(a) Atlas: proposed approach      (b) Atlas: classical approach      (c) 10-arms

1   ✉ **To Reviewer #1:** Thank you for recognising that *"a more principled method of compositional PCA would certainly*
2   *be an interesting development"*. We will add a visualization but leave alr/ilr/clr+PCA to future work.

3   ① **missing visualization:** This is mainly due to space limit, and we prefer to show quantitative comparisons. We do
4   have some potentially meaningful visualizations on the studied datasets. For example, figure (a,b) shows the 2D plots
5   on the Atlas dataset, where the colors mean the nationality and markers mean the DNA extraction method. CoDA-PCA
6   in Fig. (a) uncovers a clustering structure (sketched with ellipses) which is not presented by CLR+PCA in Fig. (b).

7   ② **alr/ilr/clr+PCA** are equivalent to each other, as alr/ilr/clr are different coordinate systems of the same affine space
8   and can be linearly transformed back-and-forth. See figure (c) for our toy dataset under **(vanilla) PCA**, where the black
9   dots indicates that the PCA reconstruction is outside of the simplex. PCA cannot be directly adapted to CoDA: the
10   projection on the principal components, although still in the affine hull, may go beyond the convex hull of the vertices.

11   ③ **nonlinear structure**: we will rephrase L52 and following: using Bregman divergences make explicit a dual affine
12   coordinate space which is in fact the log coordinates of Aitchison. It is in this space that we have affine constraints,
13   which are therefore non linear in the "primal", ambient space (we will quote Shun-Ichi Amari's latest book).

14   ④ Our comment on L91 about the Poisson distribution are misleading. Both CoDA and our proposed approach is
15   not limited to only count data and applies to also real valued data on the simplex. However we do not work out the
16   corresponding exponential family distribution when the base measure is Lebesgue (as opposed to counts).

17   ✉ **To Reviewer #3:** Thank you for agreeing that we are *"convincingly demonstrating an improvement over the current*
18   *standard"*. We will focus to make the fundamental message clearer.

19   ① **transformed vs original data:** The baselines (CLR-PCA; CLR-AE) work on CLR transformed data and learning is
20   separated in two stages: CLR transformation→PCA/AE. The proposed methods ((S)CODA-PCA; CODA-AE) work
21   on original data based on a unified cost function that fits better in the simplex geometry. We show in figure 2 that the
22   proposed methods generally outperforms the baselines.

23   ② It is straightforward to include an **"intercept" term** (e.g. changing the constraint in eq. (21) to $Y = C(U^\top B + \boldsymbol{b}e^\top)$),
24   which is already used by all the compared methods in our experiments. For simplicity we didn't show it in our
25   formulations (will make clear of this in our revision).

26   ③ **visualization of real data**: see our reply ① to Reviewer #1.

27   ④ L65 (singular covariance): being clr, our representation faces the same clr properties **but** our methods avoid the
28   computational pitfalls of these properties (Other comments will be carefully addressed in the revised version).

29   ✉ **To Reviewer #4:** Thank you for seeing that our paper *"provides interesting connections of some theoretical ideas"*,
30   and we will make links clearer between: exponential family PCA, scaled Bregman theorem, and compositional data.

31   ① **extension to ilr**: see our reply ② to Reviewer #1. One of our contributions relies on a principled formulation,
32   simplification and optimisation of clr-PCA loss (Section 4). Ilr coordinates project further the clr coordinates via
33   matrices on the Stiefel manifold: a corresponding approach for ilr-PCA should exploit its Riemannian structure; it is
34   therefore out of the scope of our paper but can be built on top of our Section 4 – we take it as an exciting direction.

35   ② **unclear formulations** will be carefully checked and polished to make sure of consistency and self-containedness.

36   ③ **misleading title**: We respectfully disagree, and see PCA as learning a representation as well (linear projection).

37   ④ In machine learning, Eq.(5) can be attributed to [9] (see for example Section 4 of Roy, Gordon, Thrun, JAIR 2005).
38   However, as rightfully observed, there is an older connection to the information geometry of Bregman divergences, see
39   our reply ③ to Reviewer #1. We will make this more explicit with proper citations to Amari.

40   ⑤ L132 the first term on RHS is the regular Bregman divergence induced by the "perspective" of $\varphi\left(\check{\varphi}(\cdot)\right)$ as in eq.(8).