

1 We thank all the reviewers for their very helpful comments and suggestions. Please find below our responses.

2 **Literature review**

3 We will rewrite the existing paragraph on Projection Pursuit (PP) and discuss PP-based clustering methods that have
 4 direct relevance. In particular, we will review criteria for cluster identification including kurtosis (Peña and Prieto, 2001
 5 and follow-up works), first absolute moment and skewness (Verzelen and Arias-Castro, 2017) and relate them to CURE.
 6 We will also discuss algorithms proposed in those papers. In addition, we will review more recent clustering methods
 7 for non-spherical data, such as the one (Kushnir et al., 2019) mentioned by Reviewer 4 based on random projections.

8 **Generalization to multi-class nonlinear clustering**

9 While our theories focus on the two-component elliptical mixture model, the idea of CURE generalizes to multi-class
 10 scenarios, allowing for nonlinear discriminant functions (i.e. feature mappings). We will revise the brief and somewhat
 11 abstract discussion in Section 2.3 and provide more intuitions to the general audience. As is pointed out by Reviewer 2,
 12 we will elaborate why the linear CURE approach is a special case. Due to space constraints, we will defer the most of
 13 these details to the supplementary.

14 We will add a new numerical example (to the supplementary) with the first 4
 15 classes in Fashion-MNIST (T-shirt/top, Trouser, Pullover, Dress). All details
 16 will be included in the final version. In short, we use Wasserstein-1 distance as
 17 the discrepancy measure D for uncoupled regression and compare two classes
 18 of feature mappings: linear functions and fully-connected neural networks
 19 with one hidden layer that has 100 nodes. The learning curves in Figure 1
 20 shows the advantage of neural network and demonstrates the flexibility of
 21 CURE with nonlinear function classes.

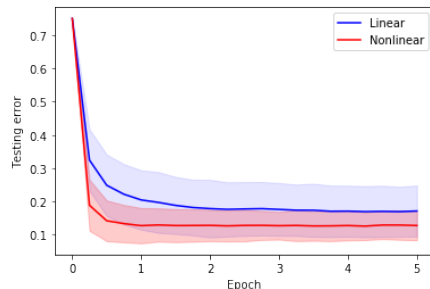


Figure 1: 4-class Fashion-MNIST.

22 Achieving Bayes optimality in multi-class clustering is indeed very challeng-
 23 ing. Under parametric models (e.g. Gaussian mixtures), one may construct
 24 suitable loss functions for CURE based on likelihood functions. We will comment on this in the discussion section.

25 **Numerical experiments**

26 A main motivation for this paper is to deal with stretched (elongated) clusters where directions with the largest
 27 variabilities of data may not be informative for clustering at all. Instead, one should aim for directions onto which the
 28 projected data exhibit cluster structures. This explains why we choose T-shirts/tops and Pullovers in the Fashion-MNIST
 29 dataset for demonstration: the bulk of a image corresponds to the belly part of clothing with different grayscales, logos
 30 and hence contributes to the most of variability. However, T-shirts and Pullovers are distinguished by sleeves. Hence
 31 the two classes can be separated by a linear function that is not related to the leading principle component of data.

32 We conducted new experiments comparing CURE with more algorithms: [1] Model-based clustering (Mclust) in Fraley
 33 and Raftery (1999); [2] Projection Pursuit (PP) in Peña and Prieto (2001); [3] alternations between linear discriminant
 34 analysis and K-means (LDA + Kmeans) in Ding and Li (2007); [4] Minimum Density Hyperplane (MDH) in Pavlidis
 35 et al. (2016). [1], [2] and [4] are implemented using open-source R packages with default settings. While there
 36 is no implementation of [3] available publicly, we did it ourselves following the instructions in the paper. Table 1
 37 shows the results. For randomized algorithms we do 50 independent runs and report means and standard deviations of
 38 misclassification rates. Again, CURE outperforms all the competitors. On a Macbook Pro it takes less than 10 seconds
 39 to converge while others usually require one minute or more. Currently, we are also exploring other datasets.

Table 1: Misclassification rates of CURE and other methods.

Method \ $N_1 : N_2$	1 : 1	2 : 1	3 : 1	4 : 1
CURE (ours)	5.2 ± 0.2%	7.1 ± 0.4%	9.3 ± 0.7%	11.3 ± 1.1%
Mclust [1]	48.7 ± 1.3%	39.1 ± 4.8%	34.1 ± 8.0%	28.2 ± 7.8%
Projection Pursuit [2]	36.9 ± 9.8%	37.4 ± 9.6%	39.7 ± 6.9%	40.6 ± 7.3%
LDA + Kmeans [3]	45.9%	49.0%	45.6%	44.3%
MDH [4]	48.6%	43.1%	38.3%	35.2%

40 It is worth pointing out that CURE learns a classification rule that readily predicts labels for any new data. This is an
 41 advantage over many existing approaches for clustering and embedding, including spectral methods and t-SNE where
 42 out-of-sample extensions are not so straightforward. We will highlight this in the paper.

43 **Other issues.** For the balanced two-class problem we have an explicit construction of f and set the regularization
 44 parameter λ to be 1. Our theory goes through as long as $\lambda \geq 1$, and our experimental results are not sensitive to the
 45 choice of λ . So we choose $\lambda = 1$ to reduce the need of tuning and simplify statements of theoretical results.