

1 We sincerely appreciate the effort given by reviewers for the time to read, review and provide insightful comments. We
 2 will fix all the typos (in particular Eq.(2.5)). We thank the reviewers for pointing out the crucial bug in the display of
 3 citations. In general, we will improve the presentation, the experimental section and provide open source code. We now
 4 address the reviewers questions in detail.

5 **Assumptions A and B are mostly satisfied in practice.** Assumption A where C is the full space and $h = \frac{1}{2} \|\cdot\|^2$ is
 6 the standard setting of Proximal Gradient Descent. Examples such as QIP [8] and Cubic Regularization [45] satisfy the
 7 assumption for non-Euclidean function h . For the generalization to strict subsets C of \mathbb{R}^d , Assumption A uses a flexible
 8 formulation that does not require the function g to be smooth on the full space (only required on an open set containing
 9 $\text{dom } h$). Also $\text{dom } f \cap C$ has to be non-empty otherwise the optimization will take place over an empty set, which is
 10 not desirable. Moreover, one is often interested in reaching global minimum and do not expect to decrease the function
 11 value to negative infinity, hence the lower-boundedness. Assumption B ensures that the updates of BPG-methods are
 12 well-defined with mild requirements as explained on page 6 in [8], for example, under a classical constraint qualification.

13 **Statistical evaluation experiments.** From the empirical results given in the paper, BPG-methods are seemingly better
 14 than (or at least as good as) PALM-methods. We believe that the trajectories taken by both the methods are different
 15 and is source of good performance of BPG-methods over PALM-methods. The theoretical justification however is an
 16 open question, which we could not justify. We agree that statistical evaluation varied with initialization is an important
 17 experiment and we will add the proposed experiments. With these additional experiments, we hope to find some
 18 new insights (similar to preliminary results below) on the convergence behavior of BPG-methods vs PALM-methods.
 19 Our new preliminary statistical evaluation results with setting used in "Figure 2 with No regularization" show that
 20 PALM-methods get stuck at large objective values (≈ 96728.941 on average) however BPG-methods do not (≈ 378.173
 21 on average). We will give detailed experiments along with precise settings in the next update.

22 **Computational benefits and insights.** We briefly remark some properties of the update steps of BPG-methods. Note
 23 that the updates are independent for \mathbf{U} and \mathbf{Z} , for example, see the simple illustration on Page 2, where updates can be
 24 done in parallel blockwise (ignoring the 1D cubic equation). This should increase the speedup in practice, in particular
 25 for large matrices. Also, note that some terms in gradients overlap, so using temporary variables in implementation can
 26 increase the speedup. We will discuss these tricks, the time and computational complexity in detail in conjunction with
 27 the timing experiments along with a brief pointer in the main paper, as per suggestion by Reviewer 3. We now provide
 28 insights on why BPG-methods are a better choice over other methods, with focus on alternating methods.

- 29 • PALM-methods estimate a Lipschitz constant with respect to a block of coordinates in each iteration, which is
 30 expensive for large block matrices. BPG-methods use a global L-smad constant, which is efficient.
- 31 • PALM-methods cannot be parallelized block wise, for example, in the two block case, the computation of the
 32 Lipschitz constant of second block must wait for the first block to be updated, hence it is inherently serial.
- 33 • Alternating minimization methods do not converge for non-smooth regularization terms and can be inefficient
 34 (for, e.g., ALS) for some matrix factorization problems (see, for example, "Tensor Decompositions and
 35 Applications" by T. G. Kolda and B. W. Bader, also see "On search directions for minimization algorithms" by
 36 M.J.D. Powell). But, BPG-methods and PALM-methods converge (due to linearization).
- 37 • PALM is not applicable to the 2D function $g(x, y) = (x^3 + y^3)^2$, because the block-wise Lipschitz continuity
 38 of the gradients fails to hold even after fixing one variable. BPG-methods are applicable here.
- 39 • PALM is not applicable to, for example, symmetric Matrix Factorization as also pointed in [20] or the following
 40 penalty method based (relaxed) orthogonal NMF problem (similar to Eq.1.1)

$$\min_{\mathbf{U} \in \mathcal{U}, \mathbf{Z} \in \mathcal{Z}} \left\{ \Psi \equiv \frac{1}{2} \|\mathbf{A} - \mathbf{UZ}\|_F^2 + \frac{\rho}{2} \|\mathbf{U}^T \mathbf{U} - \mathbf{I}\|_F^2 + \mathbf{I}_{\mathbf{U} \geq \mathbf{0}} + \mathbf{I}_{\mathbf{Z} \geq \mathbf{0}} + \mathcal{R}_1(\mathbf{U}) + \mathcal{R}_2(\mathbf{Z}) \right\}, \quad (1)$$

41 where second term does not have a block-wise Lipschitz continuous gradient for any $\rho > 0$. But, here
 42 BPG-methods are applicable (similarly also for Projective NMF) with minor changes to the Bregman distance.

- 43 • The block wise Lipschitz continuity of gradients was the primary motive for PALM-methods. Now, with the
 44 BPG-methods, we can directly tackle the original problem with L-smad property.
- 45 • BPG-methods are very general so the choice of applications will increase substantially and we believe that this
 46 will open doors to design new losses and regularizers, without restricting to Lipschitz continuous gradients.

47 We agree that an expanded discussion of competing methods is crucial for the paper. We will also discuss the points
 48 mentioned in the previous answers in detail in the expanded discussion along with state of the art matrix factorization
 49 models and new perspectives with BPG-methods.