2. Even if this vital obstacle were to be somehow overcome, the logarithmic dependence of the step-size of Ge et
38 al. on ζ implies that the induced step-size policy would have to vanish at an exponential rate in order to apply
39 Borel-Cantelli. However, it is well known from standard results in stochastic approximation that SGD with
40 summable step-size policies *does not converge* (Kushner and Yin, 1997, Chap. 4), so this approach would fail.

41 **#R2: On Bonnabel (2013)**. We thank R2 for bringing this paper to our attention, we will definitely discuss it! At
42 the same time, we should point out that **Bonnabel’s paper makes the explicit assumption that SGD remains in a
43 compact set** (cf. Theorems 1 and 2). Boundedness assumptions of this kind are prevalent in the literature, and **this is
44 precisely one of the key gaps that our paper closes**: convergence of SGD *without* implicit, unverifiable boundedness
45 assumptions. This was the main weakness identified by R2, so we hope that the above clarifies the merits of our work.

46 **#R3: On the rates of escape**. Deriving rates of escape that hold with probability 1 is a whole new paper in itself.

47 **#R3: On the size of \mathcal{U}_1** . The size of \mathcal{U}_1 only depends on the landscape of f around x^* , *not* δ ; see (E.17) and (E.18).

48 **#R3: On Dashamand et al**. Dashamand et al. refine the analysis of Ge et al. and provide positive probability results
49 for second-order stationary points. There is no overlap with our techniques or results; we will cite it to make this clear.

50 **#R4: On the hitting time to \mathcal{U}_1** . This is a very difficult global-to-local estimate. To the best of our knowledge, no one
51 has succeeded in making progress on similar questions in general non-convex settings, so we do not address this here.

52 **#R4: On the dependence on d** . Great question! The rate does not explicitly depend on d , see (1) above.