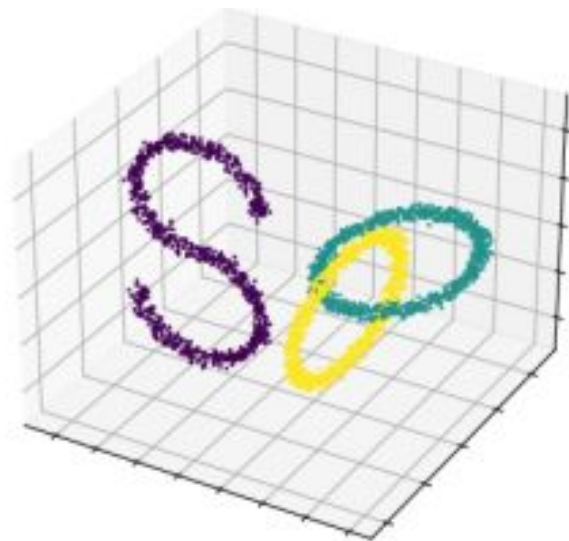# SHADE:
# Deep Density-based Clustering

December 11th, 2024

Anna Beer*, Pascal Weber*,
Lukas Miklautz, Collin Leiber, Walid Durani, Christian Böhm, Claudia Plant

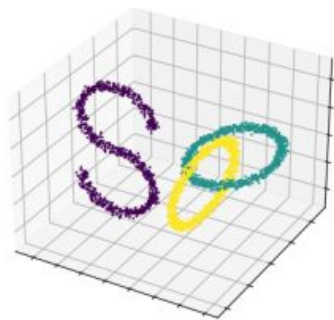ICDM 2024, 9-12 December 2024, Abu Dhabi, UAE

# Motivation



- **Density-based Clustering** is one of the main concepts of clustering

- Clusters of **arbitrary shape** are common in real world data

- Data can contain **noise points**

- Traditional methods are not optimal for **high-dimensional data**

# Motivation

- Autoencoders (AE) optimizing the reconstruction loss are not optimal for **intertwined** or non-contractible clusters



3d data



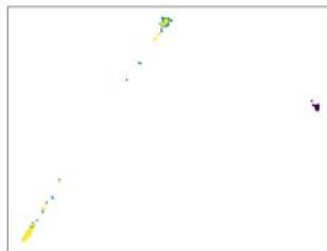2d embedding by
an AE



Desired embedding for
clustering (SHADE)

# Deep Clustering

- Deep Clustering methods usually combine an AE with some cluster loss

$$\mathcal{L} = \mathcal{L}_{rec} + \mathcal{L}_{cluster}$$

- It can be relevant to **preserve the shape** in the embedding



(d) DEC (0.59)    (e) IDEC (0.64)    (f) DCN (0.64)

(g) DipEncoder (0.60)    (h) DipDECK (0.57)    (i) DDC (0.25)

# OPTICS Plots



(a) 3d dataset (0.98)

(b) SHADE (1.00)

(c) Autoencoder (0.58)

(d) DEC (0.26)

(e) IDEC (0.26)

(f) DCN (0.58)

(g) DipEncoder (0.41)

(h) DipDECK (0.67)

(i) DDC (0.44)

# Capturing Density

- Capture the density-connectivity with **dc-distance** $d_{dc}$
- Similar to **minimax path distance**, enriched with concept of density by using **core distances**
- Density-connectivity loss $\mathcal{L}_d$ : Similarity between Euclidean distance in low-dimensional embedding and dc-distance in high-dimensional space

December 11th, 2024 – SHADE: Deep Density-based Clustering – Talk at ICDM 2024 by Anna Beer* and Pascal Weber*

# Overview of SHADE

$$\mathcal{L}_d = \frac{1}{|\mathcal{B}|^2} \sum_{x_i, x_j \in \mathcal{B}} \left( d_{dc}(x_i, x_j) - d_{eucl}(z_i, z_j) \right)^2$$

Training

Density-Connectivity $d_{dc}$  $\mathcal{L}_D$  Euclidean Distances $d_e$

Autoencoder

Data $\mathcal{X}$ → Batch $\mathcal{B}$ → Embedding $\mathcal{Z}$ Reconstr. $\hat{X}$

Clustering $\mathcal{C}$

Final Step

$\mathcal{L}_{rec}$

Density-Connectivity $d_{dc}$ → $argmax(S(\mathcal{C}))$

$$\mathcal{L}_{rec} = \frac{1}{|\mathcal{B}|} \sum_{x_i \in \mathcal{B}} \|x_i - \hat{x}_i\|_2^2$$

7

# Results

| | Dataset | SHADE | SHADE_1nn | DDC | DipDECK | DipEncoder | DCN | DEC | IDEC |
|---|---|---|---|---|---|---|---|---|---|
| **Tabular data** | Synth_low | $99.4 \pm 1.8$ | $\mathbf{98.9 \pm 2.0}$ | $56.9 \pm 5.5$ | $33.9 \pm 6.4$ | $10.1 \pm 9.9$ | $9.6 \pm 9.7$ | $40.2 \pm 3.0$ | $15.3 \pm 8.8$ |
| | Synth_high | $98.4 \pm 1.2$ | $\mathbf{97.5 \pm 1.4}$ | $33.9 \pm 11.1$ | $29.9 \pm 13.6$ | $9.3 \pm 10.7$ | $8.8 \pm 10.5$ | $30.3 \pm 3.5$ | $17.9 \pm 6.6$ |
| | letterrec. | $43.2 \pm 1.7$ | $23.0 \pm 0.9$ | $9.9 \pm 2.9$ | $7.3 \pm 3.5$ | $24.7 \pm 1.3$ | $22.6 \pm 1.0$ | $23.9 \pm 1.6$ | $\mathbf{25.2 \pm 1.8}$ |
| | htru2 | $72.8 \pm 23.4$ | $\mathbf{65.0 \pm 19.5}$ | $49.4 \pm 13.0$ | $9.7 \pm 19.4$ | $4.3 \pm 0.8$ | $49.7 \pm 2.8$ | $3.0 \pm 0.5$ | $3.2 \pm 0.6$ |
| | Mice | $32.5 \pm 3.8$ | $\mathbf{27.7 \pm 2.9}$ | $25.2 \pm 1.9$ | $22.7 \pm 4.3$ | $21.6 \pm 2.6$ | $21.7 \pm 1.4$ | $22.0 \pm 1.5$ | $21.8 \pm 1.4$ |
| | Pendigits | $85.1 \pm 1.4$ | $75.1 \pm 0.8$ | $\mathbf{76.9 \pm 2.0}$ | $74.3 \pm 1.1$ | $64.6 \pm 3.0$ | $61.6 \pm 1.9$ | $65.7 \pm 3.3$ | $64.9 \pm 2.6$ |
| **Video** | Weizmann | $57.1 \pm 5.9$ | $\mathbf{48.2 \pm 3.6}$ | $14.7 \pm 1.8$ | $12.0 \pm 1.9$ | $23.3 \pm 1.2$ | $24.6 \pm 1.1$ | $24.9 \pm 1.2$ | $24.7 \pm 1.2$ |
| | Keck | $9.3 \pm 0.5$ | $\mathbf{7.5 \pm 0.4}$ | $-0.2 \pm 1.1$ | $6.9 \pm 0.8$ | $7.1 \pm 0.3$ | $6.4 \pm 0.5$ | $6.1 \pm 0.9$ | $6.2 \pm 0.9$ |
| **Image** | COIL20 | $82.5 \pm 4.5$ | $\mathbf{68.7 \pm 3.5}$ | $62.0 \pm 5.5$ | $50.5 \pm 7.8$ | $64.0 \pm 3.0$ | $62.4 \pm 2.8$ | $63.7 \pm 2.8$ | $62.9 \pm 2.9$ |
| | COIL100 | $78.1 \pm 7.3$ | $56.8 \pm 5.0$ | $16.4 \pm 3.8$ | $21.4 \pm 3.0$ | $54.3 \pm 1.9$ | $55.9 \pm 3.0$ | $55.8 \pm 2.0$ | $\mathbf{56.9 \pm 2.0}$ |
| | cmu_faces | $38.9 \pm 7.6$ | $34.6 \pm 6.2$ | $35.0 \pm 3.5$ | $29.8 \pm 9.8$ | $37.9 \pm 2.2$ | $\mathbf{40.3 \pm 2.0}$ | $35.8 \pm 2.8$ | $39.4 \pm 3.3$ |

# Robust against Noise

# Hyperparameters - Ablation Studies

- **μ / min_points**: Similar results for the tested values μ $\in$ [3, 7]
  - Default value: μ = 5
- **batchsize**: 300 is large enough to estimate the dc-distance $d_{dc}$ good enough
  - Default value: batchsize = 500


- **No other Parameters!**

# Summary

- SHADE is the first deep density-based clustering method
- It preserves density-connected structures in low-dimensional embeddings
- SHADE finds noise, arbitrarily shaped clusters, and detects the number of clusters fully automatically

ArXiv: https://arxiv.org/abs/2410.06265

Code: https://github.com/pasiweber/SHADE

# Density-based vocabulary

- Core Distance (where $x_\mu$ is the $\mu$-th nearest neighbor of a point $x$)

$$d_{core}(x) = d_{eucl}(x, x_\mu)$$

- Mutual Reachability Distance

$$d_m(x, y) = \max(d_{eucl}(x, y), core\_dist(x), core\_dist(y))$$

- dc-distance $d_{dc}$
  - Create minimum spanning tree on mutual reachability distance
  - Get minimax path distance from that tree - longest edge on the path between x and y is the dc-distance between those points

# Video data often has density-based clusters in HD space

- Each scene of a movie is a density-connected cluster
  - every image/frame is similar to the next one, but beginning and end of a scene might be very different



Walk · Run · Jump · Gallop sideways · Bend · One-hand wave · Two-hands wave · Jump in place · Jumping Jack · Skip

# Extract clustering from the tree metric



$$S(c) = \left( \frac{1}{\mathcal{T}_d(a)} - \frac{1}{\mathcal{T}_d(p(a))} \right) \cdot |l(a)|$$

| | Dataset | $n$ | $d$ | $k$ | #noise | Source |
|---|---|---|---|---|---|---|
| **Tabular data** | Synth_low | 5000 | 100 | 10 | 500 | [18] |
| | Synth_high | 5000 | 100 | 10 | 500 | [18] |
| | HAR | 10,299 | 561 | 6 | 0 | [26] |
| | letterrecognition | 20,000 | 16 | 26 | 0 | [26] |
| | htru2 | 17,898 | 8 | 2 | 0 | [26] |
| | Mice | 1,077 | 68 | 8 | 0 | [26] |
| | TCGA | 801 | 20,264 | 5 | 0 | [26] |
| | Pendigits | 10,992 | 16 | 10 | 0 | [26] |
| **Video** | Weizmann | 5,701 | 77,760 | 90 | 0 | [4] |
| | Keck | 25,457 | 120,000 | 60 | 0 | [43] |
| **Image data** | COIL20 | 1,440 | 16,384 | 20 | 0 | [26] |
| | COIL100 | 7,200 | 49,152 | 100 | 0 | [26] |
| | cmu_faces | 624 | 960 | 20 | 0 | [26] |
| | Optdigits | 5,620 | 64 | 10 | 0 | [26] |
| | USPS | 9,298 | 256 | 10 | 0 | [16] |
| | MNIST | 70,000 | 784 | 10 | 0 | [21] |
| | FMNIST | 70,000 | 784 | 10 | 0 | [40] |
| | KMNIST | 70,000 | 784 | 10 | 0 | [8] |

| Dataset | Metric | SHADE | SHADE_1nn | DDC | DipDECK | DipEncoder | DCN | DEC | IDEC |
|---|---|---|---|---|---|---|---|---|---|
| **Tabular data** | | | | | | | | | |
| Synth_low (*noise*: $1.1 \pm 1.3$) $k = 10$ | ARI | $99.4 \pm 1.8$ | $\mathbf{98.9 \pm 2.0}$ | $56.9 \pm 5.5$ | $33.9 \pm 6.4$ | $10.1 \pm 9.9$ | $9.6 \pm 9.7$ | $40.2 \pm 3.0$ | $15.3 \pm 8.8$ |
| | NMI | $99.7 \pm 1.0$ | $\mathbf{99.2 \pm 1.2}$ | $82.9 \pm 3.2$ | $58.0 \pm 5.0$ | $31.3 \pm 13.6$ | $29.4 \pm 13.4$ | $69.4 \pm 3.0$ | $44.3 \pm 9.5$ |
| | k | $10.0 \pm 0.0$ | $10.0 \pm 0.0$ | $16.8 \pm 1.0$ | $4.1 \pm 0.5$ | - | - | - | - |
| Synth_high (*noise*: $2.4 \pm 1.2$) $k = 10$ | ARI | $98.4 \pm 1.2$ | $\mathbf{97.5 \pm 1.4}$ | $33.9 \pm 11.1$ | $29.9 \pm 13.6$ | $9.3 \pm 10.7$ | $8.8 \pm 10.5$ | $30.3 \pm 3.5$ | $17.9 \pm 6.6$ |
| | NMI | $98.1 \pm 1.2$ | $\mathbf{97.3 \pm 1.3}$ | $61.8 \pm 6.1$ | $53.2 \pm 12.0$ | $28.9 \pm 13.4$ | $26.9 \pm 13.2$ | $61.4 \pm 4.8$ | $46.7 \pm 5.0$ |
| | k | $13.0 \pm 1.5$ | $13.0 \pm 1.5$ | $10.3 \pm 2.1$ | $4.7 \pm 1.3$ | - | - | - | - |
| HAR (*noise*: $3.2 \pm 5.2$) $k = 6$ | ARI | $36.0 \pm 5.6$ | $36.4 \pm 6.4$ | $49.4 \pm 3.3$ | $51.3 \pm 4.0$ | $60.0 \pm 6.9$ | $\mathbf{66.1 \pm 1.3}$ | $63.4 \pm 2.4$ | $64.9 \pm 0.9$ |
| | NMI | $58.5 \pm 5.8$ | $58.2 \pm 5.5$ | $68.2 \pm 2.1$ | $71.3 \pm 2.0$ | $73.6 \pm 4.1$ | $\mathbf{75.4 \pm 1.2}$ | $75.0 \pm 1.8$ | $74.6 \pm 0.7$ |
| | k | $3.3 \pm 2.0$ | $3.3 \pm 2.0$ | $4.2 \pm 1.0$ | $3.1 \pm 0.3$ | - | - | - | - |
| letterrec. (*noise*: $50.3 \pm 1.8$) $k = 26$ | ARI | $43.2 \pm 1.7$ | $23.0 \pm 0.9$ | $9.9 \pm 2.9$ | $7.3 \pm 3.5$ | $24.7 \pm 1.3$ | $22.6 \pm 1.0$ | $23.9 \pm 1.6$ | $\mathbf{25.2 \pm 1.8}$ |
| | NMI | $75.6 \pm 1.0$ | $\mathbf{57.4 \pm 0.5}$ | $43.5 \pm 2.5$ | $34.4 \pm 3.8$ | $49.7 \pm 0.9$ | $46.4 \pm 1.0$ | $50.8 \pm 1.0$ | $49.7 \pm 1.2$ |
| | k | $111.4 \pm 4.9$ | $111.4 \pm 4.9$ | $14.0 \pm 1.0$ | $13.0 \pm 2.4$ | - | - | - | - |
| htru2 (*noise*: $14.0 \pm 3.8$) $k = 2$ | ARI | $72.8 \pm 23.4$ | $\mathbf{65.0 \pm 19.5}$ | $49.4 \pm 13.0$ | $9.7 \pm 19.4$ | $4.3 \pm 0.8$ | $49.7 \pm 2.8$ | $3.0 \pm 0.5$ | $3.2 \pm 0.6$ |
| | NMI | $59.9 \pm 19.6$ | $\mathbf{47.4 \pm 14.4}$ | $42.1 \pm 7.0$ | $5.5 \pm 11.0$ | $10.5 \pm 1.4$ | $31.6 \pm 7.0$ | $10.8 \pm 0.6$ | $10.7 \pm 0.6$ |
| | k | $5.9 \pm 2.1$ | $5.9 \pm 2.1$ | $3.7 \pm 0.6$ | $1.2 \pm 0.4$ | - | - | - | - |
| Mice (*noise*: $30.0 \pm 5.5$) $k = 8$ | ARI | $32.5 \pm 3.8$ | $\mathbf{27.7 \pm 2.9}$ | $25.2 \pm 1.9$ | $22.7 \pm 4.3$ | $21.6 \pm 2.6$ | $21.7 \pm 1.4$ | $22.0 \pm 1.5$ | $21.8 \pm 1.4$ |
| | NMI | $56.2 \pm 4.6$ | $48.7 \pm 2.3$ | $\mathbf{49.6 \pm 2.0}$ | $48.0 \pm 5.3$ | $38.7 \pm 2.6$ | $38.3 \pm 1.5$ | $39.4 \pm 1.8$ | $38.3 \pm 1.8$ |
| | k | $11.9 \pm 2.8$ | $11.9 \pm 2.8$ | $15.3 \pm 1.2$ | $13.4 \pm 5.3$ | - | - | - | - |
| TCGA (*noise*: $21.8 \pm 9.0$) $k = 5$ | ARI | $86.6 \pm 7.8$ | $80.0 \pm 13.7$ | $87.5 \pm 0.8$ | $88.8 \pm 4.4$ | $\mathbf{93.4 \pm 6.0}$ | $87.2 \pm 5.3$ | $85.1 \pm 2.7$ | $82.6 \pm 0.9$ |
| | NMI | $90.8 \pm 4.4$ | $87.4 \pm 7.2$ | $91.9 \pm 0.7$ | $91.9 \pm 2.1$ | $\mathbf{94.0 \pm 3.8}$ | $89.1 \pm 3.1$ | $89.5 \pm 1.2$ | $86.1 \pm 1.4$ |
| | k | $5.7 \pm 0.9$ | $5.7 \pm 0.9$ | $6.0 \pm 0.0$ | $5.9 \pm 0.5$ | - | - | - | - |
| Pendigits (*noise*: $17.5 \pm 2.2$) $k = 10$ | ARI | $85.1 \pm 1.4$ | $75.1 \pm 0.8$ | $\mathbf{76.9 \pm 2.0}$ | $74.3 \pm 1.1$ | $64.6 \pm 3.0$ | $61.6 \pm 1.9$ | $65.7 \pm 3.3$ | $64.9 \pm 2.6$ |
| | NMI | $89.3 \pm 0.9$ | $82.7 \pm 0.6$ | $\mathbf{84.1 \pm 1.0}$ | $82.0 \pm 0.8$ | $75.2 \pm 1.3$ | $72.7 \pm 0.6$ | $76.7 \pm 1.4$ | $75.9 \pm 0.9$ |
| | k | $20.0 \pm 1.8$ | $20.0 \pm 1.8$ | $13.0 \pm 0.4$ | $14.8 \pm 0.7$ | - | - | - | - |
| **Video data** | | | | | | | | | |
| Weizmann (*noise*: $14.6 \pm 2.6$) $k = 90$ | ARI | $57.1 \pm 5.9$ | $\mathbf{48.2 \pm 3.6}$ | $14.7 \pm 1.8$ | $12.0 \pm 1.9$ | $23.3 \pm 1.2$ | $24.6 \pm 1.1$ | $24.9 \pm 1.2$ | $24.7 \pm 1.2$ |
| | NMI | $86.6 \pm 0.6$ | $\mathbf{80.2 \pm 1.1}$ | $57.5 \pm 2.0$ | $52.8 \pm 1.9$ | $61.9 \pm 0.8$ | $62.3 \pm 0.9$ | $63.7 \pm 1.0$ | $63.3 \pm 1.0$ |
| | k | $57.1 \pm 2.1$ | $57.1 \pm 2.1$ | $20.8 \pm 2.4$ | $24.5 \pm 2.2$ | $88.4 \pm 1.2^*$ | - | - | - |
| Keck (*noise*: $25.4 \pm 1.1$) $k = 60$ | ARI | $9.3 \pm 0.5$ | $\mathbf{7.5 \pm 0.4}$ | $-0.2 \pm 1.1$ | $6.9 \pm 0.8$ | $7.1 \pm 0.3$ | $6.4 \pm 0.5$ | $6.1 \pm 0.9$ | $6.2 \pm 0.9$ |
| | NMI | $66.3 \pm 0.3$ | $\mathbf{61.6 \pm 0.3}$ | $18.1 \pm 4.3$ | $34.9 \pm 6.0$ | $42.4 \pm 0.6$ | $39.4 \pm 2.5$ | $39.9 \pm 3.5$ | $39.2 \pm 5.2$ |
| | k | $226.6 \pm 10.4$ | $226.6 \pm 10.4$ | $10.6 \pm 3.1$ | $32.3 \pm 7.4$ | - | - | - | - |
| **Image data** | | | | | | | | | |
| COIL20 (*noise*: $12.5 \pm 1.9$) $k = 20$ | ARI | $82.5 \pm 4.5$ | $\mathbf{68.7 \pm 3.5}$ | $62.0 \pm 5.5$ | $50.5 \pm 7.8$ | $64.0 \pm 3.0$ | $62.4 \pm 2.8$ | $63.7 \pm 2.8$ | $62.9 \pm 2.9$ |
| | NMI | $93.6 \pm 1.7$ | $\mathbf{85.6 \pm 1.3}$ | $85.5 \pm 0.9$ | $79.9 \pm 2.4$ | $80.2 \pm 1.1$ | $79.6 \pm 1.3$ | $80.6 \pm 1.0$ | $80.0 \pm 1.1$ |
| | k | $16.5 \pm 1.0$ | $16.5 \pm 1.0$ | $14.3 \pm 0.8$ | $18.9 \pm 1.0$ | - | - | - | - |
| COIL100 (*noise*: $24.9 \pm 1.7$) $k = 100$ | ARI | $78.1 \pm 7.3$ | $56.8 \pm 5.0$ | $16.4 \pm 3.8$ | $21.4 \pm 3.0$ | $54.3 \pm 1.9$ | $55.9 \pm 3.0$ | $55.8 \pm 2.0$ | $\mathbf{56.9 \pm 2.0}$ |
| | NMI | $94.4 \pm 0.9$ | $85.4 \pm 0.6$ | $69.5 \pm 3.0$ | $69.5 \pm 1.7$ | $82.6 \pm 0.6$ | $84.2 \pm 0.7$ | $\mathbf{85.8 \pm 0.5}$ | $84.7 \pm 0.5$ |
| | k | $77.8 \pm 3.2$ | $77.8 \pm 3.2$ | $18.3 \pm 3.3$ | $26.8 \pm 1.9$ | $99.8 \pm 0.4$ | - | - | - |
| cmu_faces (*noise*: $15.1 \pm 4.1$) $k = 20$ | ARI | $38.9 \pm 7.6$ | $34.6 \pm 6.2$ | $35.0 \pm 3.5$ | $29.8 \pm 9.8$ | $37.9 \pm 2.2$ | $\mathbf{40.3 \pm 2.0}$ | $35.8 \pm 2.8$ | $39.4 \pm 3.3$ |
| | NMI | $69.1 \pm 5.3$ | $65.0 \pm 4.9$ | $64.4 \pm 2.3$ | $62.3 \pm 7.6$ | $67.7 \pm 1.6$ | $\mathbf{68.5 \pm 1.0}$ | $66.1 \pm 2.0$ | $68.2 \pm 2.2$ |
| | k | $10.6 \pm 1.7$ | $10.6 \pm 1.7$ | $12.0 \pm 0.9$ | $8.9 \pm 2.2$ | - | - | - | - |

| Dataset | Metric | SHADE | SHADE_1nn | DDC | DipDECK | DipEncoder | DCN | DEC | IDEC |
|---|---|---|---|---|---|---|---|---|---|
| Optdigits | ARI | $93.3 \pm 1.9$ | $78.0 \pm 7.4$ | $\mathbf{88.9 \pm 2.3}$ | $82.2 \pm 2.3$ | $80.2 \pm 4.0$ | $77.0 \pm 3.5$ | $80.4 \pm 3.1$ | $80.7 \pm 3.6$ |
| (noise: $38.2 \pm 5.3$) | NMI | $94.6 \pm 1.0$ | $83.4 \pm 3.9$ | $\mathbf{91.7 \pm 1.3}$ | $86.2 \pm 0.9$ | $85.6 \pm 2.0$ | $82.9 \pm 1.6$ | $\underline{86.3 \pm 1.4}$ | $86.2 \pm 1.5$ |
| $k = 10$ | $k$ | $12.4 \pm 1.0$ | $12.4 \pm 1.0$ | $11.1 \pm 0.5$ | $11.0 \pm 1.2$ | - | - | - | - |
| USPS | ARI | $88.1 \pm 2.7$ | $68.1 \pm 1.5$ | $\mathbf{92.2 \pm 2.2}$ | $68.2 \pm 5.1$ | $72.5 \pm 2.0$ | $66.2 \pm 1.3$ | $73.4 \pm 0.9$ | $74.0 \pm 0.9$ |
| (noise: $45.7 \pm 2.8$) | NMI | $91.8 \pm 1.1$ | $75.9 \pm 1.1$ | $\mathbf{91.0 \pm 0.8}$ | $78.5 \pm 2.5$ | $79.9 \pm 1.6$ | $74.5 \pm 1.0$ | $81.0 \pm 0.5$ | $\underline{81.3 \pm 0.6}$ |
| $k = 10$ | $k$ | $9.5 \pm 0.9$ | $9.5 \pm 0.9$ | $9.8 \pm 0.4$ | $8.0 \pm 1.2$ | - | - | - | - |
| MNIST | ARI | $86.1 \pm 8.4$ | $54.1 \pm 2.8$ | $\mathbf{94.5 \pm 1.2}$ | $80.9 \pm 3.0$ | $79.6 \pm 2.8$ | $\underline{82.0 \pm 4.5}$ | $79.5 \pm 2.1$ | $81.9 \pm 2.0$ |
| (noise: $58.9 \pm 5.3$) | NMI | $84.5 \pm 2.4$ | $63.7 \pm 1.3$ | $\mathbf{93.7 \pm 0.8}$ | $84.5 \pm 1.5$ | $84.9 \pm 1.6$ | $84.6 \pm 1.9$ | $85.1 \pm 1.0$ | $\underline{87.5 \pm 0.8}$ |
| $k = 10$ | $k$ | $38.0 \pm 9.9$ | $38.0 \pm 9.9$ | $10.0 \pm 0.0$ | $11.3 \pm 0.8$ | - | - | - | - |
| FMNIST | ARI | $72.2 \pm 2.7$ | $35.7 \pm 1.4$ | $35.5 \pm 7.1$ | $\underline{46.5 \pm 0.9}$ | $43.5 \pm 3.6$ | $45.8 \pm 2.9$ | $42.3 \pm 4.1$ | $\mathbf{46.9 \pm 3.6}$ |
| (noise: $66.3 \pm 2.3$) | NMI | $73.4 \pm 1.0$ | $53.9 \pm 0.9$ | $64.1 \pm 3.3$ | $\mathbf{65.3 \pm 0.7}$ | $59.4 \pm 2.4$ | $62.4 \pm 2.2$ | $59.3 \pm 2.7$ | $63.7 \pm 2.5$ |
| $k = 10$ | $k$ | $68.7 \pm 9.2$ | $68.7 \pm 9.2$ | $7.1 \pm 0.9$ | $9.6 \pm 1.1$ | - | - | - | - |
| KMNIST | ARI | $58.6 \pm 9.4$ | $27.5 \pm 1.7$ | $\mathbf{60.3 \pm 3.5}$ | $33.8 \pm 9.9$ | $\underline{43.6 \pm 1.1}$ | $40.3 \pm 1.6$ | $42.4 \pm 0.9$ | $42.8 \pm 1.1$ |
| (noise: $73.4 \pm 2.2$) | NMI | $66.2 \pm 3.1$ | $46.8 \pm 0.6$ | $\mathbf{73.6 \pm 1.6}$ | $55.9 \pm 5.0$ | $\underline{56.7 \pm 1.3}$ | $53.6 \pm 1.4$ | $55.7 \pm 1.1$ | $56.7 \pm 1.2$ |
| $k = 10$ | $k$ | $64.1 \pm 10.4$ | $64.1 \pm 10.4$ | $16.2 \pm 0.7$ | $11.7 \pm 2.8$ | - | - | - | - |