

Testing and Tuning SkinnyDip: Noise-Robust Clustering

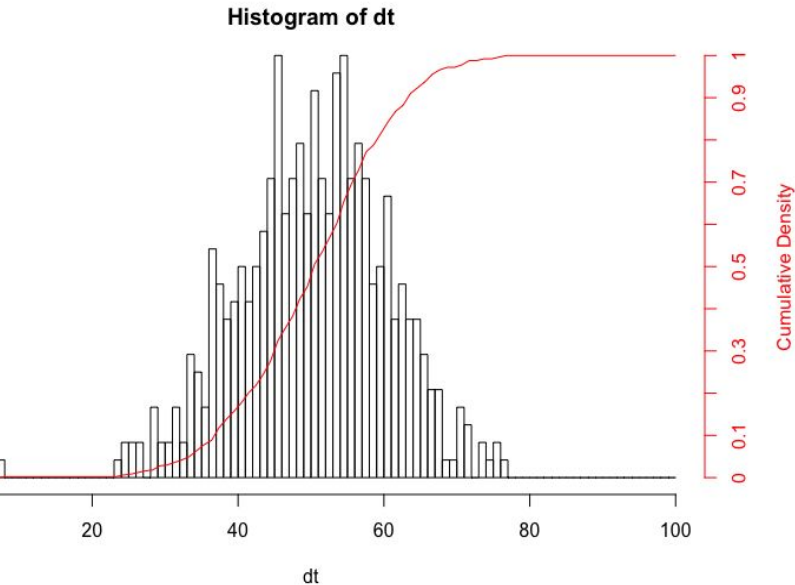
Presenter: Grant King

SkinnyDip

- Non-parametric, deterministic clustering
- Robust to very noisy data
- Based on the “dip-test” for modality

Modality and Empirical Cumulative Distribution Function

- ECDF is a monotonically increasing function



- Unimodal: ECDF is concave up to mode and convex after

Dip Test

82

HARTIGAN AND HARTIGAN

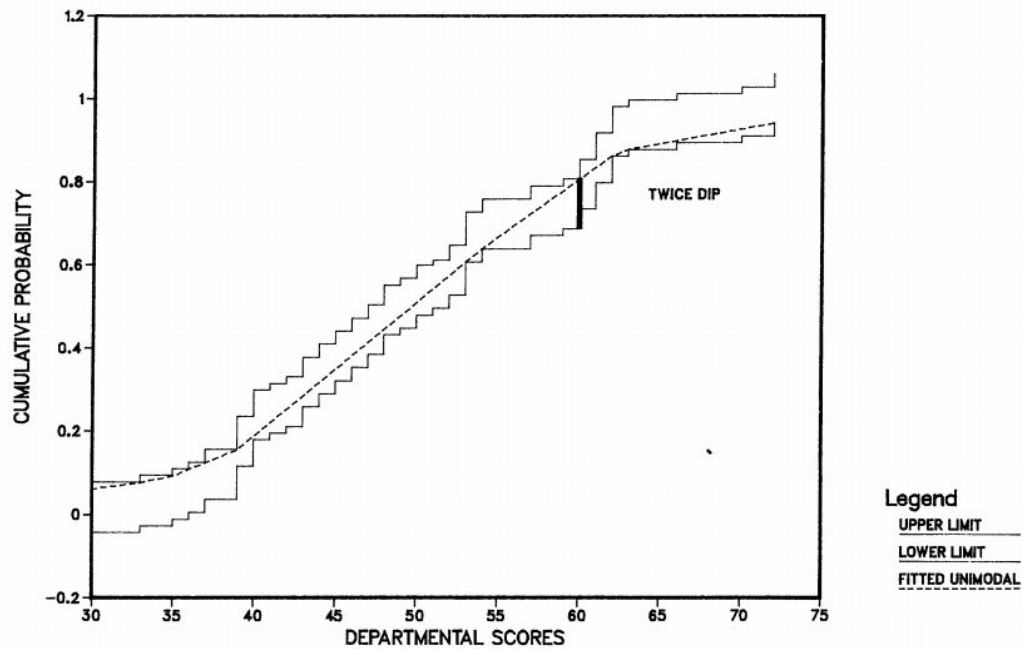
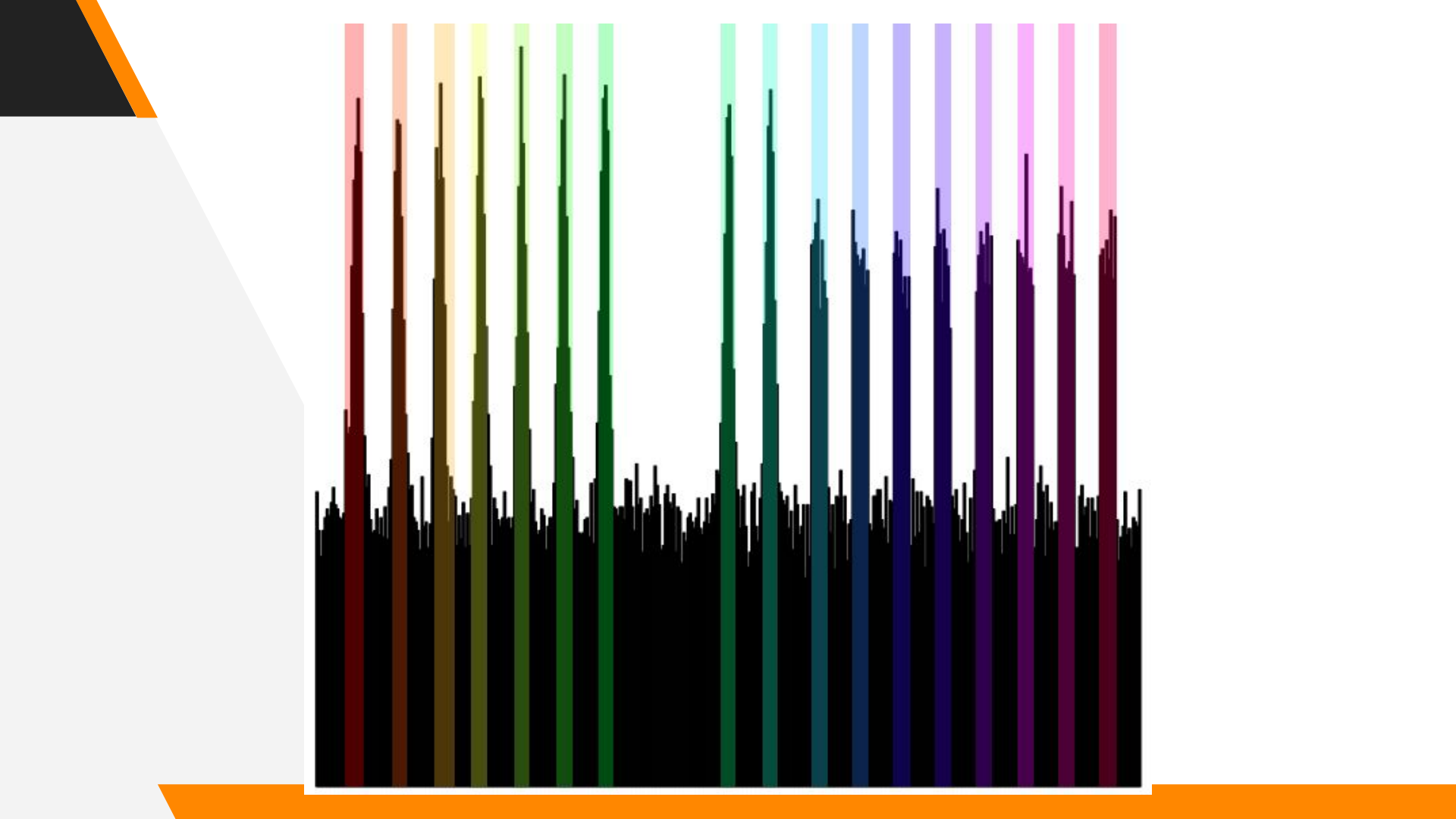
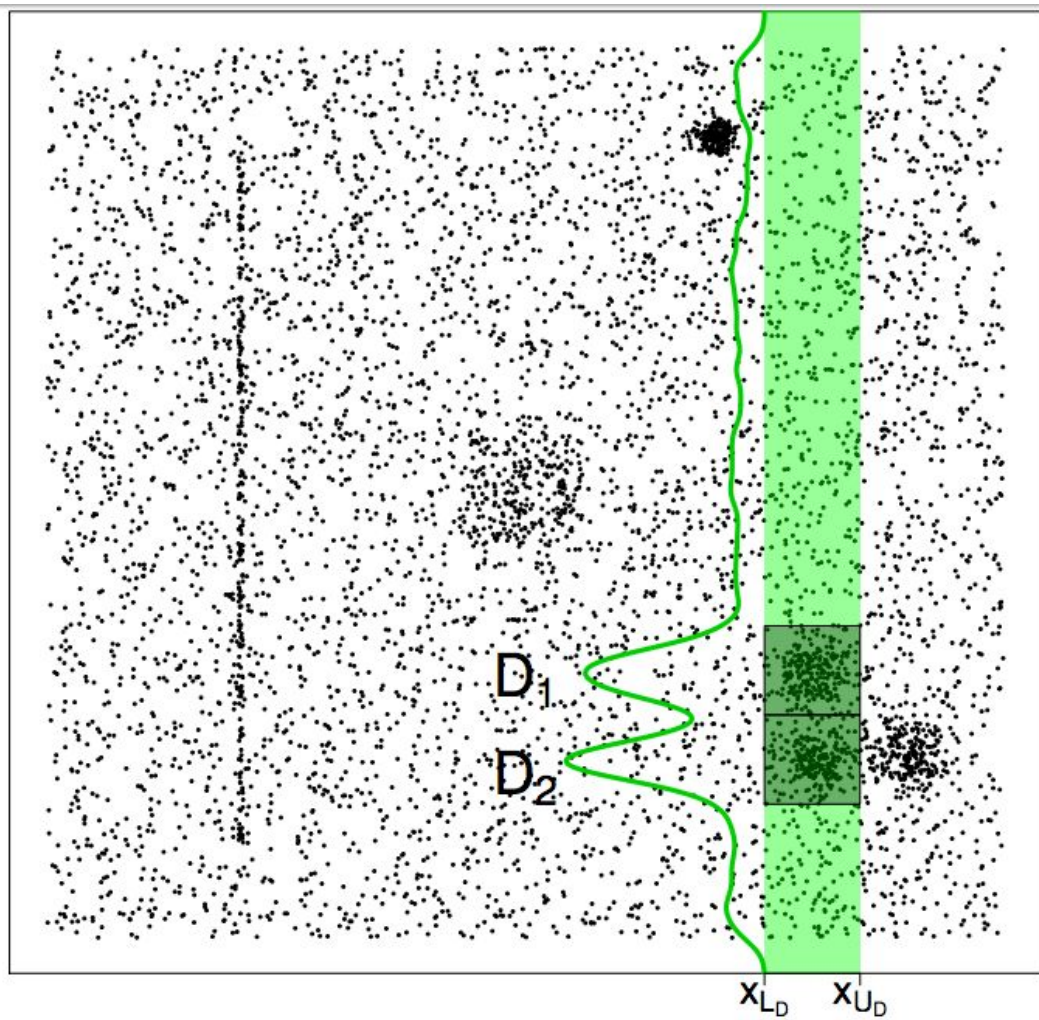


FIG. 1. Faculty quality in statistics department.

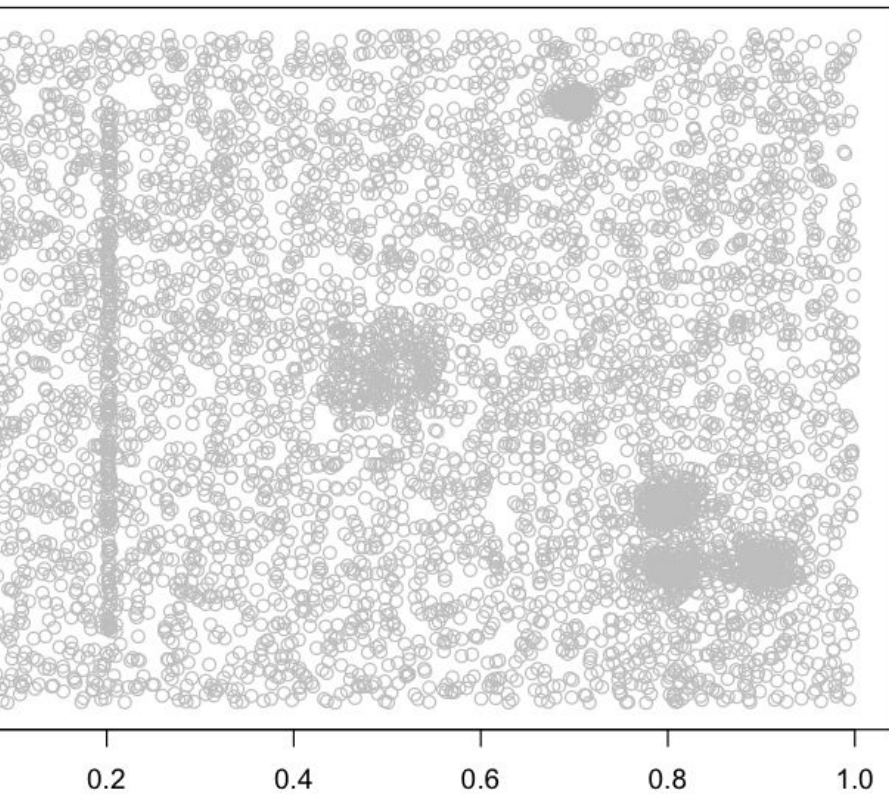




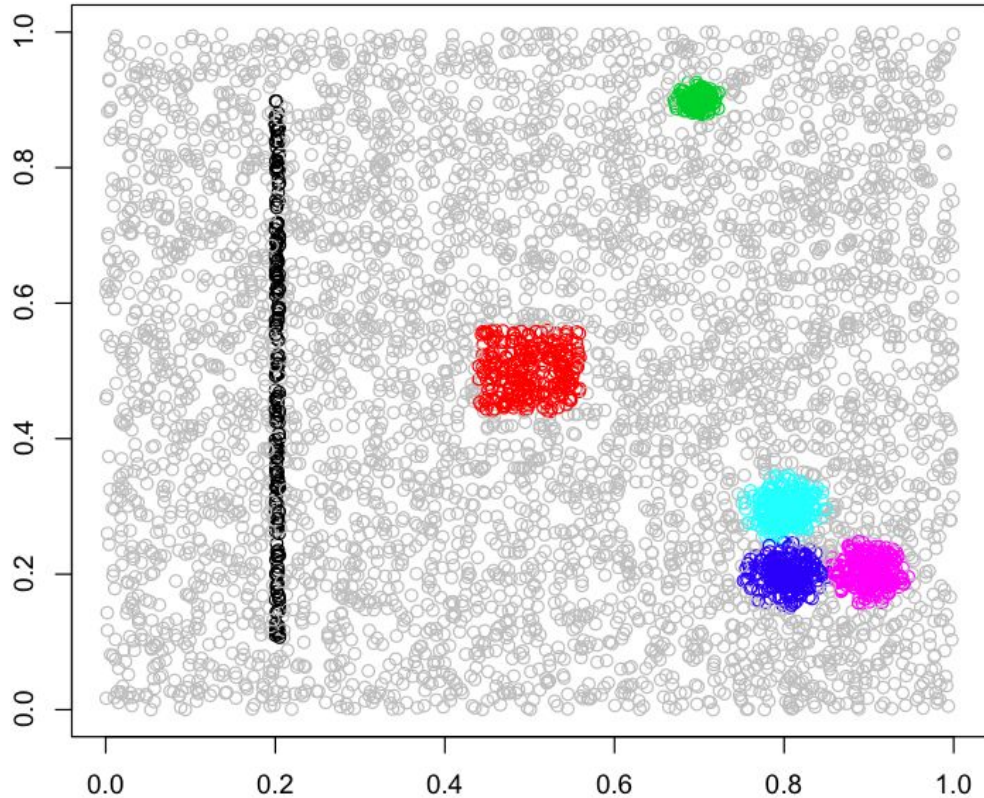
Hypercubic Regions

- SkinnyDip recurses through each dimension of the data and finds approximate unimodal regions
- These intersections of these regions are then combined to form the basis of each cluster
- Accuracy drops exponentially with increasing dimensionality

SkinnyDip Running Example Unlabeled



SkinnyDip Running Example Ground Truth



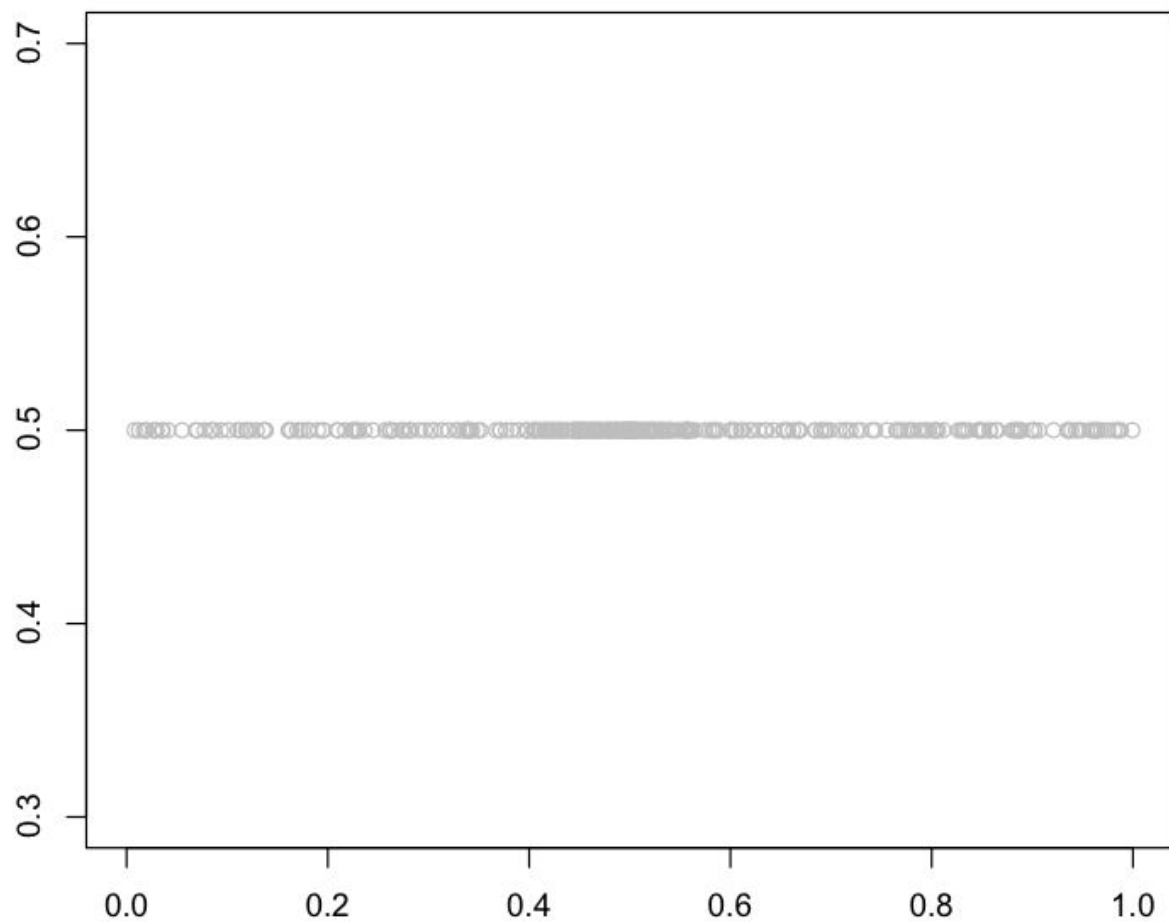
Testing

- Map unsupervised clusters to ground truth labels
- Gain a better understanding of inherent error
- Test multiple dimensions and clusters

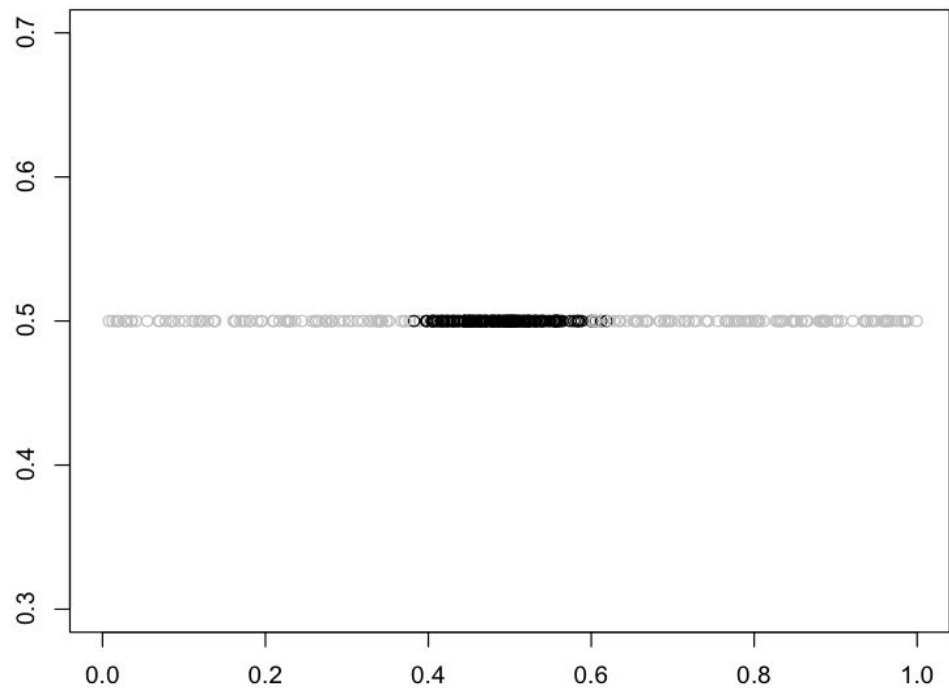
Tuning

- Use distribution of points in each cluster as evidence for underlying variance
- Remap points in the original data using empirical data
 - e.g. Mahalanobis distance using pseudo-inverse cov
- Retest

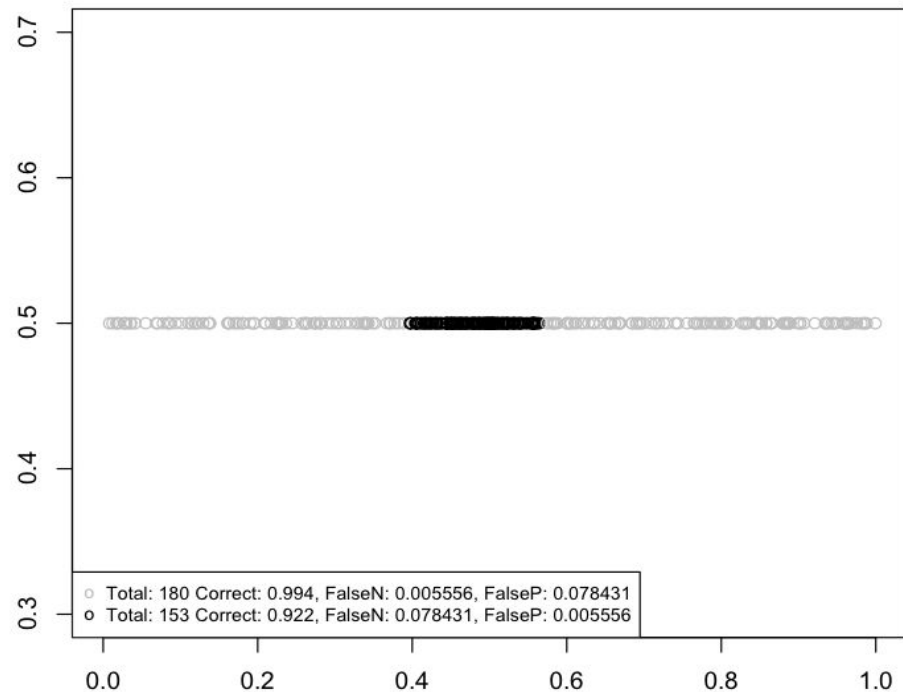
Unlabeled Single Gaussian 1D



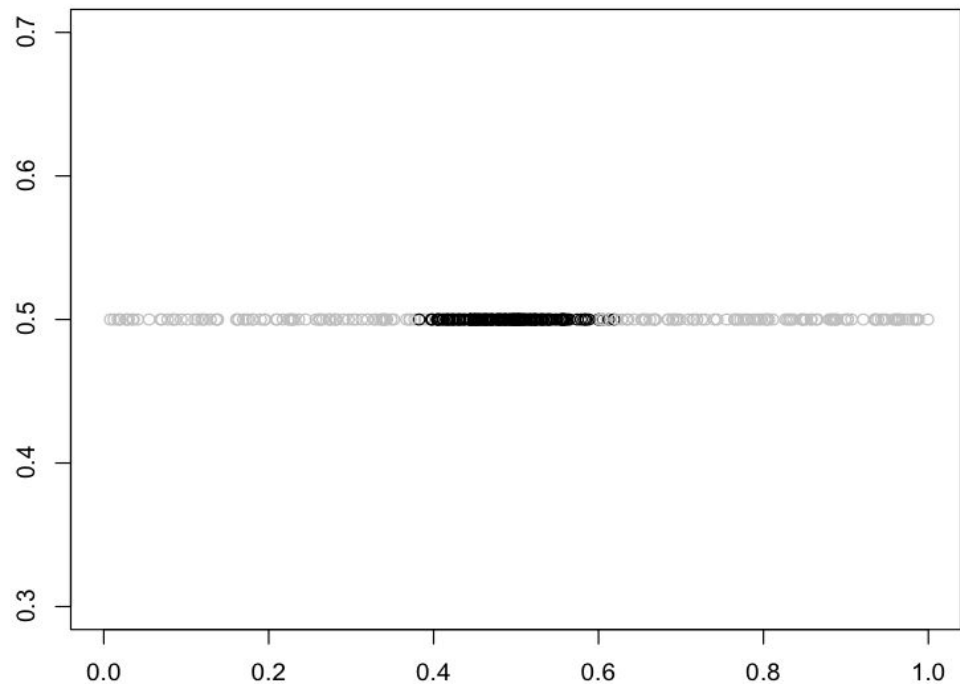
Ground Truth for Single Gaussian 1D



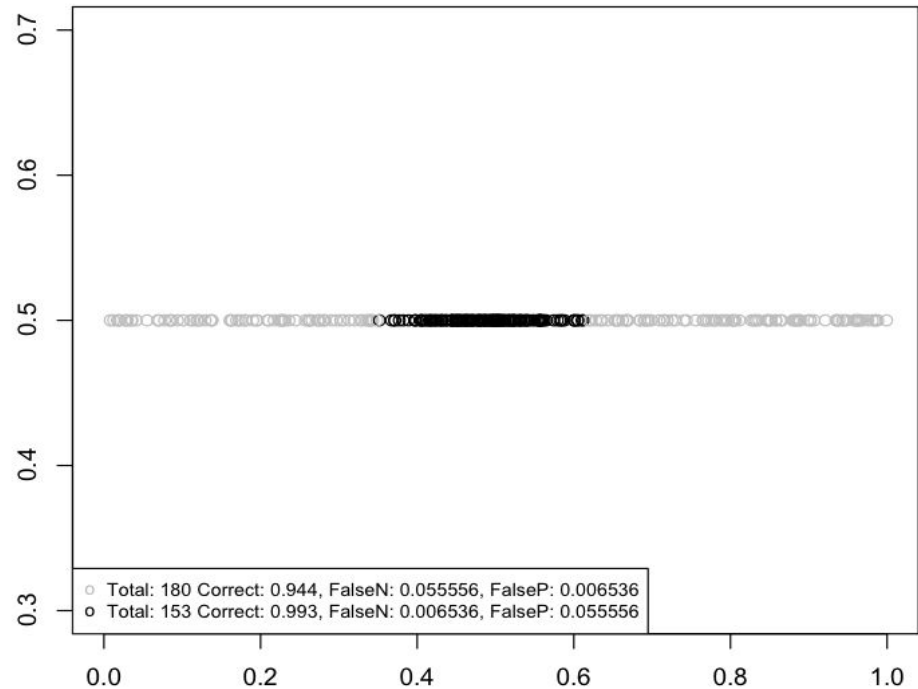
SkinnyDip for Single Gaussian 1D



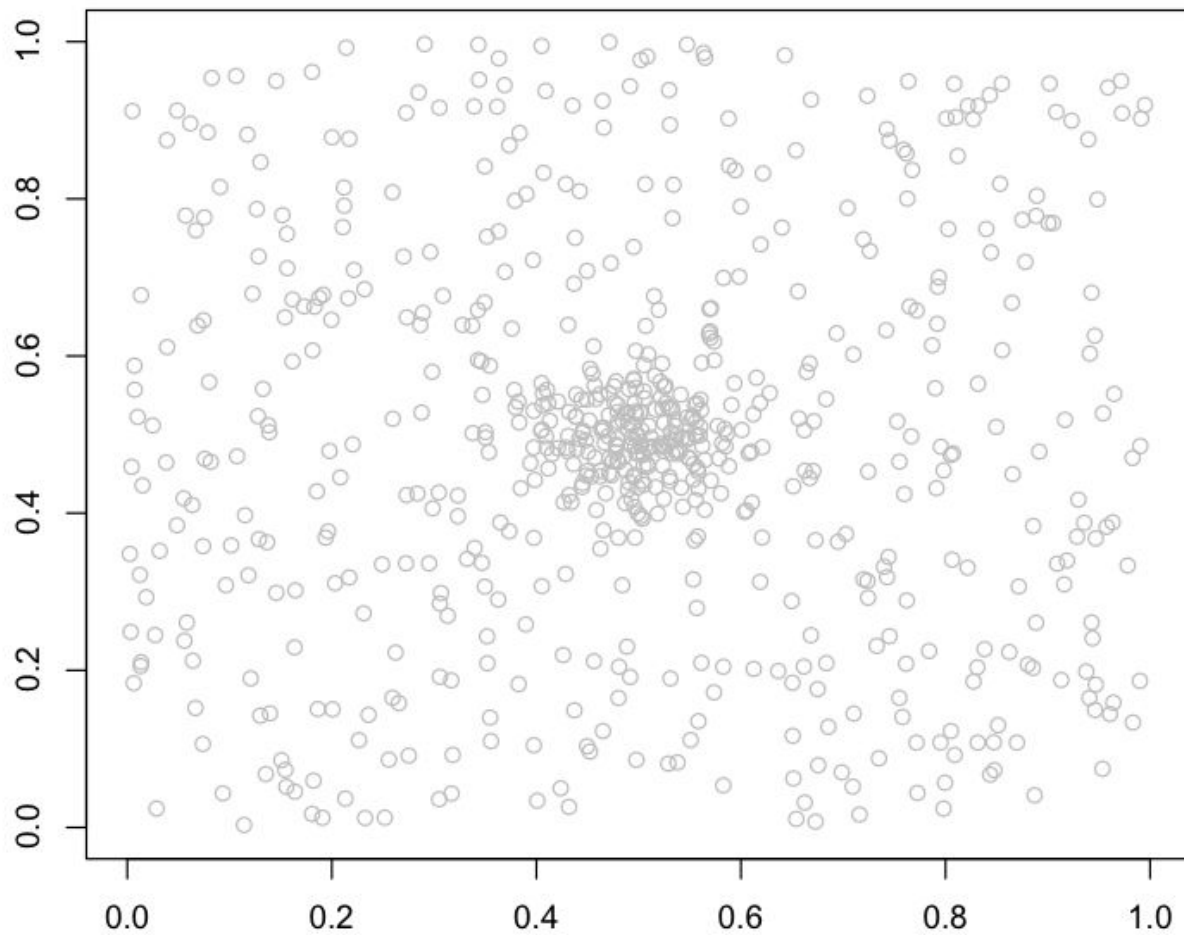
Ground Truth for Single Gaussian 1D



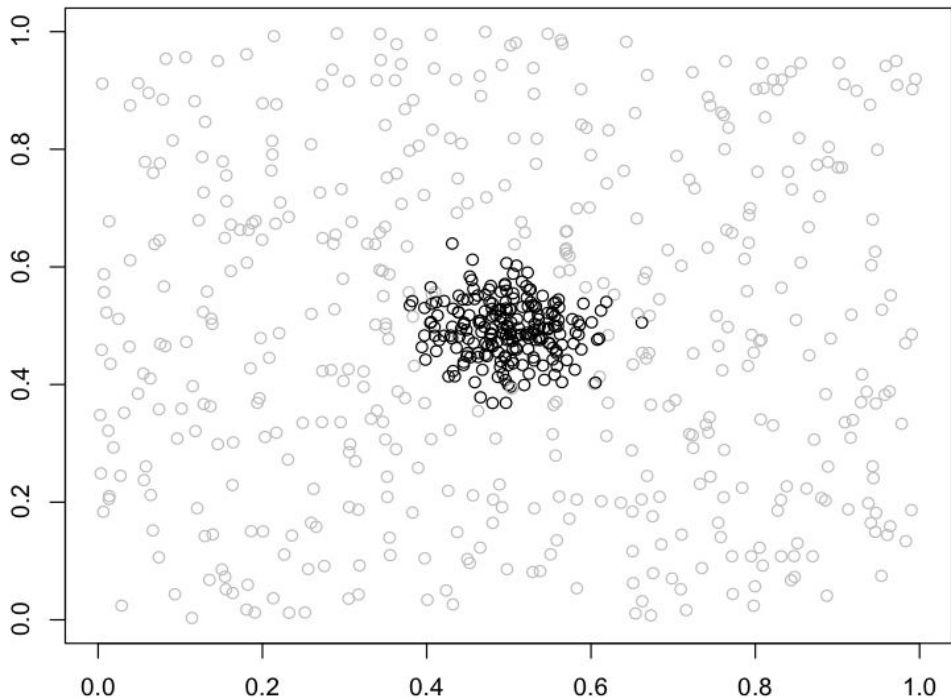
SkinnyDip+Recluster for Single Gaussian 1D



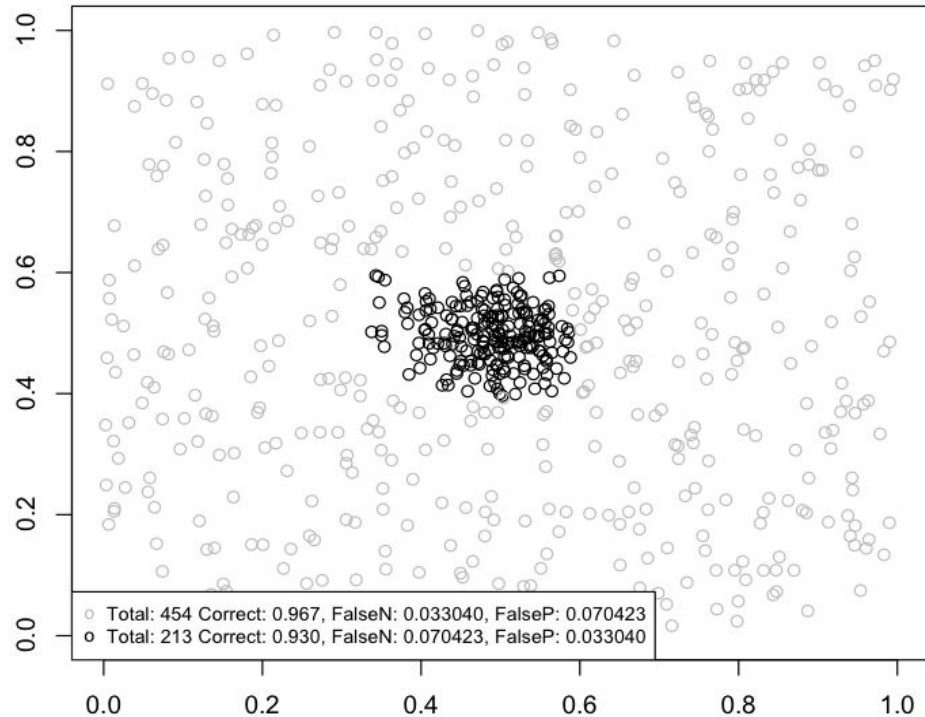
Unlabeled for Single Gaussian 2D



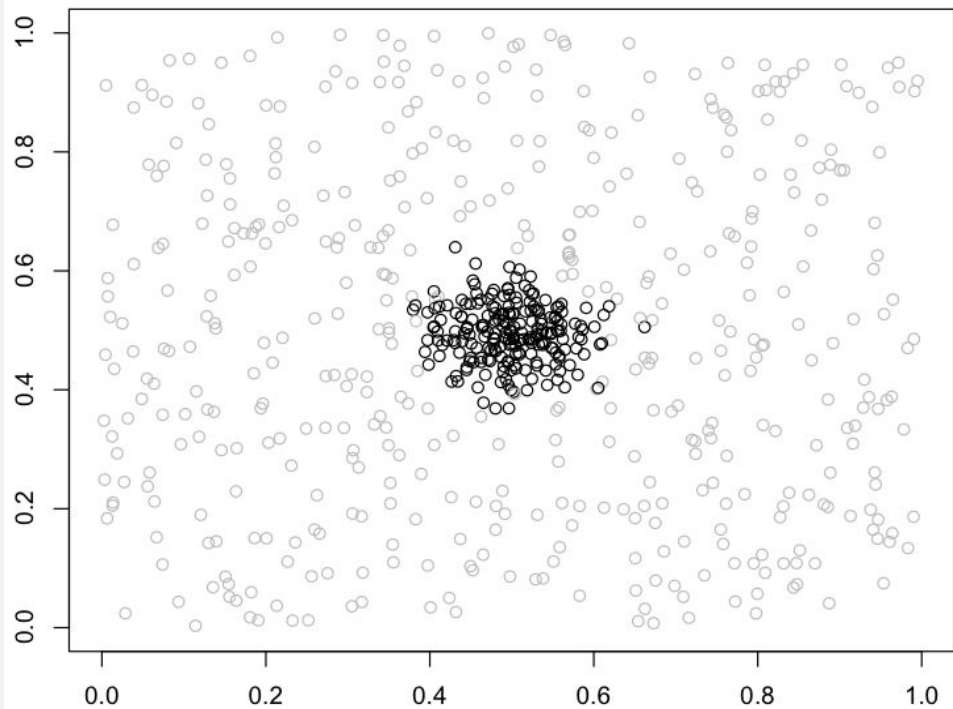
Ground Truth for Single Gaussian 2D



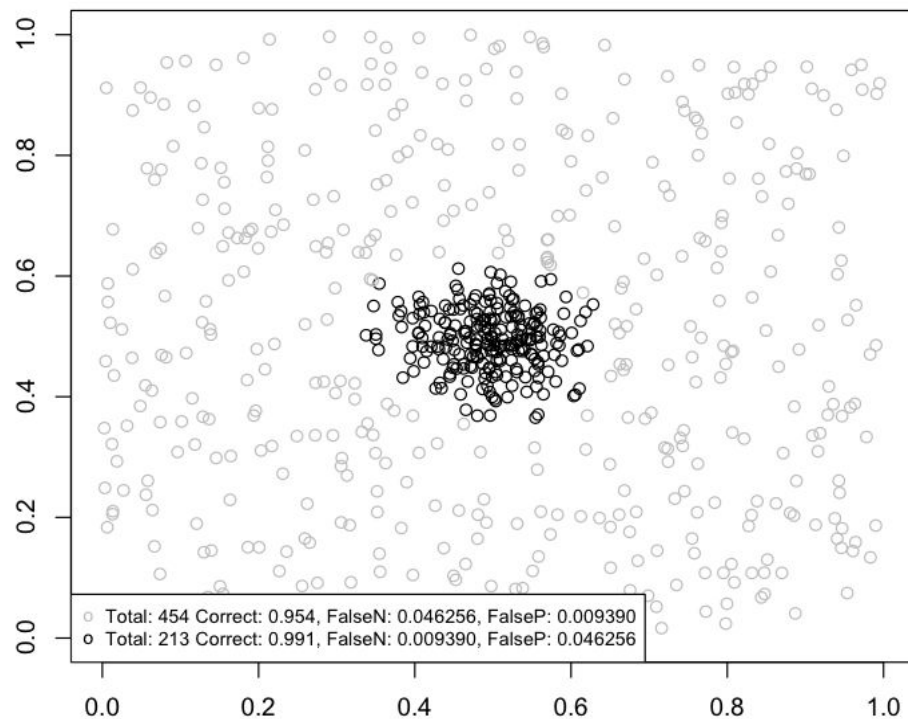
SkinnyDip for Single Gaussian 2D



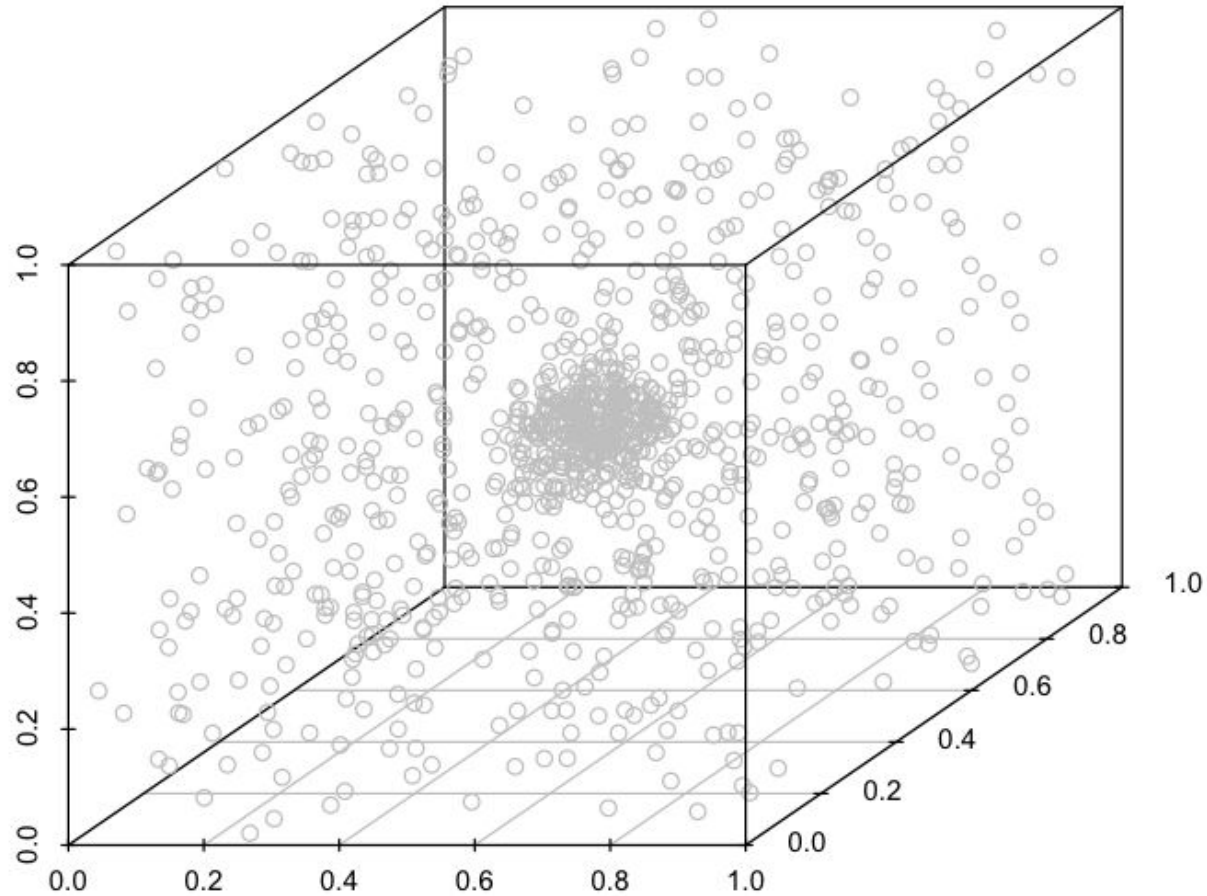
Ground Truth for Single Gaussian 2D



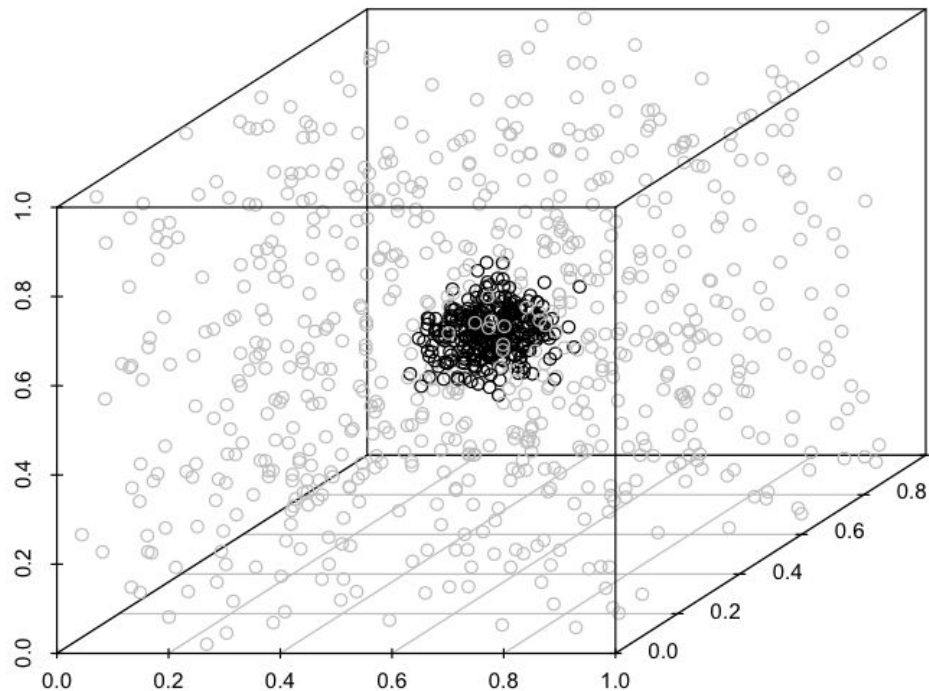
SkinnyDip+Recluster for Single Gaussian 2D



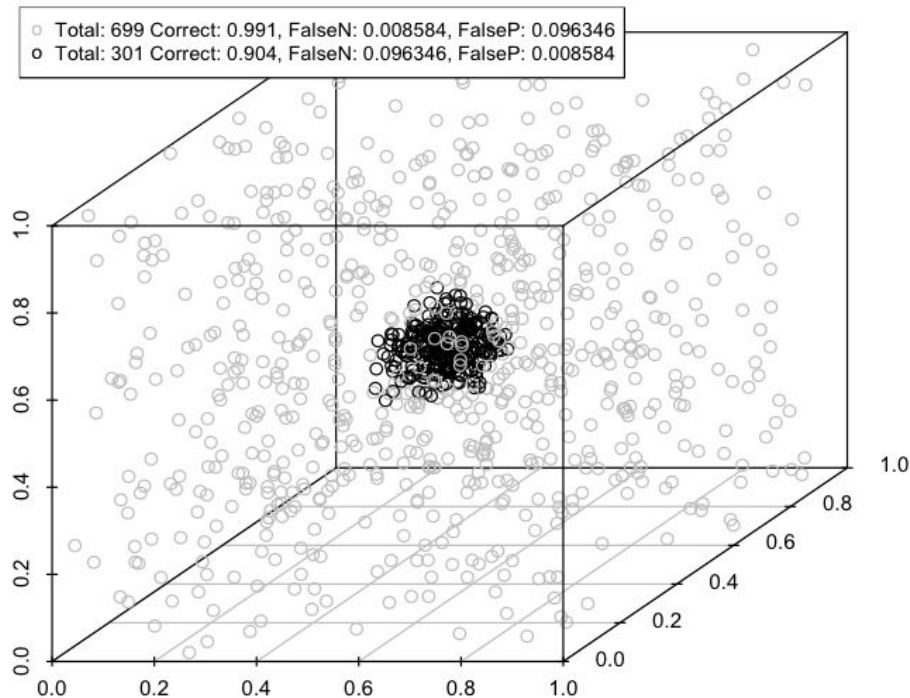
Unlabeled Single Gaussian 3D



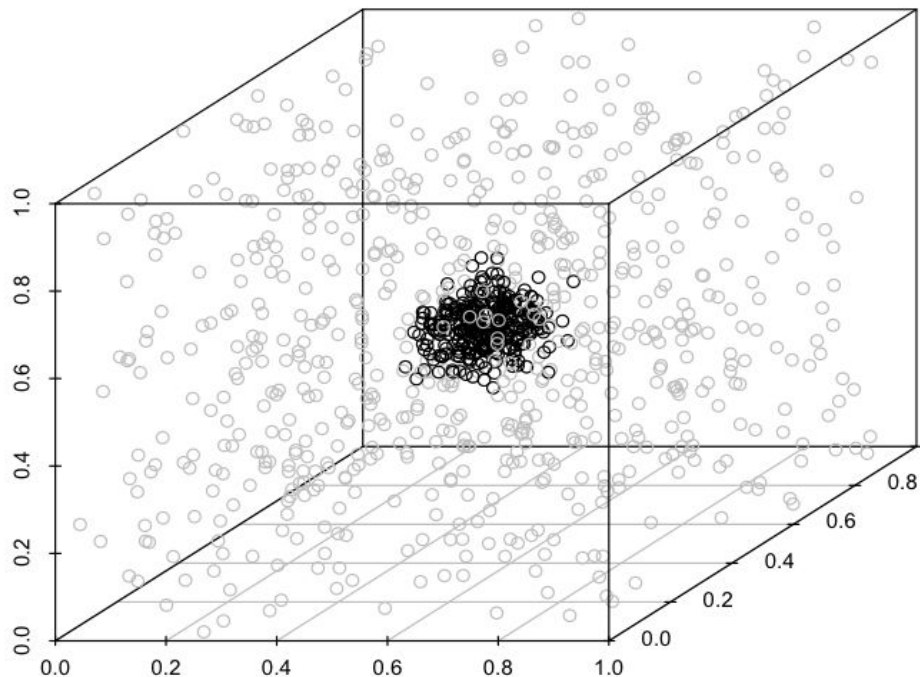
GroundTruth for Single Gaussian 3D



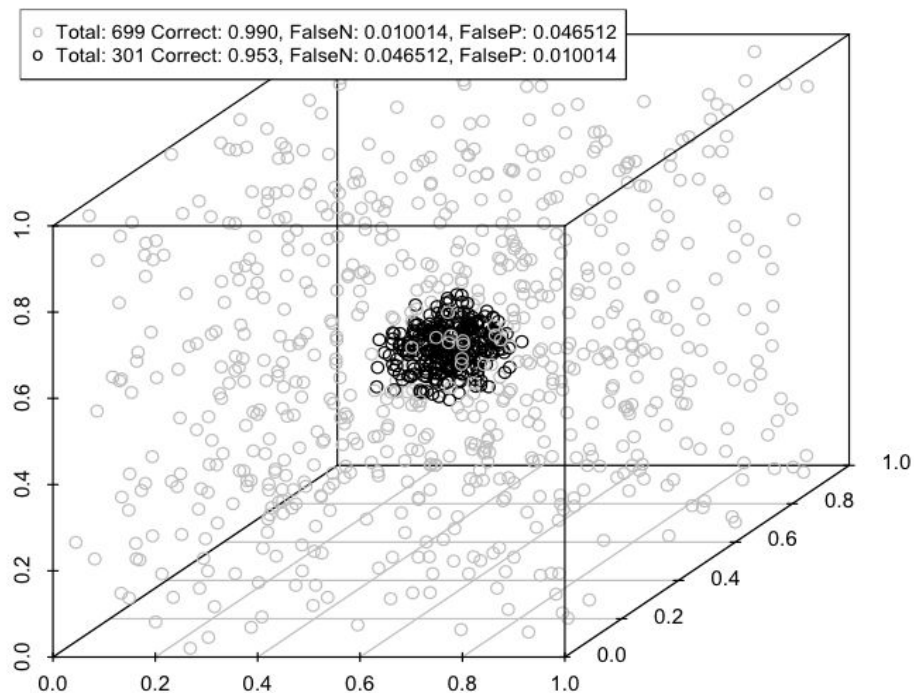
SkinnyDip for Single Gaussian 3D



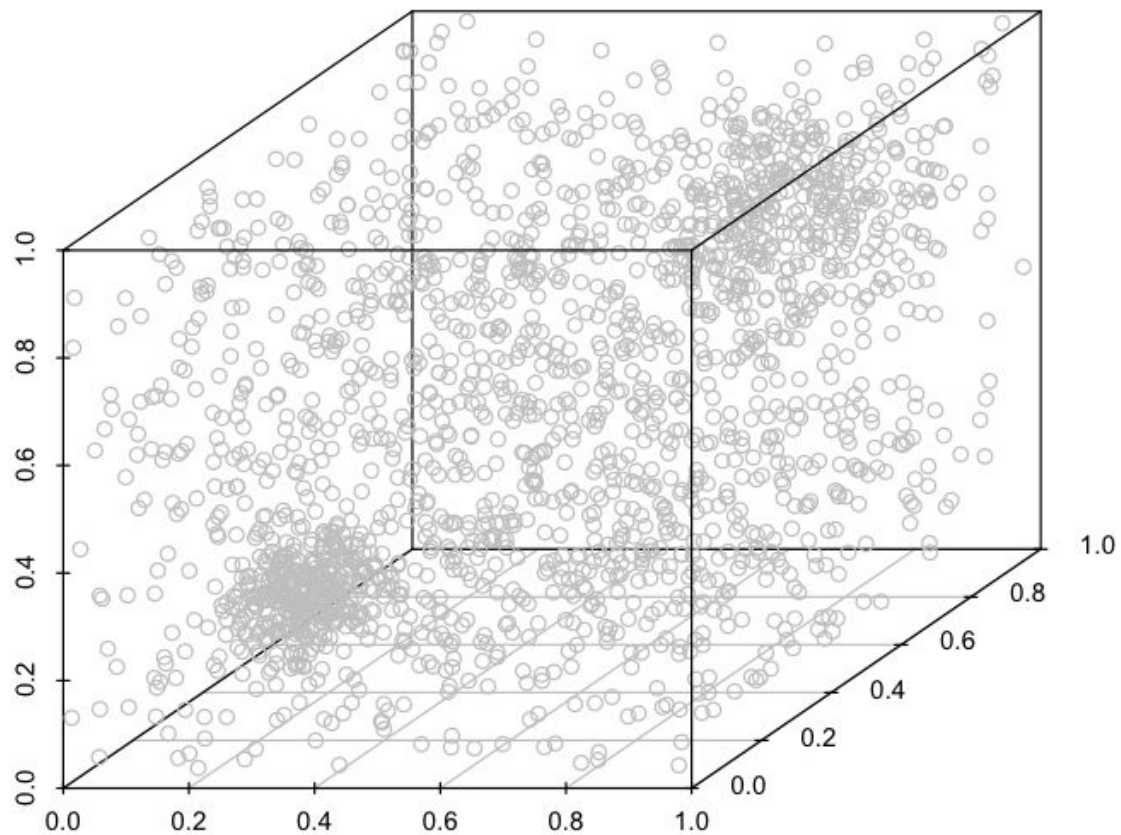
GroundTruth for Single Gaussian 3D



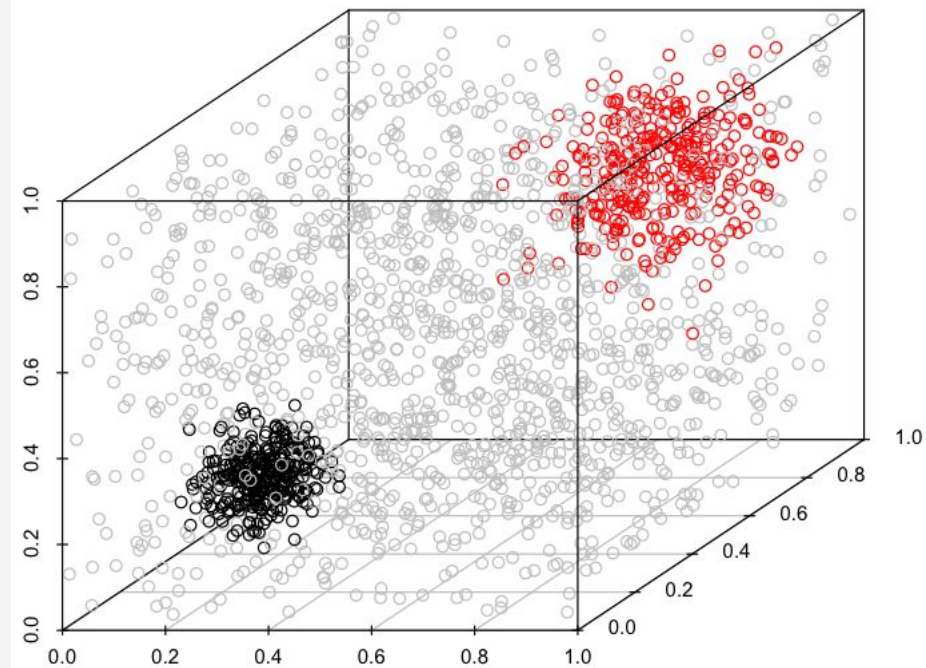
SkinnyDip+Recluster for Single Gaussian 3D



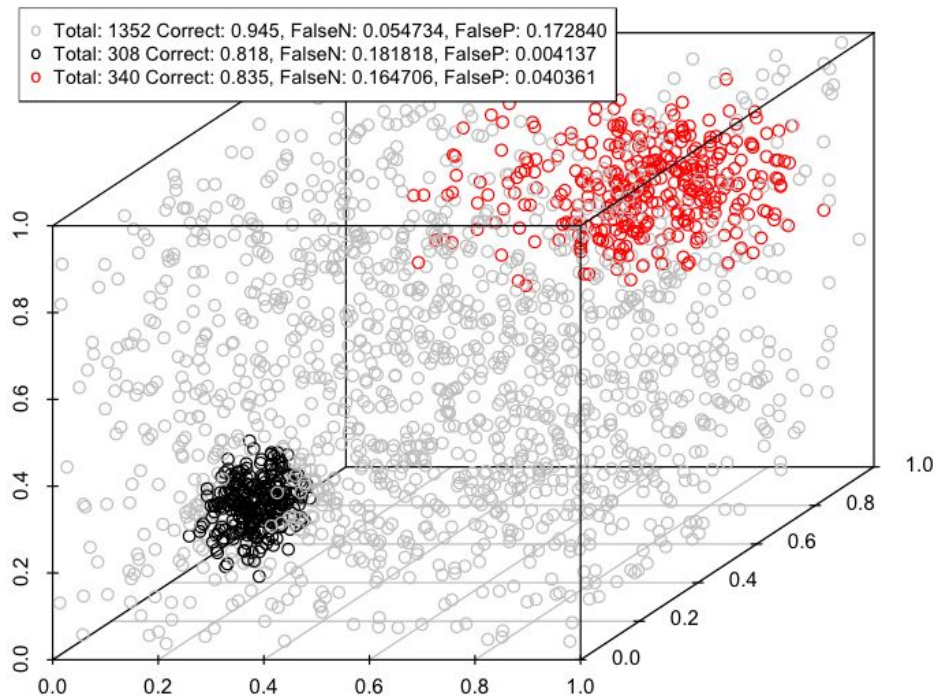
Random Data Unlabeled



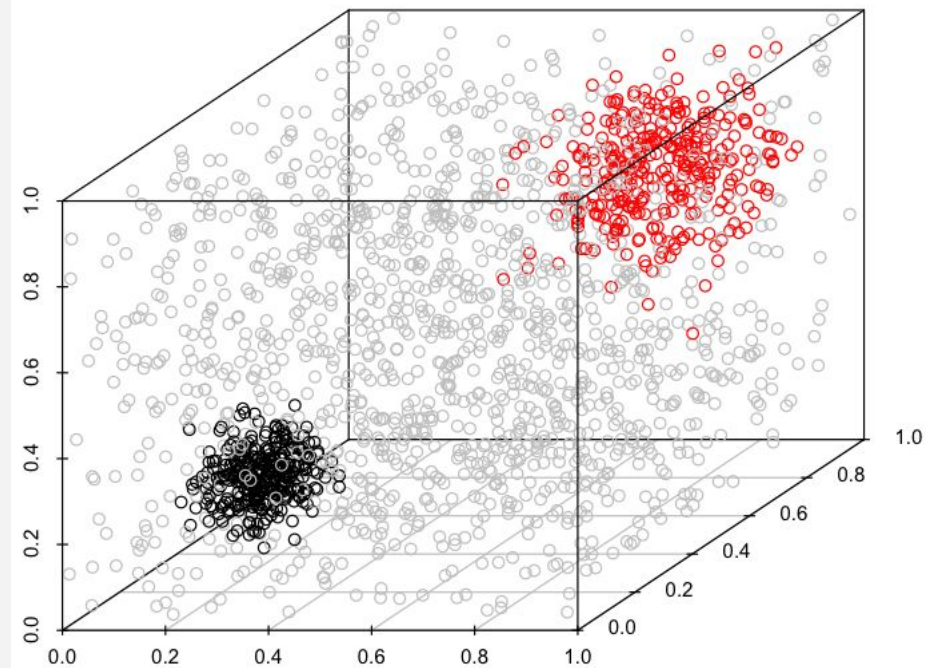
Random Data Ground Truth



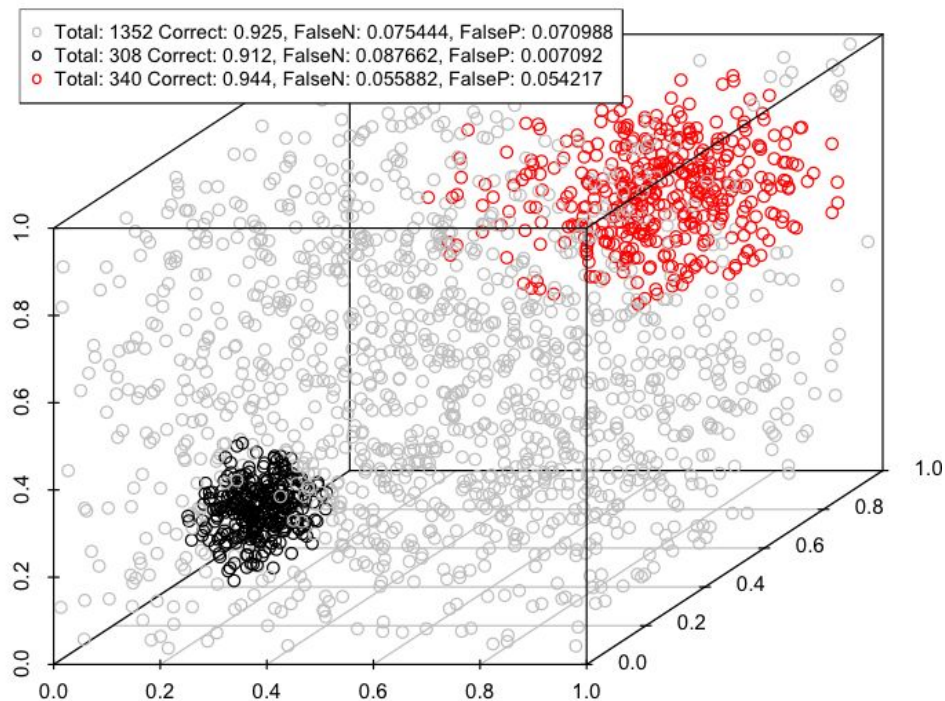
Random Data SkinnyDip



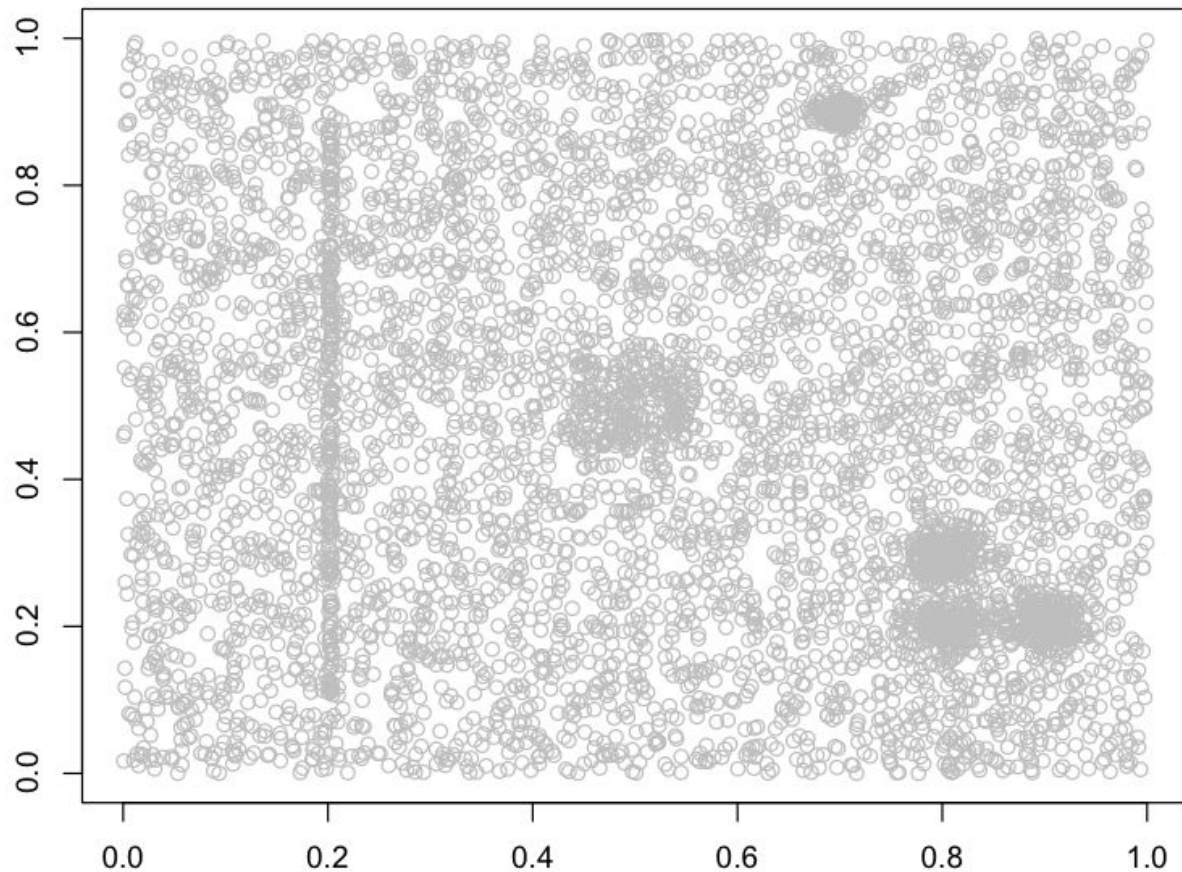
Random Data Ground Truth



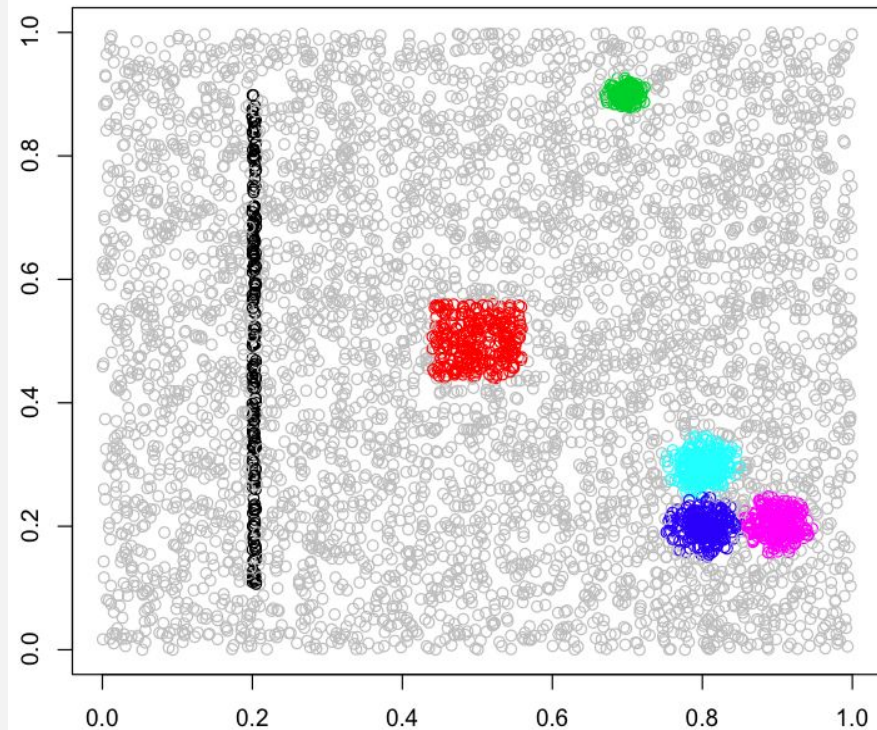
Random Data SkinnyDip+Recluster



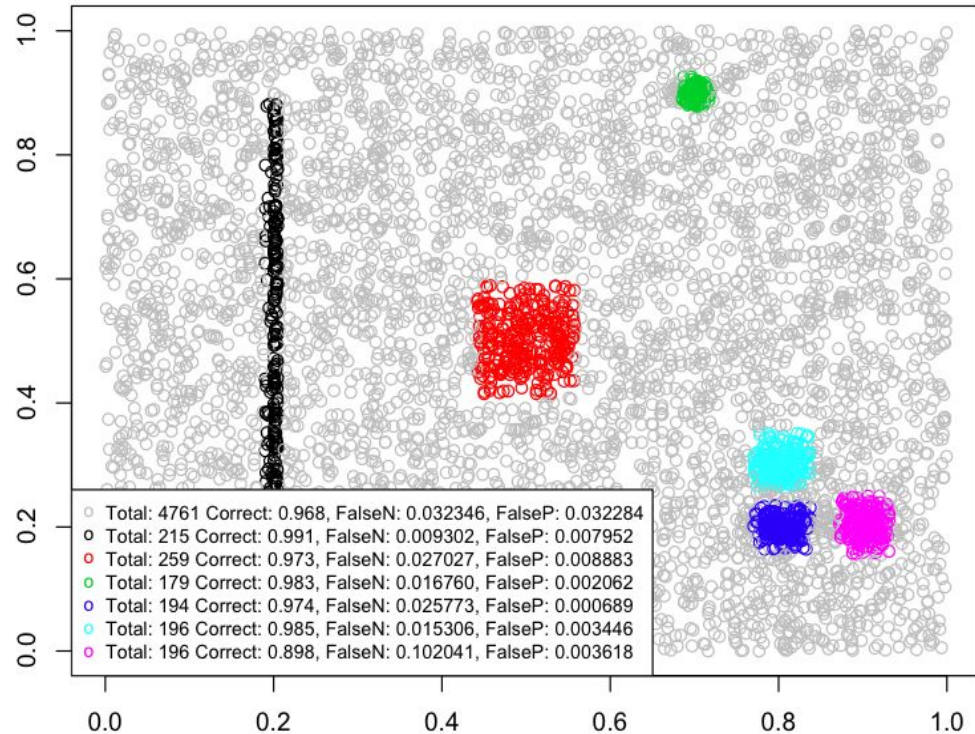
SkinnyDip Running Example Unlabeled



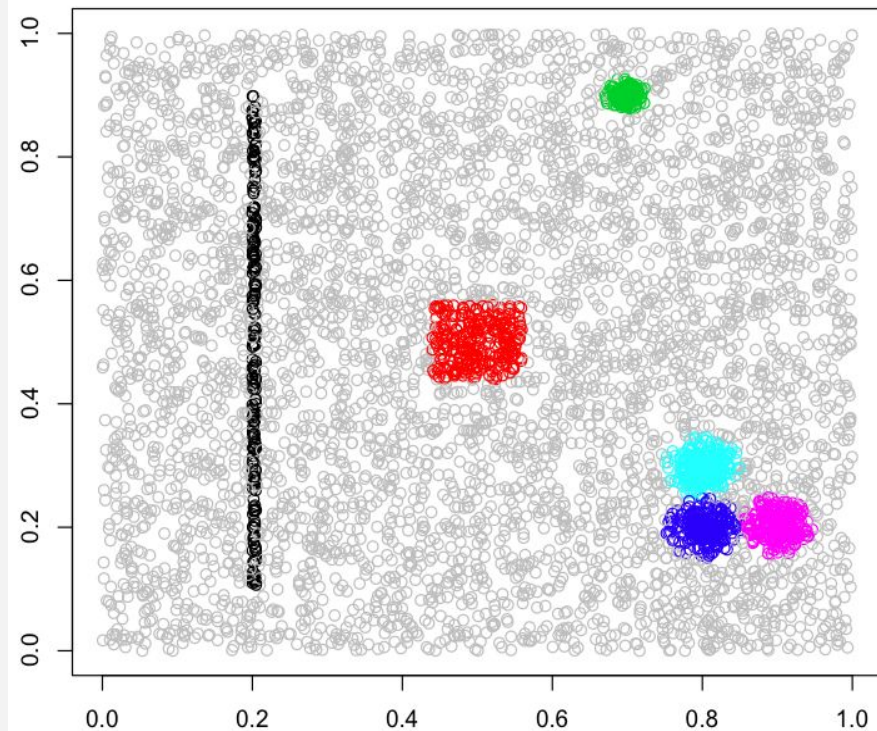
SkinnyDip Running Example Ground Truth



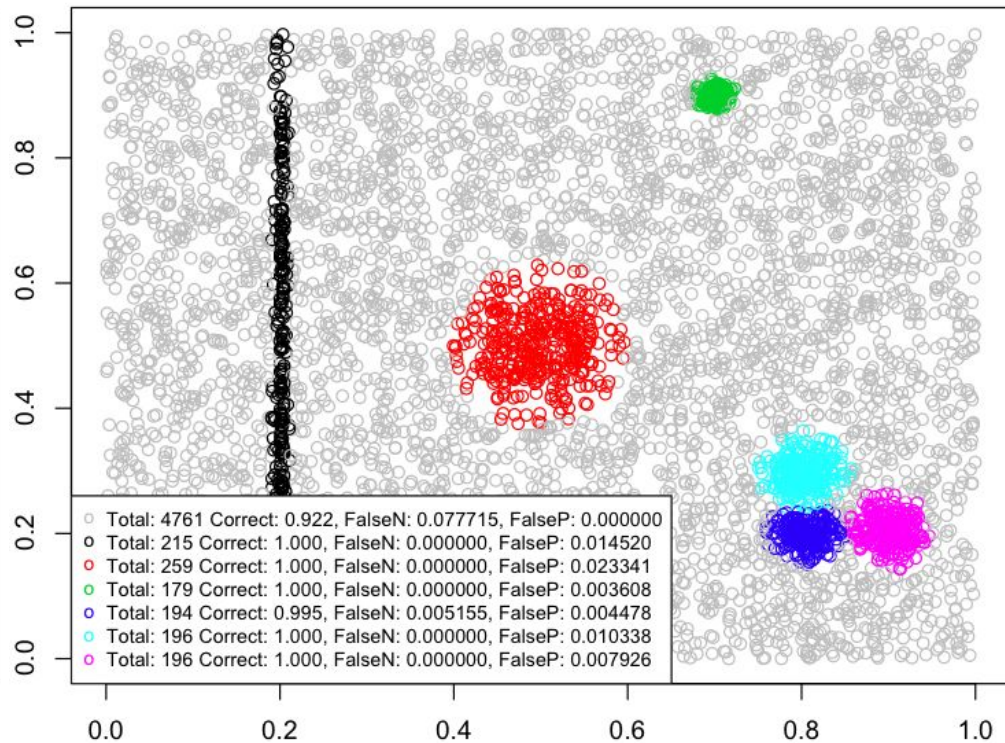
SkinnyDip Running Example



SkinnyDip Running Example Ground Truth

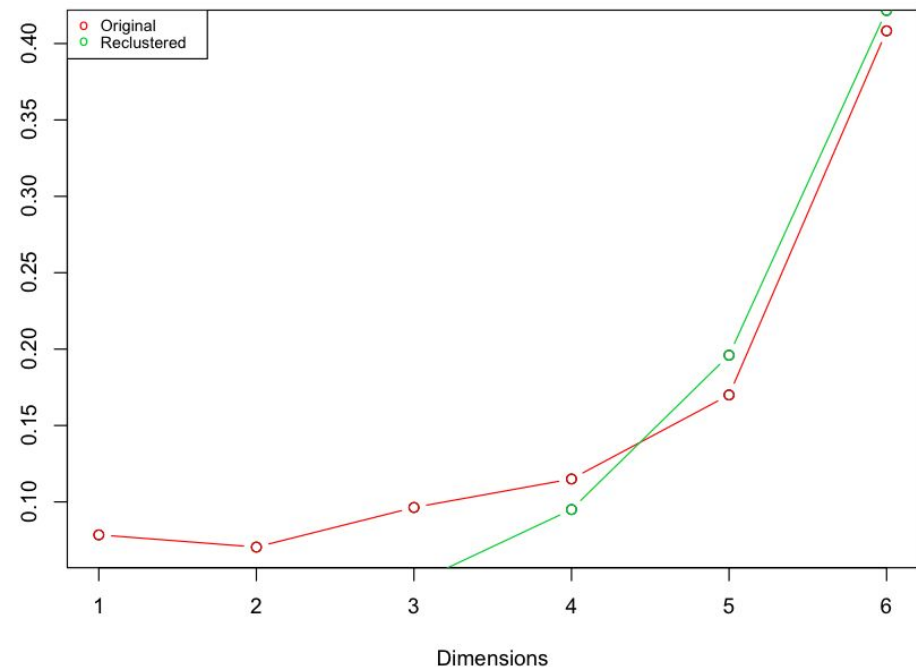


SkinnyDip Running Example with Reclustering



Conclusions

False Negative Rates for original and Reclustered versions of SkinnyDip



- Reclustering to balance out error rates is only practical in lower dimensions
- Small neglected regions still exponentially increase error in higher dimensions
- Less uniform the background noise means more prone to extraneous clusters
- Authors discuss the uses of SkinnyDip on real-world noisy data. Locational classification in 2 or 3 dimensions is still practical

References

- Samuel Maurus and Claudia Plant. 2016. Skinny-dip: Clustering in a Sea of Noise. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (KDD '16). ACM, New York, NY, USA, 1055-1064. DOI: <https://doi.org/10.1145/2939672.2939740>
- J. A. Hartigan and P. Hartigan. The dip test of unimodality. *The Annals of Statistics*, 1985.



Questions?