1 R1: Q1: It is hard to understand what Remark 6 conveys.

2 A: Yes, the error bound condition refers to the inequality in Lemma 1. Lemma 1 implies that the error bound condition

3 is weaker than the PL condition. Remark 6 means that our analysis only requires the error bound condition to hold

4 though we assume the PL condition at the beginning since it is more widely known in the deep learning literature. The

5 reason that we include this remark is that our experiments directly verifies the error bound condition.

6

7 R1: Q2: how the bound can affect or guide a specific choice of K in stagewise SGD.

8 A: The value of $K$ depends on the choice of $\epsilon$ (i.e., optimization error). Theoretically, the testing error bound (e.g.,

9 Theorem 5) allows us to find an $\epsilon$ that balances the optimization error and the generalization error. However, it only

10 affects K in a logarithmic way. In practice, it is just a small number.

11

12 R1: Q3: It might be better to use "START" in the paper title and Figure 1, instead of "SGD".

13 A: Thanks for the suggestion! We will make the change. Indeed, "stagewise SGD (Vx)" are variants of START with

14 different algorithmic choices. Their common feature is the stagewise step size scheme. stagewise SGD (V1) does not

15 use algorithmic regularization (with $\gamma = \infty$). SGD$(c/\sqrt{t})$ and SGD$(c/t)$ refer to the vanilla SGD (not covered by

16 START) using two different polynomially decreasing step sizes. The comparison between stagewise SGD (Vx) and the

17 vanilla SGD demonstrates the importance of the stagewise step size scheme, which is the key point of this paper.

18

19 R1: Q4: This paper mentions "Corollary 2" several times.

20 A: Sorry for the confusion. Corollary 2 should be Theorem 2.

21

22 R1: Q5: In Line 49, definition of $\mu$. In Line 123, the assumption $|f(w, z)| \leq 1$ is strong.

23 A: $\mu$ refers to the constant in the inequality (PL condition) just before line 49 (between line 46 and line 47). We will

24 make it clear. The boundness assumption on $f(\mathbf{w}, \mathbf{z})$ is only used in the stability analysis for non-convex loss functions.

25 This is following the analysis in [13]. Since a general upper bound $|f(w, z)| \leq M$ only affects the result by a constant

26 factor, we said without loss of generality.

27 R2: We thank R2 for all comments.

28 R2: Q1: what is the main take-away message for me after reading this paper.

29 A: In this paper, we focus on comparing two different step size schemes instead of challenging the classical framework

30 that either analyzes the optimization error convergence or the generalization error of SGD with a particular step size

31 scheme. Most existing theoretical analysis of SGD uses a polynomially decreasing step size or a small step size. However,

32 in practice people mostly use a stagewise step size for SGD, which decreases in a stagewise fashion geometrically.

33 **The main takeaway message of this paper is that we give the first theory to justify why the widely used stagewise**

34 **step size scheme gives faster convergence than a polynomially decreasing step size, i.e., the stagewise step size**

35 **scheme can adapt to the nice properties of deep neural networks.** That is why we compare the results in Theorem

36 5 and Theorem 9 (using a stagewise step size scheme) with the result in Theorem 2 (using a polynomially decreasing

37 step size), and in Figure 1 we compare stagewise SGD with SGD with a polynomially decreasing step size.

38 R2: Q2: Is $\mathcal{A}(S)$ a random variable, or a random probability measure.

39 A: We use $\mathcal{A}(S)$ to denote a randomized model returned by the algorithm $\mathcal{A}$ based on the dataset $\mathcal{S}$. Basically it is a

40 random variable. Please refer to line 89. So $f(\mathcal{A}(S), Z)$ means the loss of the randomized model found by algorithm $\mathcal{A}$

41 on a random data $Z$.

42 R2: Q3: Why is there a <= in the error decomposition after line 96?

43 A: Yes, it is indeed an equality.

44 R2: Q4: It seems that $F_w^\gamma$ is never used in this routine? What does the function O represent at stage 5 of Algorithm 2?

45 A: We will find a better way to present it. The structure of $F_w^\gamma$ that decomposes to the original objective and a quadratic

46 regularizer is used in Algorithm 2. The function $\mathcal{O}$ returns a solution given a sequence of intermediate solutions.

47 R2: Q5: Line 257: I am confused by the way you want to choose theta. I assume you want theta to be very large?

48 A: We expect a larger $\theta$ in order to explain the advantage of stagewise SGD compared with the vanilla SGD with a

49 polynomially decreasing step size, since the vanilla SGD has a complexity of $O(1/(\mu^2\epsilon))$ for reaching an $\epsilon$-level of

50 optimization error and the considered stagewise SGD has a complexity of $O(1/(\theta^2\mu\epsilon))$. Our results in Theorem 6 and

51 Theorem 9 indicate that the larger the $\theta$, the faster the convergence of optimization error and the smaller of the testing

52 error. Nevertheless, $\theta$ is a property of the function. A convex function has $\theta \geq 1$. Our experiments verify that for deep

53 neural networks $\theta$ is also around 1.

54 Response to R3: Thanks for the comments. The green curve is actually for the $\mu$ values across all iterations (the same

55 as $\theta$). We just add a number to mark its average value so that readers can have a sense how small is the $\mu$ as the curve is

56 almost on zero. We will add discussion to discuss the limitation of the presented results.