

1 We thank the reviewers for their thoughtful comments and positive feedback.

2 **Reviewer 1. Q1:** Thank you very much for raising this valid technical point. To satisfy
3 assumption 1, we can introduce box constraints on the entries of \mathbf{H} and \mathbf{W} . Our algorithm
4 will still be applicable in the presence of the box constraints. In addition, our experiments
5 show that if the box constraints are chosen large enough, they will not become active over
6 the iterates of the algorithms and hence will not change the trajectory of the algorithm. We
7 will add this discussion to the paper and share our code in the presence of constraints.

8 **Q2:** As mentioned in [1] and [2], solving the subproblems is NP-hard in the number of
9 constraints for linear constrained problems. Thus, even for a simple NMF problem, solving
10 subproblems in [1, 2] requires exponential computational complexities in the number of
11 constraints (which is the same as the number of variables in NMF).

12 **Q3:** There are different algorithms leveraging the NMF problem structure such as alternating minimization (Alter-
13 Min). Here, we further compare Alter-Min with SNAP (shown in Figure 1). This plot shows that Alter-Min behaves
14 similar to PGD-LS (but it costs less computational time numerically than PGD-LS due to its simplicity of the subprob-
15 lems). We plan to add more numerical experiments and include large-scale cases in the revised version.

16 **Reviewer 2.** Regarding the computational cost of checking SOSP1, first notice that the number of times we check
17 SOSP1 condition is reduced in our algorithm by introducing flag_α . When $\text{flag}_\alpha = \emptyset$, we do not check SOSP1
18 condition for at least r_{th} iterations. Second, to reduce the computational cost of checking SOSP1, we proposed SGDN
19 (a first-order method) for SNAP⁺ with only $\mathcal{O}(d)$ per-iteration computational complexity (instead of $\mathcal{O}(d^2)$ required
20 for other methods such as Lanczos). SGDN approximates the *smallest* eigenvalue in the interested subspace and avoids
21 the evaluation of the restricted Hessian matrix. Our numerical experiments demonstrate that the gain obtained by
22 SGDN procedure outweighs its computational cost. In particular, in most cases, SNAP⁺ outperforms classic projected
23 gradient descent in terms of computational time.

24 **Reviewer 3. Q1:** Yes, we will add the comparison. **Q2:** The line search is a simple step-size selection method with
25 only logarithmic time complexity. Also notice that our algorithm SNAP⁺ does not require Hessian estimation and it
26 only works with gradient. We will include more high level discussions to simplify the understanding of our algorithm.
27 Thank you for pointing this out.

28 **Q3:** The SC condition is only used to establish the connection between SOSP1 and SOSP2. In other words, if the
29 SC condition holds, then the obtained SOSP1 by our algorithms is also a SOSP2. However, our algorithms compute
30 SOSP1 without needing SC condition to hold.

31 **Q4:** Showing the global optimality of SOSPs (in some practical problems) is a very interesting future research direc-
32 tion. However, even when they are not globally optimal, there is still value in computing SOSPs. This is because
33 SOSPs satisfy stronger optimality conditions and are in general a subset of FOSPs. Hence they have higher chance
34 of being globally optimal. Moreover, one can always run his/her own favorite algorithm for finding FOSPs and use
35 our algorithm for escaping saddle points if needed. In addition, our *SOSP-finding* algorithms are “hessian-aware”
36 and hence they can escape saddle points much faster as demonstrated in our numerical experiments. These benefits
37 of algorithms developed for finding SOSPs are the main motivation behind many research works in the optimization
38 society for finding SOSPs. Classical algorithms such as Newton, trust-region, or cubic regularization methods, were
39 all (at least partially) motivated by these facts (long before researchers showing the global optimality of SOSPs for
40 certain problem instances).

41 **Q5:** Thanks for pointing out this issue. We will revise accordingly. **Q6:** If $q_\pi(\mathbf{x}^{(r)})^T \mathbf{v}(\mathbf{x}^{(r)}) > 0$, then we only need to
42 add a minus sign to $\mathbf{v}(\mathbf{x}^{(r)})$, otherwise, just keep it. We will revise to clarify this point further. **Q7:** Yes. The projection
43 step is used in line 18 algorithm 1. **Q8:** Good suggestion! We will include a discussion on the LICQ and/or MFCQ in
44 the revised version.

45 **Q9:** [Prop. 2] The definitions of the exact SOSP1 and SOSP2 have been used in the literature [35], but the relation
46 between these two notations is unknown. Here, we provide rigorous proof to show their equivalence under the SC
47 condition. [Prop. 3] This proposition shows that the introduced notion is continuous and hence proposed algorithms
48 for computing this stationary notion can indeed escape (strict) saddle points asymptotically. Notice that as explained
49 in [31, Remark 2.4] the continuity of the measure is necessary for escaping saddle points (where some previous works
50 could get stuck in a strict saddle point)

51 **Q10:** i) It refers to [Algorithm 1](#). We will clarify it in our revision. In the case the inactive set becomes empty and
52 the first-order stationary condition has been satisfied, which implies that Algorithm 1 has already achieved the SOSP1
53 and therefore will stop. ii) No. Even when a constraint turns active, algorithm 1 is still possible to use PGD to update
54 the iterates (when the conditions shown in line 3 of algorithm 1 are not satisfied), then the active constraint might be
55 deactivated after the PGD step.

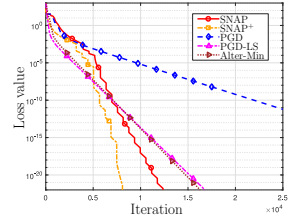


Figure 1: See Sec. E2.1 for the details of the this experiment ($c = 10^{-5}$).