1   We thank all reviewers for their comments and we will make every effort to address all comments in the revision.

2   **R1 and R5: Rates/finite sample bounds.** One of main difficulties in finding a rate of convergence of our estimator
3   is precisely characterizing how convergence of the mixture components (l. 192) depends on the convergence of the
4   mixture distribution (l. 191) in Theorem 3. Our algorithm can be seen as approximating a distribution with a low-rank
5   symmetric tensor. While matrix convergence implies convergence of spectral decompositions and is well characterized,
6   general tensors suffer from a "lack of closeness[1]." E.g. it is possible to construct a sequence of symmetric rank
7   3 tensors that converges to a symmetric rank 4 tensor[1]. This sort of phenomenon makes it difficult to characterize
8   convergence, and this is only exacerbated since we are working in tensor spaces of functions rather than Euclidean
9   vectors. Interestingly the full version of Theorem 3 (Supp.) implies that for general product spaces "lack of closeness"
10  is solved by assuming densities of the form in (l. 117) identifiable. We will look into this convergence in future work.

11  **R1 and R5: Nonconvex Optimization** The reviewers have correctly assessed that the suggested algorithm is not
12  guaranteed to solve (6) globally. We note that the class of nonconvex optimization problems known to be solvable by
13  descent methods is extremely limited, e.g., phase retrieval, PCA, matrix completion. The objective function we consider
14  is a polynomial of degree six, while the constraint set is convex. While general polynomial optimization is NP-hard[2],
15  we expect and observe that the proposed optimization problem is no more difficult than other common nonconvex
16  optimization problems in machine learning, e.g., training neural networks. Projected SGD is widely used in practice,
17  and there exists some convergence analysis in the literature [3]. Specifically, convergence to a stationary point can be
18  shown under Lipschitz-type assumptions. In the revision, we will add some discussion of the technical details to the
19  main text, making it clear that we do not claim the proposed algorithm finds a global minimum.

20  **R1, R3, and R5: Additional Discussion** It was suggested that the work would benefit from more motivation and
21  discussion of applications. As the reviewers have pointed out, the grouped sample setting is not commonly considered
22  in the literature. One exception is clustering with must-link constraints. There is quite a bit of work in this area with
23  applications in interactive visual clustering and computer vision[4]. Additionally our method can be seen as a continuous
24  version of multinomial mixture modelling, which is used in psychometrics where measurements over time are collected
25  for a group of, for example, bipolar disorder patients and used to identify subgroups within that population whose
26  condition is only evident with repeated temporal measurements [5]. We will add more discussion of applications and
27  assumptions to the revised paper by folding section 3 into section 4 and removing one or more proof sketches.

28  **R1: "NDIGO (the algorithm introduced by this work) outperforms other methods on the training sample, but**
29  **not for out-of-sample setting (where NPMIX performs the best)."** This is a typo. NDIGO performs best (see supp.).

30  **R1: "In the current version, it was not immediately clear to me why (6) and (7) are equivalent."** In the main text,
31  we will add more detailed pointers to appropriate sections in the supplement.

32  **R1: "It would be nice to see a little more experiments on synthetic data."** Before submission, we removed some
33  synthetic data experiments (e.g., variations on the moons dataset, mixtures of mixtures of Gaussians) to keep the paper
34  within eight pages, but we will add these results to the supplemental in the revision.

35  **R5: "How practical is this approach on large datasets?"** The coreset approach appears to scale very well to large
36  datasets. For example, the Twitter experiment we have $2n > 3 \cdot 10^6$.

37  **R5: "What tradeoffs and limitations exist in practice, esp. compared to related algorithms?"** In the optimization
38  problem we must tradeoff between computational complexity and clustering performance by adjusting the coreset size.
39  Another limitation of our approach is when mixture components are technically mutually irreducible but their supports
40  only differ on a set of very small measure, but any algorithm may fail in this case. It is hard to point to an algorithm that
41  will work when our's does not since other approaches rely on notions of cluster separability. If the data is well fit by
42  some parametric mixture or if components are well separated it could be better to use existing techniques.

43  **R5: "The assumption that M is known should be mentioned up front and discussed early on"** We will be sure to
44  emphasize this point in the revision. Although we did not try it, our spectral initialization should offer a heuristic for
45  selecting M by looking for the knee in the eigenvalue curve. AIC/BIC/MDL should also be applicable.

46  **R5: "In the experiments, which of the algorithms compared make use of the paired observations, and which do**
47  **not?"** On line 257 we specify that NPMIX does not utilize pair information, but we will reword things to clarify this.

---

[1]See "Symmetric tensors and symmetric tensor rank" by Comon et al.
[2]See "Complete solutions and extremality criteria to polynomial optimization problems" by Gao
[3]See "Convergence of a multi-agent projected stochastic gradient algorithm for non-convex optimization" by Bianchi and Jakubowicz, as well as "Zeroth-order Stochastic Projected Gradient Descent For Nonconvex Optimization" by Liu et al.
[4]See "Constrained Clustering: Advances in Algorithms, Theory, and Applications" by Basu et al.
[5]For example "A mixed-binomial model for Likert-type personality measures" Allik 2014 or "Cognitive psychometrics: Using multinomial processing tree models as measurement tools" Batchelder 2010